

Министерство образования Российской Федерации
Ярославский государственный университет
имени П.Г.Демидова

МОДЕЛИРОВАНИЕ И АНАЛИЗ ИНФОРМАЦИОННЫХ СИСТЕМ

Том 10 №2 2003

Основан в 1999 г.
Выходит 2 раза в год

*Свидетельство о регистрации №019209 от 16.08.99
Государственного Комитета Российской Федерации по печати*

Главный редактор
В.А.Соколов

Редакционная коллегия
О.Л.Бандман, В.А.Бондаренко, М.Г.Дмитриев, А.В.Зафиевский,
Ю.Г.Карпов, С.А.Кащенко, Ю.С.Колесов, А.Ю.Левин,
И.А.Ломазова, В.В.Майоров, В.Э.Малышкин, В.А.Непомнящий

Ответственный секретарь
Е.А.Тимофеев

Адрес редакции: 150000, Ярославль, ул. Советская, 14
E-mail: mais@uniyar.ac.ru

Научные статьи в журнал принимаются на кафедре ТИ. Статья должна содержать УДК, аннотацию и сопровождаться набором текста в редакторе LaTEX.

©Ярославский
государственный
университет, 2003

СОДЕРЖАНИЕ

Моделирование и анализ информационных систем. Т.10, №2. 2003

Сводимость задач дискретной оптимизации и соотношение плотностей их полиэдральных графов <i>Максименко А.Н.</i>	3
Моделирование протоколов TCP и ARTCP с помощью раскрашенных сетей Петри <i>Чалый Д.Ю.</i>	11
Об одной нелинейной краевой задаче, моделирующей экономические циклы <i>Косарева Е.С. Кумиков А.Н.</i>	18
Об одной модификации организации пачечной волновой активности в сетях W-нейронов <i>Овчинин А.С.</i>	22
Выбор ранговых и доменных блоков в алгоритме фрактального сжатия изображений <i>Ходаковский В.А.</i>	26
Планирование волн в кольцевых структурах ассоциаций импульсных нейронов-детекторов <i>Майоров В.В., Мышкин И.Ю., Мячин М.Л., Куксов А.Г.</i>	30
О хаотическом поведении одной модели нейронной сети <i>Богомолов Ю.В.</i>	35
Недетерминированные счётчиковые машины <i>Кузьмин Е.В.</i>	41
Об обобщении критерия однородности А.Ю. Левина <i>Янкевич А.П.</i>	50

Корректор А.А.Аладьева
Подписано в печать 10.11.2003. Формат 60x84¹/₈. Печать офсетная.
Усл.печ.л. 10,50. Уч.-изд.л. 5,95. Тираж 100 экз.

Отпечатано на ризографе. Ярославский государственный университет имени П.Г.Демидова, 150 000, Ярославль, ул.Советская, 14

Сводимость задач дискретной оптимизации и соотношение плотностей их полиэдральных графов¹

Максименко А.Н.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

получена 7 марта 2003

В работе устанавливается взаимосвязь между сводимостью задач дискретной оптимизации и соотношением плотностей их полиэдральных графов. На основе этих соображений удается получить новую экспоненциальную нижнюю оценку плотности графа многогранника для задачи коммивояжера.

Различные вопросы, связанные с анализом сложных систем, приводят к задачам дискретной оптимизации, допускающим следующую формулировку. Задано конечное множество X точек в евклидовом пространстве R^m , каждая из которых соответствует допустимому решению задачи, и вектор $c \in R^m$ исходных данных; требуется оптимизировать (найти максимум или минимум) на X функцию — скалярное произведение (x, c) , где $x \in X$. К указанному виду легко преобразуются задачи сортировки, задача коммивояжера, задача о назначениях, задача о минимальном остовном дереве и многие другие. Сформулированную выше индивидуальную задачу с фиксированным вектором c будем обозначать $[X, c]$; а совокупность таких задач, для всех $c \in R^m$, обозначим через X . Практическая необходимость создания эффективных алгоритмов решения таких задач стимулирует поиски ответов на следующий вопрос: Реально ли рассчитывать на создание эффективного алгоритма решения той или иной задачи и, если это возможно, то какими свойствами должен обладать такой алгоритм? Поставленный вопрос можно уточнить вопросом о том, какие свойства задачи являются характеристикой ее сложности.

Частичный ответ на поставленный вопрос можно получить с помощью теории С. Кука NP-полных задач [1]. Согласно этой теории задача объявляется труднорешаемой, если некоторая признанно труднорешаемая задача является частным случаем данной. Такой подход обладает следующими "полярными" особенностями: во-первых, он достаточно прост в практическом применении, во-вторых, в нем отсутствует анализ характеристик задачи.

Другой подход к исследованию сложности задач основан на изучении комбинаторных свойств многогранника $convX$. В частности, известно [2], что плотность графа многогранника является нижней оценкой сложности соответствующей задачи в широком классе алгоритмов, использующих линейные сравнения. (Напомним, что плотность графа равна максимальному числу вершин, любые две из которых смежны.) Установлено [2], что к этому классу относятся практически все классические комбинаторные алгоритмы: алгоритмы сортировки, "жадный" алгоритм, различные варианты метода ветвей и границ, метод динамического программирования, алгоритмы типа "локальный поиск", венгерский алгоритм и другие широко распространенные практические методы комбинаторной оптимизации. Ограниченностю практического применения этого подхода обусловлена, во-первых, трудностью вычисления плотности графа большого размера, и, во-вторых, для многогранников некоторых задач, в частности для задачи коммивояжера [3], уже задача распознавания смежности вершин является NP-полной.

Описываемые в настоящей работе результаты дают возможность объединить указанные подходы и позволяют, с одной стороны, сделать вычисление комбинаторных характеристик полиэдрального графа задачи столь же "обычным" занятием, как и установление ее NP-полноты, а с другой стороны, сделать вывод о том, что сложность труднорешаемой задачи в упомянутом выше широком классе алгоритмов, как правило, оказывается экспоненциальной.

Ниже на основе известного [1] алгоритма сведения задачи распознавания КЛИКА к задаче ГАМИЛЬТОНОВ ЦИКЛ строится алгоритм сведения их оптимизационных вариантов. Геометрическая интерпретация этого алгоритма позволяет, в частности, получить более высокую, чем известная ранее [2, 4],

¹Работа выполнена при финансовой поддержке Российского фонда фундаментальных исследований (грант 03-01-00822).

экспоненциальную оценку плотности графа многогранника задачи коммивояжера.

Конусное разбиение. Свойства

Пусть $X, X \subset R^m$, — массовая задача, то есть совокупность всех индивидуальных задач $[X, c], c \in R^m$, на максимум. Многогранником $M(X)$ задачи X называется выпуклая оболочка множества X : $\text{conv } X = M(X)$ [5]. Будем считать, что множество его вершин $\text{ext } M(X)$ совпадает с X (для рассматриваемых ниже задач это естественное предположение выполнено).

Совокупность многогранных множеств

$$K(x) = \{c \in R^m : (x, c) \geq (y, c) \text{ для каждого } y \in X\},$$

где $x \in X$, называется конусным разбиением [2] пространства R^m исходных данных по множеству X . Из определения следует, что если x является вершиной многогранника $M(X)$, то конус $K(x)$ — телесен (то есть $\dim K(x) = m$). Будем говорить, что конусы $K(x)$ и $K(y)$ смежны, если соответствующие вершины x и y смежны. (Легко доказывается [2], что конусы, соответствующие смежным вершинам, имеют общую гипергрань: $\dim\{K(x) \cap K(y)\} = m-1$). Поэтому граф $G(X)$ многогранника $M(X)$ будем также называть графиком конусного разбиения пространства R^m по множеству X .

Рассмотрим частный случай описанной выше задачи X , когда множество исходных данных представляет собой многогранное множество (полиэдр) $S \subset R^m$ (причем S может быть и меньшей размерности, чем R^m). Будем обозначать такую задачу (X, S) . В этом случае аналог конусного разбиения множества S по множеству X образуется полиздрами $K(x, S) = K(x) \cap S$, где $x \in X$. Граф $G_S(X)$ такого многогранного разбиения определим по аналогии с графиком многогранника $M(X)$:

- 1) Полиэдр $K(x, S)$ соответствует вершине графа $G_S(X)$, если существует $s \in S$ такой, что для каждого $y \in X$, отличного от x , выполнено $(x, s) > (y, s)$.
- 2) Вершины графа $G_S(X)$, соответствующие полиздрамам $K(x, S)$ и $K(y, S)$, соединены ребром, если найдется такой $s \in S$, что для каждого $z \in X$ отличного от x и от y , выполнено $(x, s) = (y, s) > (z, s)$. В таком случае будем говорить, что полиздры $K(x, S)$ и $K(y, S)$ смежны.

Непосредственно из определения графа $G_S(X)$ следует, что он является подграфом графа $G(X)$; будем обозначать $G_S(X) \prec G(X)$.

Рассмотрим характерный пример, когда $S = Q_m = \{x \in R^m : -1 \leq x_i \leq 1, x = (x_1, x_2, \dots, x_m)\}$ — куб с длиной ребра, равной 2, и центром в начале координат. Индивидуальная задача $[X, c]$, где $c \in R^m$, легко сводится к задаче $[X, s]$, где $s \in Q_m$. Для этого достаточно взять $s = \frac{c}{\|c\|}$.

Теорема 1. Пусть $X \subset R^m$, $X = \text{ext } M(X)$ и $Q = Q_m$, тогда $G_Q(X) = G(X)$.

Справедливость этой теоремы следует из того, что для любого вектора $c \in R^m$ из определения смежности конусов $K(x)$ и $K(y)$ всегда найдется вектор $s = \frac{c}{\|c\|}$, удовлетворяющий условию смежности полиздров $K(x, Q)$ и $K(y, Q)$.

Введем понятие "аффинной сводимости". Будем говорить, что задача (X, S) максимизации на множестве X для многогранного множества исходных данных $S \subset R^m$ аффинно сводится к задаче максимизации на множестве Y для множества исходных данных R^n (где $m \leq n$), если найдется такое невырожденное аффинное отображение $A : S \rightarrow T$ вида $t = A(s) = Cs + d$, где $t \in T \subseteq R^n$, $s \in S$, $d \in R^n$ и C — матрица размера $n \times m$, причем $\text{rang } C = m$, и найдется взаимно-однозначное соответствие $B : X \rightarrow Y'$ между множеством X и некоторым подмножеством $Y' \subseteq Y$ такие, что для каждого вектора $s \in S$ выполнено следующее условие:

$$y_0 \text{ — решение задачи } [Y, A(s)] \Leftrightarrow y_0 \in Y' \text{ и } x_0 = B^{-1}(y_0) \text{ — решение задачи } [X, s]. \quad (1)$$

Для таких задач введем обозначение: $(X, S) \propto_A Y$. Пусть $x \in X$ и $y = B(x)$. Тогда из условия (1) и невырожденности отображения A следует, что полиздры $K(x, S)$ и $K(y, T)$ — подобны. Более того, грани вида $K(x_1) \cap K(x_2)$, где $x_1, x_2 \in X$, аффинное отображение переводит в подобные грани $K(y_1) \cap K(y_2)$, где $y_1 = B(x_1)$, $y_2 = B(x_2)$. Сформулируем результат этих рассуждений в виде теоремы.

Теорема 2. Пусть $(X, S) \propto_A Y$ и $A : S \rightarrow T$, тогда $G_S(X) = G_T(Y) \prec G(Y)$.

С учетом теоремы 1 получаем

Следствие 1. Если $(X, Q_m) \propto_A Y$, то $G(X) \prec G(Y)$.

Многогранник задачи "максимальная клика". Свойства

Рассмотрим следующий вариант задачи максимальная клика (МК). Задан полный ребристо и вершинно взвешенный граф $G(V, E)$, где $V = \{v_1, v_2, \dots, v_n\}$ — множество вершин, $E = \{e_{ij} = (v_i, v_j), 1 \leq i < j \leq n\}$

— множество ребер. Его ребрам e_{ij} приписаны веса $c_{ij} \in R$, а вершинам v_i приписаны веса $c_{ii} \in R$. Задача: отыскать такое подмножество вершин $U \subseteq V$, где $0 \leq |U| \leq n$, что полный подграф графа $G(V, E)$, построенный на этих вершинах, будет максимальным по суммарному весу входящих в него ребер и вершин:

$$U : \forall W \subseteq V \quad \sum_{v_i, v_j \in U} c_{ij} \geq \sum_{v_i, v_j \in W} c_{ij},$$

где $1 \leq i \leq j \leq n$.

Чтобы построить множество X для этой задачи, сопоставим каждому ребру e_{ij} координату x_{ij} , а каждой вершине v_i координату x_{ii} для произвольной точки x в пространстве R^m , $m = \frac{n(n+1)}{2}$. Каждому полному подграфу графа $G(V, E)$, построеному на множестве вершин U , поставим в соответствие точку x с координатами

$$x_{ij} = \begin{cases} 1, & \text{если } v_i, v_j \in U \\ 0, & \text{в остальных случаях} \end{cases}$$

Множество X всех таких точек и будет множеством допустимых решений задачи максимальная клика (МК).

Рассмотрим многогранник $M(X)$ задачи МК и покажем, что все точки множества X являются вершинами этого многогранника. Следуя определению вершины, покажем, что для любого $x \in X$ найдется такой $c \in R^m$, что для всех $y \in X \setminus \{x\}$ выполнено $(x, c) > (y, c)$. Зафиксируем x . Пусть U_x — множество вершин соответствующего подграфа. Очевидно $|U_x| = \sum_{1 \leq i \leq n} x_{ii}$. Пусть $k = |U_x|$. Выберем координаты вектора c :

$$c_{ij} = \begin{cases} 0, & \text{если } i \neq j \\ -1, & \text{если } x_{ii} = 0, i = j \\ \frac{1}{k}, & \text{если } x_{ii} = 1, i = j \end{cases}$$

Возьмем $y \in X \setminus \{x\}$. Пусть U_y — множество вершин соответствующего подграфа. Тогда возможны два варианта:

- 1) $U_y \subset U_x$, причем возможно, что $U_y = \emptyset$. Тогда, очевидно, $(y, c) < (x, c) = 1$,
- 2) $U_y \setminus U_x \neq \emptyset$, тогда найдется $v_i \in U_y$ такое, что $v_i \notin U_x$. Но тогда $y_{ii} * c_{ii} = -1$ и $(y, c) \leq (x, c) - 1$.

Теорема 3. Любые две вершины многогранника $M(X)$ задачи максимальная клика смежны.

Доказательство. Схема доказательства заимствована из монографии [2], где доказан аналогичный результат для задачи максимальная клика на графе с не взвешенными вершинами.

Пусть x и y — две произвольные вершины многогранника $M(X)$. Для доказательства их смежности определим такой вектор $c \in R^m$, для которого

$$(x, c) = (y, c) > (z, c) \tag{2}$$

для всех $z \in X \setminus \{x, y\}$. Обозначим через U_x множество вершин подграфа, соответствующему точке x , а через U_y — множество вершин подграфа, соответствующему y . И пусть $k_x = |U_x|$, $k_y = |U_y|$.

Рассмотрим четыре варианта.

- 1) $|U_y \cap U_x| = 0$, $U_x \neq \emptyset$, $U_y \neq \emptyset$. Положим

$$c_{ij} = \begin{cases} -2, & \text{если } x_{ij} = 0 = y_{ij} \\ 0, & \text{если } x_{ij} = 1 \text{ или } y_{ij} = 1, i \neq j \\ \frac{1}{k_x}, & \text{если } x_{ij} = 1 \text{ и } y_{ij} = 0, i = j \\ \frac{1}{k_y}, & \text{если } x_{ij} = 0 \text{ и } y_{ij} = 1, i = j \end{cases}$$

Легко проверяется, что $(x, c) = (y, c) = 1$. Обозначим через U_z множество вершин подграфа, соответствующему точке $z \in X \setminus \{x, y\}$. Возможны два случая.

- a) $U_z \not\subseteq U_x$ и $U_z \not\subseteq U_y$. Тогда найдутся $v_i, v_j \in U_z$ такие, что $v_i \notin U_x$ и $v_j \notin U_y$. Но тогда $z_{ij} = 1$, в то время как $c_{ij} = -2$. Следовательно, $(z, c) \leq 0 < (x, c)$.
- б) $U_z \subset U_x$ или $U_z \subset U_y$. Допустим, для определенности, что $U_z \subset U_x$. Но тогда, очевидно, $(x, c) > (z, c)$.

И условие (2) выполнено.

2) $|U_y \cap U_x| = k > 0$ и ни одно из множеств U_x и U_y не является частью другого. Положим

$$c_{ij} = \begin{cases} -2, & \text{если } x_{ij} = 0 = y_{ij} \\ 0, & \text{если } x_{ij} = 1 \text{ или } y_{ij} = 1, i \neq j \\ \frac{1}{2k}, & \text{если } x_{ij} = 1 = y_{ij}, i = j \\ \frac{1}{2(k_x - k)}, & \text{если } x_{ij} = 1 \text{ и } y_{ij} = 0, i = j \\ \frac{1}{2(k_y - k)}, & \text{если } x_{ij} = 0 \text{ и } y_{ij} = 1, i = j \end{cases}$$

Действуя так же, как и в первом варианте, непосредственно устанавливается, что $(x, c) = (y, c) = 1$ и $(z, c) < 1$, при $z \in X \setminus \{x, y\}$.

3) Одно из множеств U_x, U_y является подмножеством другого и $||U_x| - |U_y|| = 1$. Допустим, для определенности, что $U_x \subset U_y$, причем возможно, что $U_x = \emptyset$. Тогда возьмем

$$c_{ij} = \begin{cases} -2, & \text{если } x_{ij} = 0 = y_{ij} \\ 0, & \text{если } y_{ij} = 1, i \neq j \\ \frac{1}{k_x}, & \text{если } x_{ij} = 1 = y_{ij}, i = j \\ 0, & \text{если } x_{ij} = 0 \text{ и } y_{ij} = 1, i = j \end{cases}$$

Очевидно, что $(x, c) = (y, c) = 1$, если $|U_x| > 0$. И $(z, c) = 1$ только если $z = x$, или $z = y$. Если же $|U_x| = 0$, то $|U_y| = 1$ и $(x, c) = (y, c) = 0 > (z, c)$ для всех $z \in X \setminus \{x, y\}$.

4) Одно из множеств U_x, U_y является подмножеством другого и $||U_x| - |U_y|| > 1$. Допустим, для определенности, что $U_x \subset U_y$, причем возможно, что $U_x = \emptyset$. Положим

$$c_{ij} = \begin{cases} -2, & \text{если } x_{ij} = 0 = y_{ij} \\ \frac{1}{k_x}, & \text{если } x_{ij} = 1 = y_{ij}, i = j \\ -\frac{1}{k_y - k_x}, & \text{если } x_{ij} = 0 \text{ и } y_{ij} = 1, i = j \\ \frac{2}{(k_y - k_x)(k_y - k_x - 1)}, & \text{если } x_{ii} = 0 = x_{jj} \text{ и } y_{ij} = 1, i \neq j \\ 0, & \text{в остальных случаях} \end{cases}$$

Не трудно проверить, что $(x, c) = (y, c) = 1$, если $k_x > 0$. Обозначим через U_z множество вершин подграфа, соответствующего точке $z \in X \setminus \{x, y\}$. Очевидно, что если $U_z \not\subset U_y$, то $(z, c) < (y, c)$. Допустим $U_z \subset U_y$, тогда U_z можно разбить на два непересекающихся подмножества $U_z^x = U_z \cap U_x$ и $U_z^y = U_z \setminus U_x$. Пусть $a = |U_z^x|$, $b = |U_z^y|$, причем $a \leq k_x$ и $b \leq k_y - k_x$. Тогда

$$(z, c) = \frac{a}{k_x} - \frac{b}{k_y - k_x} + \frac{b(b-1)}{(k_y - k_x)(k_y - k_x - 1)} = \frac{a}{k_x} - \frac{b(k_y - k_x - b)}{(k_y - k_x)(k_y - k_x - 1)}$$

Так как $z \neq x$ и $z \neq y$, то либо $a < k_x$, либо $0 < b < k_y - k_x$. Очевидно, что в первом случае $(z, c) \leq \frac{a}{k_x} < 1$, а во втором $(z, c) < \frac{a}{k_x} \leq 1$. Аналогично разбирается случай, когда $U_x = \emptyset$ и $(c, x) = (c, y) = 0$.

Теорема доказана.

Следствие 2. Плотность $p(X_n)$ графа $G(X_n)$ многогранника задачи максимальная клика равна мощности множества допустимых решений X_n : $p(X_n) = 2^n$, где n — число вершин графа $G(V, E)$ задачи МК.

Оценка плотности полиэдрального графа задачи коммивояжера

Основное внимание в этом пункте уделяется описанию алгоритма сведения оптимизационной задачи максимальная клика к несимметричной задаче коммивояжера, т. е. случаю, соответствующему ориентированному графу. Далее показано, что этот алгоритм удовлетворяет условию аффинной сводимости и, как следствие теоремы 2, выводится новая экспоненциальная оценка плотности полиэдрального графа для несимметричной задачи коммивояжера. В конце пункта осуществляется редукция результата для симметричной задачи.

Сформулируем несимметричную задачу коммивояжера (задача НК). Задан полный ориентированный реберно взвешенный граф $G'(V', A)$, где $V' = \{v'_1, v'_2, \dots, v'_k\}$ — множество вершин, $A = \{a_{ij} = (v'_i, v'_j), i \neq j, i, j \in \{1, k\}\}$ — множество дуг. Его дугам a_{ij} приписаны длины $c'_{ij} \in R$. Обозначим через \hat{Y} множество

всех гамильтоновых контуров этого графа (напомним, что гамильтонов контур — это контур, проходящий каждую вершину графа ровно один раз). Задача коммивояжера заключается в отыскании такого гамильтонова контура $\tilde{y}_0 \in \tilde{Y}$, суммарная длина дуг которого максимальна.

Чтобы построить многогранник этой задачи, сопоставим каждой дуге a_{ij} координату y_{ij} произвольной точки y в пространстве R^l , $l = k(k-1)$. Каждому гамильтонову контуру графа $G'(V', A)$ поставим в соответствие его характеристический вектор в R^l , положив равными единице те координаты, для которых соответствующие дуги входят в контур, остальные — равными нулю. Обозначим через Y_k множество всех таких векторов. Это множество допустимых решений задачи НК.

Пусть (X_n, Q_m) , где $m = n(n+1)/2$, — задача нахождения клики максимального веса (задача МК) в реберно и вершинно взвешенном неориентированном графе на n вершинах вес c_{ij} каждого ребра и каждой вершине которого удовлетворяет условию $-1 \leq c_{ij} \leq 1$.

Теорема 4. Задача (X_n, Q_m) аффинно сводится к несимметричной задаче коммивояжера Y_k для графа на $k = 2n^2 - n$ вершинах: $(X_n, Q_m) \propto_A Y_k$.

Доказательство. Пусть $G(V, E)$ — граф задачи МК, где $V = \{v_1, v_2, \dots, v_n\}$ — множество вершин, $E = \{e_{ij} = (v_i, v_j), 1 \leq i < j \leq n\}$. Его ребрам e_{ij} приписаны веса $c_{ij} \in R$, а вершинам v_i приписаны веса $c_{ii} \in R$. И пусть \tilde{X} — множество всех возможных клик этого графа, включающее в себя пустое множество. Построим ориентированный реберно-взвешенный граф $G'(V', A)$ для задачи НК, такой, что каждой клике $\tilde{x} \in \tilde{X}$ соответствует гамильтонов контур $\tilde{y} = B(\tilde{x})$ и гамильтонов контур \tilde{y}_0 является решением задачи НК тогда и только тогда, когда клика $\tilde{x}_0 = B^{-1}(\tilde{y}_0)$ является решением задачи МК.

Основная идея для построения графа G' заимствована из широко известного [1] алгоритма сведения *NP*-полной задачи ВЕРШИННОЕ ПОКРЫТИЕ к задаче ГАМИЛЬТОНОВ ЦИКЛ. Как и в монографии [1], наша конструкция будет состоять из набора компонент и вершин соединенных специально выбранными дугами.

Во-первых, граф G' имеет n "выбирающих" вершин a_1, a_2, \dots, a_n , которые будут использованы для выбора множества вершин $U \subseteq V$ графа G . Из каждой вершины a_i ($1 \leq i \leq n$) выходит ровно две дуги, назовем их "положительной" и "нулевой". "Положительную" дугу, начинающуюся в вершине a_i , обозначим b_{ii} , а множество всех таких дуг — $A^+ = \{b_{ii} : 1 \leq i \leq n\}$. Исходя из определения гамильтонова контура он может содержать только одну из дуг, берущих начало в общей вершине: либо b_{ii} , либо "нулевую". Далее будет показано, что гамильтонов контур, охватывающий все вершины конструируемого графа, однозначно определяется набором $P \subseteq A^+$ принадлежащих этому контуру "положительных" дуг.

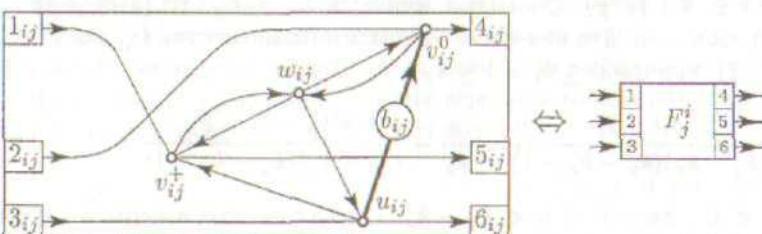


Рис. 1. Компонента F_j^i , соответствующая ребру e_{ij} из графа G

Во-вторых, для каждого ребра $e_{ij} \in E$ графа G' содержит компоненту F_j^i (см. рис. 1). Эта компонента состоит из четырех вершин $v_{ij}^+, v_{ij}^0, w_{ij}, u_{ij}$, которые соединены дугами

$$(v_{ij}^+, w_{ij}), (w_{ij}, v_{ij}^+), (v_{ij}^0, w_{ij}), (w_{ij}, v_{ij}^0), (w_{ij}, u_{ij}), (u_{ij}, v_{ij}^+), (u_{ij}, v_{ij}^0).$$

На рисунке дуга (u_{ij}, v_{ij}^0) обозначена b_{ij} . Эта компонента F_j^i соединена с другими компонентами и, возможно, с выбирающими вершинами дополнительными, внешними для данной компоненты дугами. На рисунке показаны части этих дуг:

- 1) 1_{ij} — конец дуги, входящей в вершину v_{ij}^+ ,
- 2) 2_{ij} — конец дуги, входящей в вершину v_{ij}^0 ,
- 3) 3_{ij} — конец дуги, входящей в вершину u_{ij} ,
- 4) 4_{ij} — начало дуги, выходящей из вершины v_{ij}^0 ,
- 5) 5_{ij} — начало дуги, выходящей из вершины v_{ij}^+ ,

6) b_{ij} — начало дуги, выходящей из вершины u_{ij} .

Причем в графе G' отсутствуют дуги, непосредственно соединяющие вершину w_{ij} с не принадлежащими компоненте F_j^i вершинами.

Эта компонента обладает следующими свойствами. Предположим, что для пары вершин v_{ij}^+, v_{ij}^0 в искомом гамильтоновом контуре присутствует ровно одна входящая в одну из этих вершин внешняя дуга. Тогда возможны четыре варианта прохождения гамильтонова контура через компоненту F_j^i (см. рис. 2). Перечислим дуги, входящие в гамильтонов контур в каждом из четырех случаев:

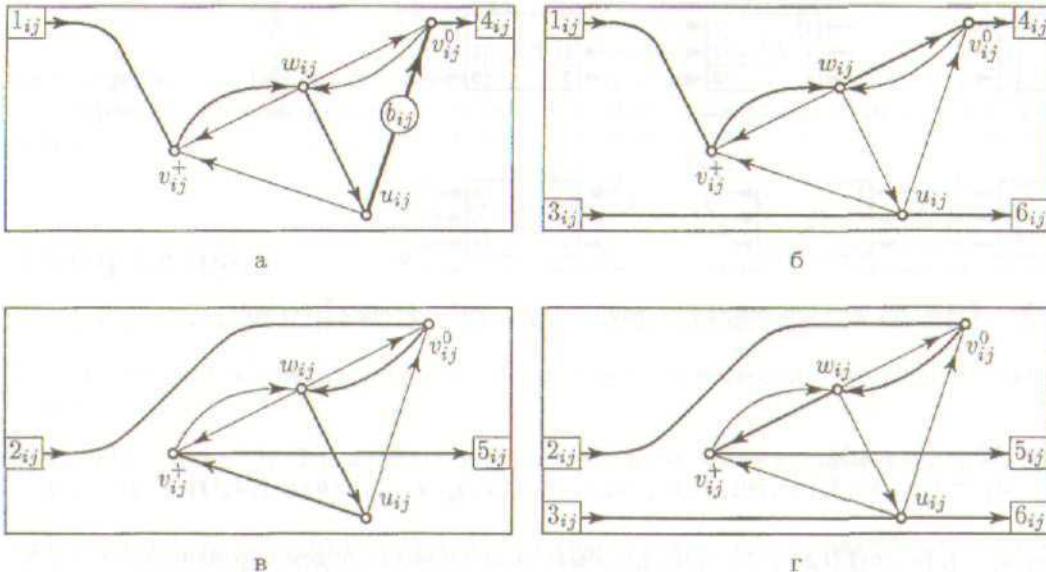


Рис. 2. Четыре возможные конфигурации прохождения гамильтонова контура через компоненту F_j^i

- a) $1_{ij}, (v_{ij}^+, w_{ij}), (w_{ij}, u_{ij}), b_{ij}, 4_{ij};$
- б) $1_{ij}, (v_{ij}^+, w_{ij}), (w_{ij}, v_{ij}^0), 4_{ij}, 3_{ij}, 6_{ij};$
- в) $2_{ij}, (v_{ij}^0, w_{ij}), (w_{ij}, u_{ij}), (u_{ij}, v_{ij}^+), 5_{ij};$
- г) $2_{ij}, (v_{ij}^0, w_{ij}), (w_{ij}, v_{ij}^+), 5_{ij}, 3_{ij}, 6_{ij};$

Причем дуга b_{ij} входит в гамильтонов контур тогда и только тогда, когда в нем отсутствуют дуги 2_{ij} и 3_{ij} . Если изобразить компоненту F_j^i в виде прямоугольника (см. рис. 1), с тремя расположенными по порядку сверху вниз входами 1, 2, 3, и тремя аналогично расположенным выходами 4, 5, 6, то в любом из четырех описанных вариантов набор входящих извне в компоненту дуг будет симметричен набору выходящих.

Рассмотрим все компоненты F_j^l , соответствующие ребрам e_{ij} графа G , инцидентным одной и той же вершине v_i (т. е. $l = i$ или $j = i$). Два набора дуг A_i^+ и A_i^0 соединяют эти компоненты в "гирлянду", начинающуюся в дополнительной вершине a_i и заканчивающуюся в вершине a_{i+1} ; в случае $i = n$ начало "гирлянды" — вершина a_n , а конец — вершина a_1 (рис. 3):

1) Если $i = 1$, то

$$\begin{aligned} A_1^+ &= \{(a_1, v_{12}^+), (v_{12}^0, v_{13}^+), (v_{13}^0, v_{14}^+), \dots, (v_{1_{n-1}}^0, v_{1_n}^+), (v_{1_n}^0, a_2)\}, \\ A_1^0 &= \{(a_1, v_{12}^0), (v_{12}^+, v_{13}^0), (v_{13}^+, v_{14}^0), \dots, (v_{1_{n-1}}^+, v_{1_n}^0), (v_{1_n}^+, a_2)\}. \end{aligned}$$

2) Если $1 < i < n$, то

$$\begin{aligned} A_i^+ &= \{(a_i, v_{i+1}^+), (v_{i+1}^0, v_{i+2}^+), (v_{i+2}^0, v_{i+3}^+), \dots, (v_{i_{n-1}}^0, v_{i_n}^+), (v_{i_n}^0, a_{i+1})\}, \\ A_i^0 &= \{(a_i, u_{1_i}), (u_{1_i}, u_{2_i}), (u_{2_i}, u_{3_i}), \dots, (u_{i-2_i}, u_{i-1_i}), (u_{i-1_i}, v_{i+1}^0), \\ &\quad (v_{i+1}^+, v_{i+2}^0), (v_{i+2}^+, v_{i+3}^0), \dots, (v_{i_{n-1}}^+, v_{i_n}^0), (v_{i_n}^+, a_{i+1})\}. \end{aligned}$$

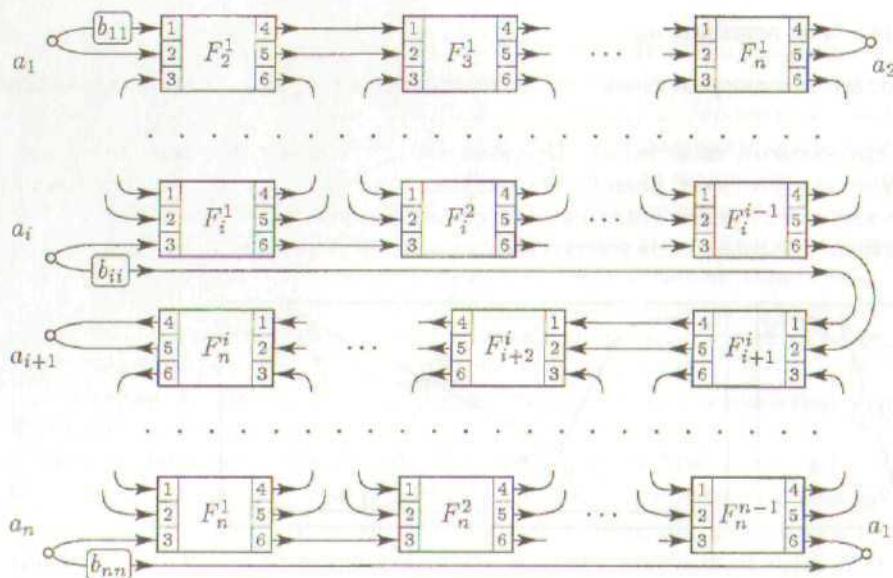


Рис. 3. Схема соединения всех компонент и "выбирающих" вершин

3) Если $i = n$, то

$$\begin{aligned} A_n^+ &= \{(a_n, a_1)\}, \\ A_n^0 &= \{(a_n, u_{1_n}), (u_{1_n}, u_{2_n}), (u_{2_n}, u_{3_n}), \dots, (u_{n-2_n}, u_{n-1_n}), (u_{n-1_n}, a_1)\}. \end{aligned}$$

Дуги $(a_1, v_{1_2}^+), (a_i, v_{i+1}^+), (a_n, a_1)$ (где $1 < i < n$), как было сказано выше, будем обозначать b_{11}, b_{ii}, b_{nn} . Учитывая особенности конструкции компонент F_j^i , описанные выше, несложно проверяется, что дуга b_{ii} ($1 \leq i \leq n$) входит в гамильтонов контур тогда и только тогда, когда в него входят все дуги из набора A_i^+ и не входит ни одна дуга из набора A_i^0 . Таким образом, гамильтонов контур, обходящий все вершины указанного графа, однозначно задается набором дуг $P \subseteq A^+$. Причем, как было отмечено выше, дуга b_{ij} из компоненты F_j^i входит в искомый контур только в том случае, когда в этом контуре отсутствуют дуги, входящие "извне" в вершины u_{ij} и v_{ij}^0 , т. е. дуги, принадлежащие наборам A_i^0 и A_j^0 . А это возможно тогда и только тогда, когда в "задающий" набор P входят дуги b_{ii} и b_{jj} . Следовательно, дуга b_{ij} входит в гамильтонов контур, если и только если в нем входят дуги b_{ii} и b_{jj} .

Допустимое решение задачи МК однозначно определяется набором вершин $U \subseteq V$. Гамильтонов контур в описанной конструкции однозначно задается набором "положительных" дуг $P \subseteq A^+$. Взаимно однозначное соответствие B между двумя этими наборами определим следующим образом. Вершине v_i из графа G будет соответствовать дуга b_{ii} из графа G' . Тогда аффинный оператор A из определения аффинной сводимости задается следующей системой уравнений:

- 1) Координата целевого вектора c' , соответствующая дуге b_{ij} ($1 \leq i \leq j \leq n$) графа G' , равна координате c_{ij} (весу ребра e_{ij} , если $i \neq j$, или весу вершины v_i , если $i = j$).
- 2) Координата целевого вектора c' , соответствующая дуге, входящей в нашу конструкцию, но не равной b_{ij} ($1 \leq i \leq j \leq n$), равна нулю.
- 3) Все остальные координаты вектора c' равны $-n^2$. Эта величина гарантирует отсутствие этих дуг в гамильтоновом контуре максимального веса, т. к. по условию $-1 \leq c_{ij} \leq 1$.

Очевидно, уравнения первого типа образуют единичную матрицу C из определения аффинной сводимости, причем $\text{rang } C = m$. А уравнения второго и третьего типа задают вектор d , координаты которого равны нулю, либо $-n^2$. Итак, взаимно однозначное отображение B и невырожденное аффинное отображение A из определения аффинной сводимости известны. Условие (1) выполнено. Число вершин k графа G' равно учетверенному числу компонент F_j^i , где $1 \leq i < j \leq n$, плюс n "выбирающих" вершин a_i . Теорема доказана.

Теперь, пользуясь следствием 2 о том, что плотность полиэдрального графа задачи МК равна 2^n , получаем более высокую, чем известная ранее [2], нижнюю оценку плотности полиэдрального графа для

несимметричной задачи коммивояжера на k вершинах:

$$p(Y_k) > 2^{\sqrt{k/2}-1}.$$

Используя этот результат, можно установить новую экспоненциальную нижнюю оценку плотности $p(Y'_t)$ полиздрального графа симметричной задачи коммивояжера Y'_t для t городов. Для этого достаточно воспользоваться известным [2] соотношением плотностей симметричной и несимметричной задач:

$$p(Y'_t) \geq p(Y_k), \text{ при } k = \left[\frac{t}{2} \right],$$

где квадратные скобки означают выделение целой части.

Теорема 5. Для плотности графа многогранника симметричной задачи коммивояжера справедлива оценка

$$p(Y'_t) > 2^{\sqrt{t/2}-2}.$$

Литература

- [1] Гэри М., Джонсон Д. Вычислительные машины и труднорешаемые задачи. М.: Мир, 1982. 416 с.
- [2] Бондаренко В.А. Полиздральные графы и сложность в комбинаторной оптимизации. Ярославль, 1995. 126 с.
- [3] Papadimitriou C. H. The adjancency relation on the traveling salesman polytope is NP-complete // Math. Prog. 1978. P. 312-324.
- [4] Бондаренко В.А. Неполиномиальная нижняя оценка сложности задачи коммивояжера в одном классе алгоритмов // Автоматика и телемеханика. 1983. N 9. С. 45-50.
- [5] Емеличев В.А., Ковалев М.М., Кравцов М.К. Многогранники, графы, оптимизация. М.: Наука, 1981. 344 с.

Моделирование протоколов TCP и ARTCP с помощью раскрашенных сетей Петри

Чалый Д.Ю.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

E-mail: chaly@uniyar.ac.ru

получена 14 марта 2003

В работе рассматривается проблема моделирования транспортных протоколов TCP (Transmission Control Protocol) и ARTCP (Adaptive Rate TCP). В качестве средства моделирования предлагается использовать раскрашенные сети Петри. В работе представлена оригинальная модель и рассмотрено несколько примеров того, как могут быть проверены качественные (корректность) и количественные (производительность) свойства протоколов.

1. Введение

Набор протоколов TCP/IP является важнейшим элементом для обеспечения коммуникации между устройствами, подключенными к сети Интернет. Этот набор протоколов позволяет обмениваться данными между компьютерами, имеющими разную архитектуру, и на которых выполняются различные операционные системы. Протокол Управления Передачей (Transmission Control Protocol – TCP) является важнейшим протоколом транспортного уровня в системе протоколов TCP/IP. Он обеспечивает надёжную дуплексную передачу данных от одного хоста к другому с учётом механизма управления потоком. С 1981 года, когда была опубликована первая спецификация протокола TCP [1], было сделано много усовершенствований протокола и исправлен ряд ошибок. Самыми важными спецификационными документами, в которых описаны эти изменения, являются: [2], содержащий ряд исправлений и устанавливающий стандарт протокола; [3], посвящённый производительности протокола на высокоскоростных каналах связи; [4], описывающий алгоритм выборочных подтверждений (Selective Acknowledgements – SACK), чтобы более эффективно справляться с потерями сегментов; [5], где описаны стандартные алгоритмы управления потоком; [7], предлагающий алгоритм Ограниченней Передачи (Limited Transmit algorithm) – более эффективное средство для восстановления от потерь сегментов.

Большое число научных работ было посвящено исследованию различных аспектов работы протокола TCP. В работе [20] Кумар использовал стохастическую модель для исследования производительности различных версий протокола TCP, предполагая наличие случайных потерь сегментов в беспроводном канале связи. В работе [19] Фолл и Флойд с помощью имитационной модели исследовали преимущества алгоритма выборочных подтверждений. Для исследования протокола TCP в работе [16] использовались раскрашенные сети Петри, однако представленная там модель была достаточно упрощенной, т.е. требовалась более точная реализация некоторых стандартных алгоритмов работы протокола TCP (например, расчёт значения таймера повторной передачи). Другим недостатком этой модели является неспособность моделировать систему, состоящую из нескольких по-разному сконфигурированных соединений, без существенной реконструкции модели.

Для построения модели протокола TCP мы использовали иерархические раскрашенные сети Петри. Представляемая оригинальная модель реализует алгоритмы, которые были описаны в документах, представленных выше. Необходимо отметить, что мы моделировали спецификацию протокола, а не какую-то реализацию в отдельно взятой операционной системе, т.к. реализации довольно существенно отличаются друг от друга и не всегда полностью отвечают требованиям стандартов. Кроме моделирования протокола TCP, модель была расширена с целью моделирования протокола ARTCP [8, 9, 10].

Для построения и исследования модели мы использовали систему моделирования Design/CPN [15, 13]. Раскрашенные сети Петри зарекомендовали себя как хороший формализм для моделирования и анализа свойств распределённых систем (в том числе и коммуникационных протоколов) и были использованы в ряде успешных проектов, например [17, 18, 16]. В дальнейшем мы будем полагать, что читатель знаком

с базовыми концепциями иерархических раскрашенных сетей Петри [11, 14, 12] и основными аспектами работы протоколов TCP и ARTCP.

2. Модель протокола TCP

В данной статье мы не будем рассматривать модель в деталях, так как протокол имеет довольно сложную структуру и, как следствие, модель тоже. На рисунке 1 показан фрагмент нашей модели – подсеть Processing, которая моделирует обработку пришедших сегментов. Подсеть представляет собой двудольный ориентированный граф с двумя типами вершин – позициями (изображены овалами) и переходами (единственный переход на рисунке изображён прямоугольником). Рассматриваемая подсеть содержит позиции, где находятся маркеры, моделирующие различные служебные структуры протокола. Например, позиция TCB содержит маркеры, задающие контрольные блоки соединений – структуры, которая во многом отвечает за работу протокола; позиция SegBuffer содержит маркеры, которые задают буферы, в которых содержатся принесенные сегменты. Кроме позиций, подсеть содержит переход Process, который может сработать, когда существует набор маркеров во входных позициях (т.е. позициях, от которых идёт дуга к переходу), которые удовлетворяют охранному выражению (на рисунке – курсивом в квадратных скобках). При срабатывании перехода, из входных позиций перехода убираются соответствующие маркеры, а в выходные позиции (те, к которым идёт дуга от перехода), в зависимости от значений выражений на выходных дугах перехода, помещаются новые маркеры. Эти новые маркеры моделируют изменённые служебные структуры. Так как вычисление новых значений может быть довольно сложным, мы вынесли процесс вычисления в кодовый сегмент перехода (изображён на рисунке пунктиром прямоугольником с буквой С внутри). В кодовом сегменте вычисляются значения переменных, которые используются на выходных дугах перехода.

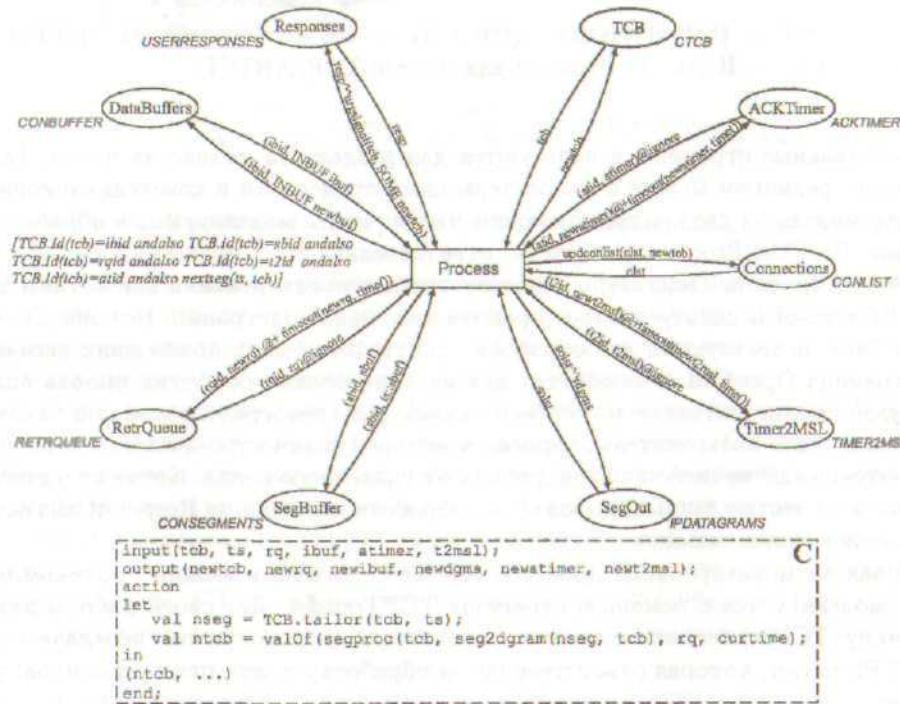


Рис. 1. Подсеть Processing – моделирование обработки полученных сегментов

Наша модель состоит из двух базовых частей. Первая часть представляет собой собственно структуру сети Петри. Выражения на дугах и охраны переходов могут содержать довольно сложные выражения, написанные на языке CPN ML. Ряд выражений мы оформили как функции и вынесли во внешние файлы, которые подключаются к модели (например, на рисунке 1 это функция procseg в кодовом сегменте). Объявления цветов и эти файлы представляют собой вторую часть нашей модели. Структура модели может быть представлена как иерархия подсетей, каждая из которых моделирует отдельный аспект работы протокола и расположена на отдельной странице. Иерархия страниц нашей модели показана на рисунке 2. Она представляет собой дерево, где листьевые страницы моделируют некоторые базовые аспекты

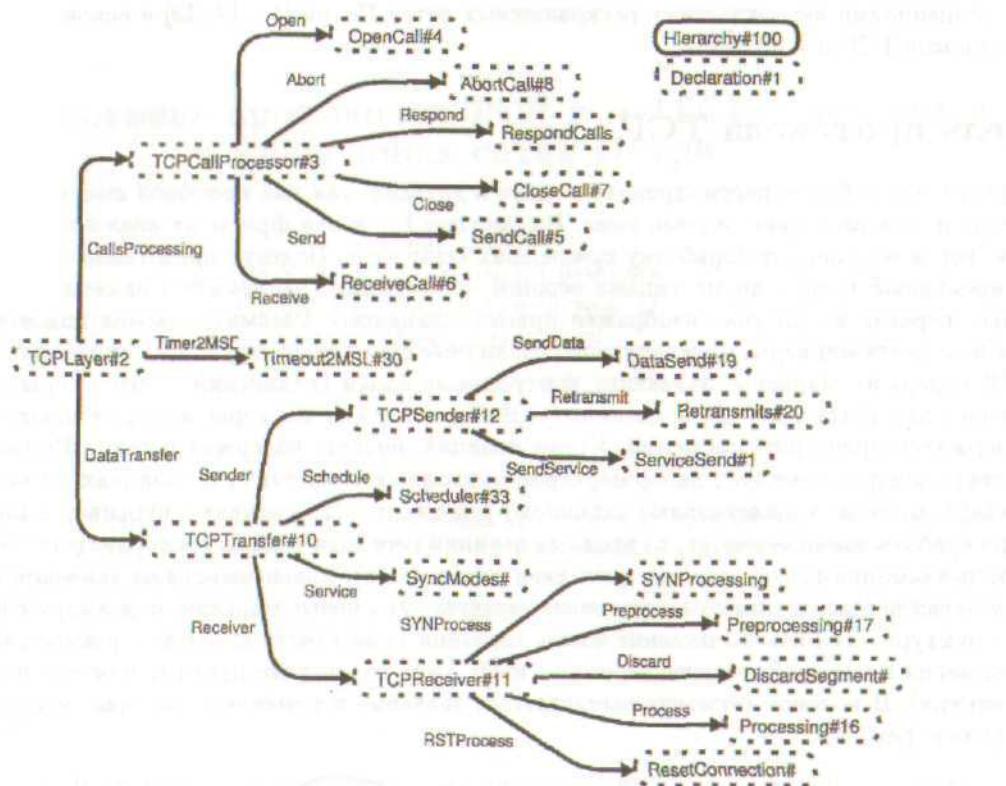


Рис. 2. Иерархия страниц для модели TCP/ARTCP

работы протокола, а остальные страницы используются для разделения модели на части. Так как протокол TCP выступает посредником между пользовательскими процессами и коммуникационной средой, было решено разделить модель на две соответствующие части: часть, моделирующая обработку вызовов пользователя (страница **TCPCallProcessor**), и часть, ответственная за отправление и приём сегментов (страница **TCPTransfer**). Страница **Timer2MSL** используется для моделирования 2MSL-тайм-аута.

Страница **TCPCallProcessor** использует в своей работе несколько подстраниц. Все они, за исключением страницы **RespondCalls**, используются для обработки соответствующих пришедших вызовов пользователя. Например, страница **OpenCall** используется для моделирования обработки вызова пользователя **OPEN**. Некоторые вызовы пользователя не могут быть обработаны непосредственно по приходу. Это может произойти, например, если пользователь запросил некоторое количество данных с помощью вызова **RECEIVE**, однако протокол ещё не получил этих данных от удалённого конца. Когда от удалённого конца придёт необходимое количество данных, вызов будет обработан. Страница **RespondCalls** используется для обработки таких отложенных вызовов.

Часть, ответственная за моделирование процесса обмена сегментами между протоколом и коммуникационной средой, моделируется с помощью страницы **TCPTransfer**. Для своей работы эта страница использует: подстраницу **TCPSender**, которая используется для моделирования передачи сегментов в сеть; подстраницу **TCPReceiver**, которая ответственная за обработку пришедших сегментов; служебную подстраницу **Scheduler**.

Страница **TCPSender** для своей работы моделирует отправление сегментов с данными (подстраница **DataSend**), передачу служебных сегментов – подтверждений или SYN-сегментов для установки соединений (подстраница **ServiceSend**) и повторную передачу сегментов, которые потеряны сетью (подстраница **Retransmits**).

Обработка пришедших сегментов моделируется с помощью страницы **TCPReceiver**. Для своей работы она использует следующие страницы: **SYNProcessing**, которая моделирует обработку пришедших SYN-сегментов; **Preprocessing**, которая моделирует начальную обработку сегментов, например, тесты по прибытию; **ResetConnection**, которая используется для обработки пришедших RST-сегментов; **Processing**, которая используется для обработки корректных сегментов. Сервисная страница **SyncModes** используется для установления корректных приоритетов различным действиям, совершаемым протоколом. Например, если обнаруживается что в данный момент времени может быть отправлен сегмент с новыми данными и

должна быть произведена повторная передача сегмента, то выполниться должно второе действие. Сервисная страница Scheduler используется протоколом ARTCP для моделирования отправления сегментов с заданной скоростью.

Представляемая модель отвечает требованиям, накладываемым на неё спецификационными документами, которые были процитированы в предыдущем разделе. Модель была разработана с учётом возможностей простой конфигурации и настройки. Существует возможность задания значений различных параметров, которые используются протоколом TCP. Это сделано путём создания отдельного конфигурационного файла, который содержит набор этих параметров. Другая возможность настройки модели состоит в том, что существует возможность создать модель системы, в которой одновременно работают несколько соединений с различной конфигурацией. Например, возможно промоделировать систему, в которой ряд соединений используют стандартный алгоритм управления потоком TCP, а остальные используют алгоритм управления потоком ARTCP.

Другой особенностью модели является относительно простая возможность изменять её поведение без существенной реконструкции сетевой структуры. Поясним способ, которым это можно сделать на примере уже рассмотренной сети Processing. Допустим, если какая-то новая модификация требует иной обработки сегмента, то необходимо лишь изменить код соответствующих функций, например `procseg`. В итоге мы получим модель, которая функционирует иным образом, чем исходная и моделирует уже необходимый нам алгоритм (конечно, если изменения, произведённые в CPN ML-коде функций, адекватны).

По заданной иерархической сети Петри можно построить эквивалентную сеть без иерархических конструкций. Это построение было проведено (с помощью средств пакета Design/CPN). Получилось, что эта сеть без иерархических конструкций состоит из 15 позиций и 30 переходов. Таким образом, модель является довольно компактной, несмотря на то, что моделируемая система довольно сложна.

3. Некоторые результаты исследования модели

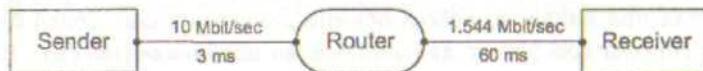


Рис. 3. Схема сетевой структуры, использовавшейся для экспериментов.

В этом разделе мы представим некоторые экспериментальные результаты для демонстрации возможностей анализа. Для анализа мы использовали прикладной пакет Design/CPN версии 4.0.5. Чтобы проанализировать протокол, мы построили подсети, которые моделируют канал связи и простейший маршрутизатор. Используя эти подсети как компоненты, можно построить довольно сложную сетевую структуру. Сетевая структура, использовавшаяся для экспериментов, показана на рисунке 3. На этой схеме представлены отправитель данных (Sender) и получатель данных (Receiver). Отправитель и получатель связаны друг с другом каналами связи через маршрутизатор (Router). Значение скоростей передачи каналов и задержки передачи соответствуют значениям для Ethernet (канал от отправителя до маршрутизатора) и телефонной линии T1 (канал от маршрутизатора до получателя). Предполагается, что данные на линии не искажаются и не теряются. Маршрутизатор для каждого канала связи имеет отдельный буфер (характер буфера – очередь FIFO) конечного размера, равный 32000 байтов. Если буфер переполняется, то пришедшие сегменты отбрасываются. Во всех экспериментах максимальный размер сегмента равен 1000 байтов.

Рассмотрим сначала обнаруженную в протоколе тупиковую ситуацию (deadlock - дедлок). Для простоты, оба конца соединения не используют никакие-либо расширения протокола, например опцию временной метки или алгоритм SACK. Для предотвращения перегрузки получателя протокол TCP использует стандартный алгоритм скользящего окна. Протокол может отправить данные, если они "помещаются" в это окно. Протокол имеет рекомендуемый алгоритм для определения размера данных, которые должны быть отправлены, описанный в [2]. Согласно этому алгоритму протокол определяет количество данных, которые можно отправить следующим образом:

- может быть отправлен сегмент максимального размера;
- данные помечены флагом PUSH и все отправленные данные могут быть отправлены;
- по крайней мере доля максимального окна может быть отправлена (рекомендуемое значение доли равно 1/2);

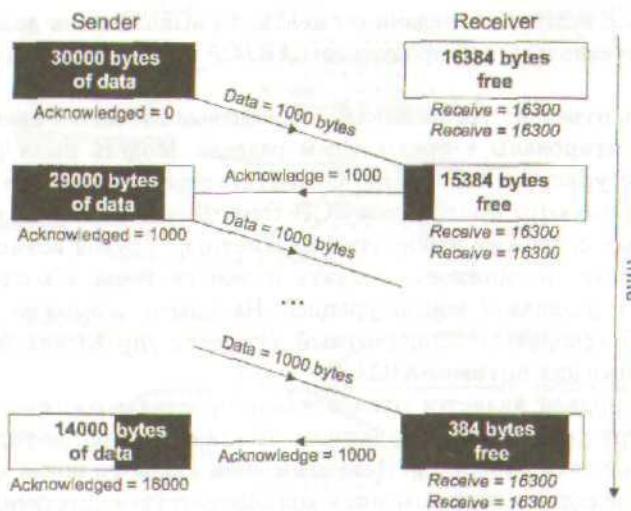


Рис. 4. Процесс передачи данных

- данные помечены флагом PUSH и срабатывает "перекрывающий" (override) тайм-аут;

Так как "перекрывающий" тайм-аут нигде не рассматривается в спецификации протокола TCP, то мы его не реализовали в нашей модели. Однако он не влияет на предлагаемый пример, так как отправляемые данные не помечаются флагом PUSH. Отправитель пытается переслать 30000 байтов получателю. Получатель имеет буфер размером 16384 байта (максимальный размер окна получателя равен этому же значению). Пришедшие данные помещаются в этот буфер а потом передаются пользовательскому процессу. Окно получателя при начале процесса обмена данными равно 16384 байта. Пользовательский процесс с получающей стороны совершает два запроса на получение данных. Каждый запрос требует 16300 байтов данных. Процесс передачи данных в рассматриваемом примере показан на рисунке 4. Отправитель прекратит передавать сегменты с новыми данными после отправления 16 полноразмерных сегментов получателю, так как ни одно из условий для определения количества новых данных, которое можно отправить, не выполняется. Максимальная порция данных, равная максимальному размеру сегмента, не может быть отправлена, так как получатель не имеет достаточно буферного пространства, чтобы поместить её (порция, равная доле окна не отправляется по этой же причине). Условия, которые имеют дело с данными, помеченными флагом PUSH, не срабатывают, так как отправляемые данные им не помечены. Таким образом, отправитель не может определить количество данных, которые необходимо отправить получателю, а получатель не может передать данные из буфера пользовательскому процессу, так как в буфер помещено недостаточное количество данных. Для избегания данного дедлока предлагается добавить к существующим следующее условие: количество данных для отправления равно размеру удалённого окна, если все отправленные ранее данные были подтверждены получателем, размер удалённого окна меньше максимального размера сегмента и количество данных, помещенных в буфер отправителя, больше или равно удалённому окну. Изначально данный дедлок был обнаружен с помощью симуляции (выполнения) модели. После его нахождения было построено множество достижимых состояний системы (всего у данной системы 364 состояния), с помощью которого было установлено, что тупиковое состояние является единственным финальным состоянием системы. То есть для любого пути исполнения за конечное число шагов наступит данное тупиковое состояние. Далее в экспериментах мы будем рассматривать подкорректированную модель.

Рассмотрим теперь эксперимент, который показывает как могут быть определены некоторые количественные характеристики рассматриваемой системы. Для увеличения производительности увеличим размер буфера получателя до 60000 байтов, а также включим алгоритмы временной метки и выборочных подтверждений. Объём данных, которые пытается передать отправитель, равен 250000 байтов. Полученный результат показан на рисунке 5 (на первых двух парах рисунков по оси у показаны значения порядковых номеров данных в сегментах по модулю 60000). Первая пара графиков показывает процесс приёма сегментов получателем. Вторая пара показывает процесс передачи сегментов отправителем. Третья пара графиков показывает сколько места в буфере маршрутизатора занимают сегменты с данными, которые передаёт отправитель.

На графиках видно, что при использовании протокола TCP в данной системе потерянся 4 сегмента

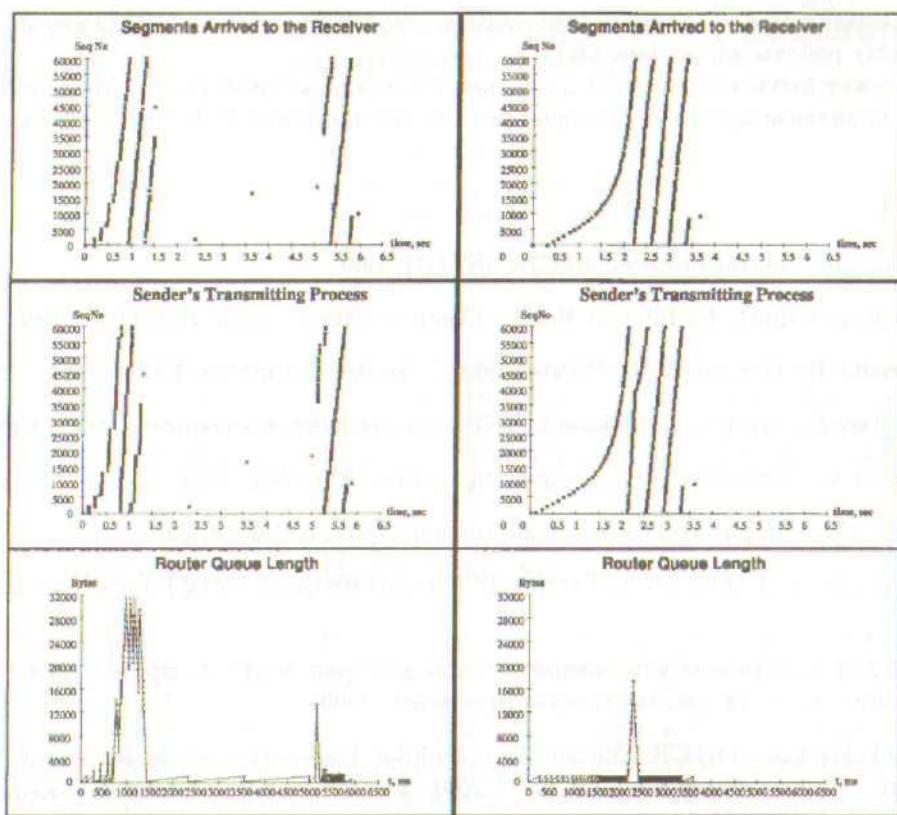


Рис. 5. Качественные характеристики системы для TCP(слева) и ARTCP(справа).

(это те сегменты, на которых передача прерывается). Первый и четвертый из потерянных сегментов будут переданы повторно с помощью механизма быстрой ретрансляции (fast retransmit). Остальные два будут переданы с помощью стандартного механизма повторной передачи. На третьем графике (слева) видно, что в моменты времени, когда происходят потери сегментов, буфер маршрутизатора заполнен. Видно, что протокол передаст нужное количество данных примерно за 6 секунд.

Рассмотрим теперь работу протокола ARTCP. Как видно, этот протокол не создаёт перегрузку буфера маршрутизатора, следовательно, потеря сегментов по этой причине не происходит. В начале своей работы протокол ARTCP резко увеличивает скорость отправления сегментов в сеть. Когда эта скорость больше свободной пропускной способности сети, происходит резкое заполнение буфера маршрутизатора (пик на третьем справа рисунке). Все остальное время буфер маршрутизатора содержит не более одного сегмента. Процесс передачи данных завершится примерно за 3.5 секунды.

Таким образом видно, что на данном примере механизм управления потоком ARTCP работает эффективнее стандартного механизма TCP. Другим преимуществом ARTCP над TCP является значительно меньшая потребность в буферном пространстве маршрутизатора (если сравнивать пики – то примерно в два раза).

4. Заключение

Мы представили смешанную модель протоколов TCP/ARTCP. Стоит отметить, что производилось моделирование спецификации протокола, а не какой-то избранной реализации. Ряд реализаций не соответствуют стандартной спецификации, однако наша модель может быть реконфигурирована, чтобы представлять их. Также модель может быть полезной при исследовании различных свойств будущих модификаций протоколов TCP и ARTCP. Также возможно использовать модель для анализа других механизмов управления потоком, которые, как и протокол ARTCP, управляют скоростью передачи данных (например, механизм TCP Friendly Rate Control).

В работе приведены несколько примеров того, как различные виды свойств протокола – качественные (корректность) и количественные (производительность) могут быть проанализированы. Будущая

работа будет посвящена более детальному исследованию свойств протокола TCP (и его модификаций), в особенности, анализу работы алгоритма ARTCP.

Наша модель может быть использована не только для анализа семейства протоколов TCP, но и как составная часть для анализа программ, которые используют протокол TCP для коммуникации.

Литература

- [1] Postel J. Transmission Control Protocol. RFC793 (STD7), 1981.
- [2] Braden R. (ed.) Requirements for Internet Hosts – Communication Layers. RFC1122, 1989.
- [3] Jacobson V., Braden R., Borman D. TCP Extensions for High Performance. RFC1323, 1992.
- [4] Mathis M., Mahdavi J., Floyd S., Romanow A. TCP Selective Acknowledgement Option RFC2018, 1996.
- [5] Allman M., Paxson V., Stevens W. TCP Congestion Control. RFC2581, 1999.
- [6] Paxson V., Allman M. Computing TCP's Retransmission Timer. RFC2988, 2000.
- [7] Allman M., Balakrishnan H., Floyd S. Enhancing TCP's Loss Recovery Using Limited Transmit. RFC3042, 2001.
- [8] Алексеев И.В. Адаптивная схема управления потоком для транспортного протокола в сетях с коммуникацией пакетов: Дисс. ... канд. физ.-мат. наук. Ярославль, 2000.
- [9] Alekseev I.V., Sokolov V.A. ARTCP: Efficient Algorithm for Transport Protocol for Packet Switched Networks //Malyshkin V. (ed.) Proceedings of PaCT'2001. Lecture Notes in Computer Science, Vol. 2127. Springer-Verlag, 2001. P.159–174.
- [10] Alekseev I.V., Sokolov V.A. Modelling and Traffic Analysis of the Adaptive Rate Transport Protocol //Future Generation Computer Systems. Vol. 18. No.6. 2002. P.813–827
- [11] Jensen K. Coloured Petri Nets. Basic Concepts, Analysis Methods and Practical Use. Vol 1, 2, 3. Springer-Verlag, 1992.
- [12] Jensen K., Rozenberg G. (eds.) High-Level Petri Nets. Springer-Verlag, 1991.
- [13] Christensen S., Jørgensen J.B., Kristensen L.M. Design/CPN - A Computer Tool for Coloured Petri Nets //Brinksma E. (ed.) Proceedings of TACAS'97. Lecture Notes in Computer Science. Vol. 1217. Springer-Verlag, 1997. P.209–223.
- [14] Coloured Petri Nets. World-Wide Web. <http://www.daimi.aau.dk/CPnets>.
- [15] Design/CPN Online. World-Wide Web. <http://www.daimi.au.dk/designCPN/>.
- [16] de Figueiredo J.C.A., Kristensen L.M. Using Coloured Petri Nets to Investigate Behavioural and Performance Issues of TCP Protocols //Jensen K. (ed.) Proceedings of the Second Workshop on Practical Use of Coloured Petri Nets and Design. 1999. P.21–40.
- [17] Clausen H., Jensen P.R. Validation and Performance Ananlysis of Network Algorithms by Coloured Petri Nets //Proceedings of PNPM'93. IEEE Computer Society Press. 1993. P.280–289.
- [18] Clausen H., Jensen P.R. Analysis of Usage Parameter Control Algorithm for ATM Networks //Tohmè S. and Casada A. (eds.) Broadband Communications II (C-24). Elsevier Science Publishers. 1994. P.297–310.
- [19] Fall K., Floyd S. Simulation-Based Comparisons of Tahoe, Reno, and SACK TCP //Computer Communication Review. V.26(3). 1996. P.5–21.
- [20] Kumar A. Comparative Performance Ananlysis of Versions of TCP in a Local Network with a Lossy Link // IEEE/ACM Transactions on Networking. V.6(4). 1998. P.485–498.

УДК 517.9 + 519.86

Об одной нелинейной краевой задаче, моделирующей экономические циклы

Косарева Е.С., Куликов А.Н.

Ярославский государственный университет
150 000, Ярославль, Советская, 14

получена 14 апреля 2003

В работе рассмотрена одна краевая задача, предложенная Т. Пу в качестве модели экономических циклов с учетом пространственного взаимодействия. Показано, что учет пространственных факторов существенным образом меняет динамику экономических циклов известной модели мультипликатора – акселератора. В частности, обнаружена возможность реализации экономических циклов разного периода.

1. Постановка задачи. Выбор модели

Предлагаемая ниже математическая модель приведена в монографии Т. Пу [1] как одна из возможных модификаций широко известной математической модели циклов деловой активности. Этот вариант или достаточно близкие к нему широко известны под названием модель мультипликатора – акселератора. Ниже приводится вариант, предложенный Алленом [2] в 1956 году и, в значительной мере считающийся среди экономистов достаточно устоявшимся, если не классическим. Следует отметить, что математические модели мультипликатора – акселератора при любых вариантах основываются на макроэкономическом подходе Кейнса [3].

Обозначим доход через $Y = Y(t)$. Сбережения $S = S(t)$ находятся в заданном отношении s к доходу. Сохраняется неизменное отношение между основным капиталом $K(t)$ и доходом, v обозначает коэффициент пропорциональности. Инвестиции I по определению являются темпами изменения основного капитала. Таким образом, $I = v\dot{Y}$ и $S = s\dot{Y}$. Эти еще достаточно простые и иногда не бесспорные положения привели к тому, что динамика дохода, как основного показателя экономического процесса, может быть описана дифференциальным уравнением

$$\ddot{Y} - (v - 1 - s)\dot{Y} + sY = -\frac{v}{3}\dot{Y}^3. \quad (1)$$

Последнее уравнение – это широко известное уравнение Релея. Так, например, при $0 < v - 1 - s \ll 1$ уравнение (1) имеет единственный асимптотически устойчивый предельный цикл. Этот цикл и интерпретируется обычно в экономике как цикл "деловой активности".

Наиболее очевидным недостатком этой модели является то обстоятельство, что она является "точечной", не учитывает пространственное взаимодействие, например, межрегиональную торговлю.

Простейшая модель уже с учетом пространственного взаимодействия имеет вид

$$\frac{\partial^2 u}{\partial t^2} - 2\varepsilon \frac{\partial u}{\partial t} - a^2 \frac{\partial^2 u}{\partial x^2} + \omega^2 u = -\left(\frac{\partial u}{\partial t}\right)^3, \quad (2)$$

где $u = u(t, x)$, $x \in [0, \pi]$. Уравнение (1) (см.[1,4]) записано уже в перенормированном виде $t \rightarrow \alpha t$, $x \rightarrow \beta x$, $Y \rightarrow \gamma u$. Коэффициент a^2 отражает склонность к импорту [1,4]. Иаконец, 2ε – величина пропорциональная разности $v - 1 - s$. В рамках данной работы будем предполагать, что $0 < \varepsilon \ll 1$.

Уравнение в частных производных следует дополнить краевыми условиями.

В [1,4] предлагаются к рассмотрению, как обычно, три их варианта

$$u(t, 0) = u(t, \pi) = 0, \quad (3)$$

$$\left. \frac{\partial u}{\partial x} \right|_{x=0} = \left. \frac{\partial u}{\partial x} \right|_{x=\pi} = 0 \quad (4)$$

и, наконец, считая, что уже $x \in [0, 2\pi]$

$$u(t, x + 2\pi) = u(t, x); \quad \frac{\partial u(t, x + 2\pi)}{\partial x} = \frac{\partial u(t, x)}{\partial x} \quad (5)$$

В предлагаемой работе более детально остановимся на первом варианте: на исследовании краевой задачи (2), (3).

Произведение $\dot{W}_2^2 \times \dot{W}_2^1$ соболевских пространств – естественное фазовое пространство краевой задачи (2), (3). Напомним, что \dot{W}_2^k это замыкание достаточно гладких функций, удовлетворяющих краевым условиям по норме пространства $L_2(0, \pi)(k \in \mathbb{N})$.

2. Квазинормальная форма краевой задачи

Для исследования динамики краевой задачи (2), (3) полезно построить квазинормальную форму. Напомним способ ее построения (см. [5], а также работы, ссылки на которые содержатся в этой работе). Для этого следует прежде всего ввести в рассмотрение функцию

$$u_0(t, \tau, x) = \sum_{n=1}^{\infty} (z_n(\tau) \exp(i\omega_n t) + \bar{z}_n(\tau) \exp(-i\omega_n t)) \sin nx, \quad (6)$$

где $\tau = \varepsilon t$, $\omega_n^2 = a^2 n^2 + \omega^2$ – собственные числа дифференциального оператора $-a^2 \frac{d^2}{dx^2} + \omega^2 I$ с граничными условиями (3). Затем положим

$$u = \varepsilon^{\frac{1}{2}} u_0 + \varepsilon^{\frac{3}{2}} u_1 + \dots \quad (7)$$

где $u_1 = u_1(t, \tau, x)$ и т.д. Подставим (7) в краевую задачу (2), (3). Приравнивая в получившихся формальных равенствах коэффициенты при $\varepsilon^{\frac{3}{2}}$, приходим к линейной неоднородной краевой задаче относительно $u_1(t, \tau, x)$, в которой считаем медленное время τ "параметром":

$$\frac{\partial^2 u_1}{\partial t^2} - a^2 \frac{\partial^2 u_1}{\partial x^2} + \omega^2 u_1 = F(t, \tau, x), \quad (8)$$

$$u_1(t, \tau, 0) = u_1(t, \tau, \pi) = 0. \quad (9)$$

Здесь

$$\begin{aligned} F(t, \tau, x) = & -2 \sum_{n=1}^{\infty} (i\omega_n z'_n(\tau) \exp(i\omega_n t) - i\omega_n \bar{z}'_n(\tau) \exp(-i\omega_n t)) \sin nx + \\ & + 2 \sum_{n=1}^{\infty} (i\omega_n z_n(\tau) \exp(i\omega_n t) - i\omega_n \bar{z}_n(\tau) \exp(-i\omega_n t)) \sin nx - \left(\frac{\partial u_0(t, \tau, x)}{\partial t} \right)^3. \end{aligned}$$

Условия разрешимости в классе тригонометрических по t многочленов приводят к счетной системе обыкновенных дифференциальных уравнений относительно $z_n(\tau)$, $\bar{z}_n(\tau)$, которую и принято называть квазинормальной формой. В общем случае структура квазинормальной формы существенно зависит от резонансов третьего порядка, когда при некоторых числах r_1, r_2, r_3 любого знака, связанных равенством $|r_1| + |r_2| + |r_3| = 3$, имеют место равенства

$$\omega_n = r_1 \omega_k + r_2 \omega_m + r_3 \omega_p.$$

Используя ортогональность функций $\{\sin nx\}$ ($n = 1, 2, \dots$), можно показать, что возможны лишь "тождественные" резонансы: $\omega_n = \omega_n + \omega_k - \omega_k$. Этот вопрос был рассмотрен в работе [5], где для его решения использовались соображения, связанные с тем обстоятельством, что функция $f(v) = (1 + \frac{a^2}{v^2})^{\frac{1}{2}}$ строго убывает и выпукла вниз. Следует отметить, что этот факт можно проверить и иначе: путем использования алгебраических преобразований с использованием свойств элементарных неравенств. Но второй путь имеет и недостатки: приходится рассматривать слишком большое количество случаев.

Итак, квазинормальная форма имеет вид

$$z'_n = z_n \left\{ 1 - \frac{9}{8} \omega_n^2 |z_n|^2 - \frac{3}{2} \sum_{k=1, k \neq n}^{\infty} \omega_n^2 |z_k|^2 \right\}.$$

После замены $x_n = \frac{8}{9\omega_n^2}y_n$ для квадратов амплитуд $\rho_n = |y_n|^2$ окончательно имеем систему

$$\rho'_n = \rho_n \left\{ 1 - \rho_n - \frac{4}{3} \sum_{k=1, k \neq n}^{\infty} \rho_k \right\}. \quad (10)$$

Здесь штрихом обозначена, естественно, производная по t .

3. Основные результаты

Рассмотрим сначала систему дифференциальных уравнений (10). Нетрудно проверить, что эта система имеет счетное множество состояний равновесия

$$R_1 : \rho_n = 1, \rho_k = 0 \ (k \neq n), \ k \in \mathbb{N}.$$

Линеаризация на нем приводит к системе

$$u'_n = -u_n, u'_k = -\frac{1}{3}u_k \ (k \neq n).$$

Последнее означает, что состояния равновесия типа R_1 асимптотически устойчивы.

Привлекая результаты работы [5], нетрудно аналогичными методами доказать следующее утверждение.

Теорема 1. По любому фиксированному $n_0 \in \mathbb{N}$ можно выбрать такое ε_0 , что при $0 < \varepsilon < \varepsilon_0$ краевая задача (2), (3) имеет по периодических решений вида

$$u_n = \frac{2}{3\omega_n} \sqrt{2\varepsilon} (\exp(i\omega_n t) + \exp(-i\omega_n t)) \sin nx + \dots,$$

где $n \leq n_0$, а точками обозначены слагаемые, имеющие более высокий порядок малости по ε . Все эти решения орбитально асимптотически устойчивы.

Рассмотрим теперь вопрос о состояниях равновесия системы дифференциальных уравнений (10), иного типа, отличных от R_1 . В качестве таковых можно предложить состояния равновесия

$$R_2 : \rho_{n_j} = \frac{3}{4m-1}, \rho_k = 0, j = 1, \dots, m, k \in \mathbb{N}.$$

Нетрудно проверить, что все эти состояния неустойчивы. Действительно, положив $\rho_{n_j} = \frac{3}{4m-1} + \omega_{n_j}$, $\rho_k = \omega_k$, после линеаризации и отбрасывания членов экспоненциально стремящихся к нулю получаем:

$$\omega'_k = -\frac{1}{4m-1}\omega_k \quad k \neq n_j,$$

$$\omega' = -\frac{3}{4m-1}B\omega,$$

где $\omega = (\omega_{n_1}, \dots, \omega_{n_m})$, а

$$B = \begin{pmatrix} 1 & \frac{4}{3} & \dots & \dots & \dots & \frac{4}{3} \\ \frac{4}{3} & 1 & \dots & \dots & \dots & \frac{4}{3} \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ \frac{4}{3} & \frac{4}{3} & \dots & \dots & \frac{4}{3} & 1 \end{pmatrix},$$

т.е. квадратная матрица порядка m , у которой есть по крайней мере одно собственное значение λ_B , лежащее в правой части комплексной полуплоскости.

Понятно, что состояниям равновесия типа R_2 соответствуют инвариантные торы, но эти торы неустойчивы.

Достаточно подробно была рассмотрена краевая задача (2), (3). Анализ краевых задач (2), (4) и (2), (5) приводит в целом к аналогичным результатам. Понятно, что учет пространственного взаимодействия, например, межрегиональной торговли следует учитывать более содержательным образом, считая, что пространственная переменная двумерна $x = (x_1, x_2)$. Тем не менее даже анализ чисто модельного случая приводит нас к выводам, что даже простейший учет межрегиональной торговли в значительной степени усложняет динамику. В частности, полученные результаты показывают на возможность экономических циклов разных периодов. Ясно, что реализуется один из них, но выбор экономической системой цикла, который реализуется, существенно зависит от начальных условий при запуске экономического механизма, начальных инвестиций и других экономических параметров, таких как размещение начальных капиталов по регионам.

Литература

- [1] Ну Т. Нелинейная экономическая динамика. Ижевск: Издательский дом "Удмуртский университет". 2000.
- [2] Allen R.G.D. Mathematical Economics. New York : MacMillan. 1956.
- [3] Кейнс Дж. Избранные произведения. М: Экономика, 1993.
- [4] Bechmann M.J. Puu T. Spectial Economics: Density, Potential and Flow. Amsterdam : North - Holland Publishing Company, 1985.
- [5] Колесов Ю.С. Свойства устойчивости циклов и торов пространственного нерезонансного уравнения волнового типа // Математические заметки. 1997. Т. 62. В 5. С. 744 – 750.

Об одной модификации организации пачечной волновой активности в сетях W-нейронов

Овечкин А.С.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

получена 23 марта 2003

В работах [1],[2] и [3] исследован вопрос об организации пачечной волновой активности в сети W-нейронов. В данной работе предлагается модифицировать W-нейроны, образующие сеть, путем введения альтернативных синаптических связей [4] (AW-нейроны), и рассмотреть характеристики полученной в результате этого сети на примере модели организации пачечной волновой активности.

1. Описание AW-нейрона

Для начала напомним введенное в работе [4] описание AW-нейрона.

Пусть $x = (x_1, \dots, x_n)$ - вектор входов, причем все $x_i \in \{-1; 1\}$, $w = (w^+, w^-)$ - вектор синаптических весов AW-нейрона, где $w^+ = (w_1^+, \dots, w_n^+)$, $w^- = (w_1^-, \dots, w_n^-)$, и при этом устройство нейронного сумматора таково, что положительные компоненты входных сигналов взвешиваются соответствующими весовыми компонентами из w^+ , а отрицательные - из w^- . Тогда мембранный потенциал AW-нейрона выглядит следующим образом:

$$u(t) = q\delta_0(s(t-1))u(t-1) + \left[\sum_{i=1}^N x_i(t)(w_i^+ \theta(x_i(t)) + w_i^- \theta(-x_i(t))) \right] \delta_0(s(t)) \quad (1)$$

$$\theta(z) = \begin{cases} 1, & \text{если } z \geq 0, \\ 0, & \text{если } z < 0. \end{cases} \quad (2)$$

Параметр q , как и в случае W-нейрона, определяется как $0 \leq q < 1$. Также, аналогично W-нейрону, вводится число $u_0 > 0$ - пороговое значение мембранныго потенциала. Если в момент времени t значение $u(t) \geq u_0$, то в следующий момент времени $s(t+1) = 1$, т.е. AW-нейрон переходит в возбужденное состояние и формирует выходной сигнал $x(t+1) = 1$. В противном случае AW-нейрон формирует выходной сигнал $x(t+1) = -1$ (в отличие от W-нейрона, формирующего на выходе 0).

Также напомним приведенное в [4] описание процедуры обучения AW-нейрона. Пусть имеется набор из M эталонов - векторов из B^n т.е. набор $(x^{(1)}, \dots, x^{(M)})$, где $x^{(k)} = (x_1^{(k)}, \dots, x_n^{(k)})$, а его компоненты $x_i^{(k)} \in \{-1; +1\}$. Пусть $y^{(1)}, \dots, y^{(M)}$ - требуемые выходы нейрона для каждого входного эталона соответственно, $y^{(i)} \in \{-1; 1\}$. Тогда алгоритм обучения AW-нейрона заключается в следующем:

1) Организуем из данных эталонов бесконечную последовательность $x^{(1)}, \dots, x^{(M)}, x^{(1)}, \dots, x^{(M)}, \dots$, полученную из исходного множества эталонов его циклическим повторением. Период обучения, в течение которого на вход сети были поданы $x^{(1)}, \dots, x^{(M)}$ будем называть макроитерацией обучения. Начальные (в момент времени $t = 0$) значения компонент векторов $w_i^+(0)$ и $w_i^-(0)$ выбираем произвольно.

2) Макроитерация, т.е. выполнение для каждого синапса на исходном множестве эталонов следующей процедуры модификации весовых коэффициентов:

$$w_i^+(t+1) = w_i^+(t) + delta(t), \text{ если } x_i^{(j)} = +1 \quad (3)$$

$$w_i^-(t+1) = w_i^-(t) + delta(t), \text{ если } x_i^{(j)} = -1 \quad (4)$$

где $delta$ вычисляется следующим образом:

$$delta(t) = x_i^{(j)} \left[y^{(j)} - sign \left(\left[\sum_{i=1}^N x_i(w_i^+ \theta(x_i) + w_i^- \theta(-x_i)) \right] - u_0 \right) \right] \quad (5)$$

3) Если в результате макроитерации не выполнено условие останова, то переход на 2). Условием останова может, например, служить отсутствие изменений весовых векторов наряду с ограничением общего числа макроитераций.

2. Описание модифицированной модели организации пачечной волновой активности

Для организации режима пачечной волновой активности в кольцевой нейроструктуре, описанной в [1], [2] и [3], потребуется, согласно указанным работам, опустить в формуле (1) сомножитель $\delta_0(s(t-1))$. В этом случае мембранный потенциал и нейрона не будет "сбрасываться" после перехода в состояние возбуждения. Однако в случае с биполярными входами этого недостаточно. Потребует также, чтобы мембранный потенциал нейрона в любой момент времени оставался неотрицательным, что согласуется и с биологическими соображениями. Таким образом формула (1) переписывается как:

$$u(t) = \begin{cases} u_1(t), & \text{если } u_1(t) \geq 0, \\ 0, & \text{если } u_1(t) < 0. \end{cases} \quad (6)$$

$$u_1(t) = q u(t-1) + \left[\sum_{i=1}^N x_i(t) (w_i^+ \theta(x_i(t)) + w_i^- \theta(-x_i(t))) \right] \delta_0(s(t)) \quad (7)$$

Пусть в некоторый момент времени мембранный потенциал нейрона становится больше порогового значения, и элемент генерирует единичный выходной сигнал. На следующий такт AW-нейрон переходит в рефрактерное состояние, т.е. нейрон становится невосприимчивым к внешнему воздействию. Однако значение мембранных потенциала остается больше порогового и нейрон генерирует еще один импульс. Таким образом, AW-нейрон будет учащенно генерировать единичные выходные сигналы - пачку импульсов. Выбирая силу синаптических связей для одиночного нейрона, можно управлять длиной пачки (количеством импульсов в пачке). Обозначим длину пачки через γ .

Для генерации импульса AW-нейроном в течение ровно двух последовательных тактов t и $t+1$ ($\gamma = 2$), как и в случае с W-нейроном, необходимо и достаточно, чтобы значение мембранных потенциала удовлетворяло неравенству:

$$q^{-1} u_0 < u(t) < q^{-2} u_0 \quad (8)$$

Пусть AW-нейрон в момент времени $t_0 + 1$ перешел в рефрактерное состояние. Переход мог быть вызван либо генерацией импульса на предыдущем такте, либо поступлением ненулевого сигнала на безусловно тормозящий синапс в момент t_0 . Обозначим $\sum_{i=1}^N x_i(t_0 + r_0 + 1) (w_i^+ \theta(x_i(t_0 + r_0 + 1)) + w_i^- \theta(-x_i(t_0 + r_0 + 1))) = M > 0$. После выхода нейрона из рефрактерного состояния, когда на его синапсы поступают сигналы, значение мембранных потенциала

$$u(t_0 + r_0 + 1) = u(t_0) q^{r_0 + 1} + M \quad (9)$$

Учитывая, что нейрон мог не генерировать импульс в момент t_0 , имеем ограничение снизу $u(t_0 + r_0 + 1) \geq M$. Если же в моменты времени t_0 и $t_0 + 1$ нейрон генерировал импульсы, то в силу ограничения (8) ($u(t_0) < u_0 q^{-2}$) получим ограничение сверху: $u(t_0 + r_0 + 1) < u_0 q^{r_0 - 1} + M$. Для выполнения условия двойной генерации импульса в моменты времени $t_0 + r_0 + 2$ и $t_0 + r_0 + 3$ необходимо потребовать $q^{-1} u_0 < u(t_0 + r_0 + 1) < q^{-2} u_0$, а из ограничений сверху и снизу на мембранный потенциал в момент $t_0 + r_0 + 1$ следует, что необходимо потребовать

$$q^{-1} u_0 < M \quad (10)$$

$$u_0 q^{r_0 - 1} + M < q^{-2} u_0 \quad (11)$$

Отсюда имеем неравенство

$$\frac{M q^2}{1 - q^{r_0 + 1}} < u_0 < M q \quad (12)$$

Рассмотрим, аналогично работе [3], сеть, состоящую из p модулей, каждый из которых содержит N AW-нейронов. Параметр $\gamma = 2$. Внутри модуля все нейроны связаны безусловно тормозящими синапсами, при поступлении сигнала на которые нейрон становится невосприимчивым к внешнему воздействию.

Обозначим через $W_{k,j} = (W_{k,j}^+, W_{k,j}^-)$ матрицу, в i -й строке которой ($i = 1, \dots, N$) расположена пара векторов (w^+, w^-), представляющих собой синаптические веса воздействия AW-нейронов j -го модуля на i -й нейрон k -го модуля.

Опишем интересующий нас колебательный режим. Пусть в нулевой момент времени часть AW-нейронов одного из модулей переходит в возбужденное состояние и генерирует единичные выходные сигналы. В следующий момент времени наблюдаются неунловые выходные сигналы у нейронов следующего модуля и вторичная импульсация у нейронов первого модуля. Эти сигналы индуцируют генерацию импульсов элементами третьего модуля. Волна возбуждения последовательно обходит все модули. Следующий тakt прохождения волны возбуждения по сети открывается генерацией единичных выходных сигналов частью нейронов исходного модуля, (не совпадающей с множеством нейронов, возбужденных на первом такте). После ряда тактов прохождения волны в возбужденное состояние перейдут те нейроны первого и второго модулей, с которых начинался процесс распространения волны. В дальнейшем процесс периодически повторяется. В каждый момент времени в возбужденном состоянии оказываются элементы двух модулей.

Пусть продолжительность рефракторного состояния AW-нейронов $r_0 = p - 2$. Рассмотрим наборы биполярных (каждая координата принадлежит множеству $\{+1; -1\}$) векторов $X_1^i, X_2^i, \dots, X_p^i$ ($i = 1, \dots, m; m \leq 2N + 1$). Будем считать, что векторы $(X_{k-1}^i, X_k^i) \in B^{2N}$ линейно независимы. Синаптические веса, обозначенные матрицами $W_{k,k-2}$ и $W_{k,k-1}$ ($k = 3, \dots, p$), выберем по правилам обучения, описанным в [4] так, чтобы на входной вектор (X_{k-2}^i, X_{k-1}^i) нейроны k -го модуля отвечали выходным вектором X_k^i .

Матрицы $W_{1,p-1}$ и $W_{1,p}$ выберем так, чтобы для векторов (X_{p-1}^i, X_p^i) выходными векторами первого модуля служили векторы X_i^{i+1} ($i = 1, \dots, m - 1$), а для вектора (X_{p-1}^m, X_p^m) - вектор X_1^1 .

Матрицы синаптических весов выбираются с учетом неравенства (12).

Аналогично указанному в работе [3], можно сформулировать утверждение:

Утверждение. При указанном выборе синаптических весов (если продолжительность рефрактерного состояния $r_0 = p - 2$) существует периодический режим функционирования сети, в котором в последовательные моменты времени генерируются выходные сигналы:

$$(X_1^1, X_2^1), \dots, (X_p^1, X_1^2), (X_1^2, X_2^2), \dots, (X_1^m, X_2^m), \dots, (X_p^m, X_1^1), (X_1^1, X_2^1), \dots \quad (13)$$

Пара (X_{k-1}^i, X_k^i) формируется $k - 1$ -ым и k -ым модулями соответственно.

Инициализация описанного режима также полностью аналогична инициализации режима, описанной в [3].

3. Результаты сравнений характеристик исходной и модифицированной сетей

Для экспериментального изучения и сравнения описанных выше моделей нейронов были рассмотрены модели кольцевых структур из 3 модулей размером 30 x 30 элементов на базе W-нейронов и AW-нейронов. Каждая сеть была обучена на наборе из 15 черно-белых образов (по 5 эталонов на модуль). В качестве образов использовались буквы латинского алфавита, выбранные произвольно. Инициализация режимов производилась в соответствии с процедурой, описанной в [3], при этом "зародыш" (а точнее говоря, пара "зародышей" (см. процесс инициализации)) выбирался из обучающей последовательности и искажался случайнym шумом заданной вероятности (см. Табл.1). Искажения заключались в зачернении с вероятностью p пикселов "зародыша". После прохождения первой волны на первом модуле измерялось расстояние Хэмминга D_h между состоянием модуля и планируемым эталонным значением:

$$D_h = \sum_{i=1}^{30 \times 30} sign|x_i - u_i| \quad (14)$$

Для каждого из 5 "зародышей" (изучалось только поведение первого модуля, т.к. поведение остальных аналогично с точностью до выбора эталонов) для различных уровней шума (см. Табл.1) проводилась серия из 20 испытаний, в результате этих испытаний вычислялось среднее значение расстояний Хэмминга и их среднеквадратичное отклонение.

В таблице приведены полученные в результате экспериментов средние по всем 5 "зародышам", значения средних расстояний Хэмминга и среднеквадратичных отклонений.

Модель	Параметр	5%	10%	20%	30%
W-модель	Расстояние Хэмминга	1.18	4.02	10.36	21.66
	Ср.-кв. отклонение	2.04	4.07	6.85	12.42
AW-модель	Расстояние Хэмминга	0.24	1.06	2.5	7.64
	Ср.-кв. отклонение	0.62	1.75	2.71	6.07

Табл.1

Таким образом, данные проведенных экспериментов наглядно показывают, что использование AW-нейронов вместо W-нейронов приводит к значительному улучшению параметров сети при организации пачечной волновой активности в кольцевой нейроструктуре.

Литература

- [1] Майоров В.В., Шабаршина Г.В. Сети W-нейронов в задаче ассоциативной памяти // Журнал Вычислительной Математики и Математической Физики, 2001. Том 41, N8. С.1289
- [2] Майоров В.В., Шабаршина Г.В. Сообщение о сетях W-нейронов // Моделирование и анализ информационных систем. Ярославль, 1997. Вып. 4. С.37-50
- [3] Майоров В.В., Шабаршина Г.В. О проблеме хранения и воспроизведения последовательностей бинарных образов в сетях W-нейронов // Сборник обзорных статей к 25-летию Математического ф-та ЯрГУ. Ярославль,2001. С.195-211
- [4] Овечкин А.С. О модификации W-нейрона и алгоритме его обучения // Современные проблемы математики и информатики. Ярославль, 2002. Вып. 5. С.95-99
- [5] Короткин А.А., Панкратов В.А. Классифицирующие свойства нейронов с альтернативными синапсами // Моделирование и анализ информационных систем. Ярославль, 1997. Вып. 4. С.118-123
- [6] Уоссермен Ф. Нейрокомпьютерная техника: теория и практика. М., 1992.

УДК 681.3

Выбор ранговых и доменных блоков в алгоритме фрактального сжатия изображений

Ходаковский В.А.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

получена 19 мая 2003

В статье рассматриваются варианты выбора ранговых и доменных блоков в алгоритме фрактального сжатия. Описаны варианты с выбором фиксированных и прогрессивных квадратов. Предлагается использование в качестве ранговых и доменных блоков 2-тайлингов.

1. Введение

Напомним, что алгоритм фрактального сжатия изображений основан на самоподобии в изображениях. Задача алгоритма следующая: Пусть изображение представлено в виде функции u , заданной на множестве X . Необходимо найти систему сжимающих преобразований, таких, чтобы их аттрактор отличался от функции u не более чем на заданный ε .

Последовательность шагов в алгоритме следующая:

Шаг 1. Выберем некоторое разбиение X на N областей D_1, \dots, D_N , которые будем называть доменными блоками. Подчеркнем, что для этих блоков должны выполняться следующие условия:

- $\bigcup_{i=1}^N D_i = X$
- $D_i \cap D_j = \emptyset$, при $i \neq j$

Шаг 2. Выберем некоторое покрытие X M областями R_1, \dots, R_N , которые назовем ранговыми блоками. Заметим, что $\bigcup_{i=1}^N R_i = X$, но R_i могут пересекаться.

Шаг 3. Для каждого D_i будем искать сжимающее преобразование f_i , переводящее в D_i один из R_j наилучшим образом.

Шаг 4. Если для всех f_i отклонение меньше, чем ε , то набор преобразований считаем найденным. В противном случае необходимо рассмотреть другие разбиения X на R и D .

Более подробно с алгоритмом и его теоретическим обоснованием можно ознакомиться в [1, 2, 3, 6].

2. Выбор разбиения X

В качестве доменных и ранговых блоков можно рассматривать любые по форме фигуры. Однако обычно, для программной реализации, на блоки накладываются различные условия. Делается это для упрощения нахождения соответствующих преобразований. Рассмотрим различные разбиения X , которые можно использовать для фрактального кодирования.

2.1. Фиксированные квадраты

Рассмотрим один из простых способов разбиения, описанный в [4], который тем не менее обеспечивает неплохие результаты.

Предполагаем, что $X = [0, 1]^2$

1 Все области D и R являются квадратами со сторонами, параллельными сторонам изображения. Это ограничение достаточно жесткое. Фактически мы собираемся приближать все многообразие геометрических фигур лишь квадратами.

2 При переводе ранговой области в доменную уменьшение размеров производится ровно в два раза. Это существенно упрощает как компрессор, так и декомпрессор, т.к. задача масштабирования небольших областей является нетривиальной.

3 Все доменные блоки - квадраты и имеют фиксированный размер. Изображение равномерной сеткой разбивается на набор доменных блоков.

4 При переводе ранговой области в доменную поворот куба возможен только на $0^\circ, 90^\circ, 180^\circ$ или 270° . Также допускается зеркальное отражение. Общее число возможных преобразований (считая тождественное) - 8.

Эти ограничения позволяют:

1 Построить алгоритм, для которого требуется сравнительно малое число операций даже на достаточно больших изображениях.

2 Очень компактно представить данные для записи в файл. Нам требуется на каждый фрактальный компонент:

- два числа для того, чтобы задать смещение рангового блока.
- одно число, для того, чтобы задать преобразование симметрии при переводе рангового блока в доменный.
- два числа, для того, чтобы задать сжатие и сдвиг по яркости при переводе.

Отрицательные стороны предложенных ограничений:

1 Поскольку все области являются квадратами, невозможно воспользоваться подобием объектов, по форме далеких от квадратов (которые встречаются в реальных изображениях достаточно часто.)

2 Аналогично мы не сможем воспользоваться подобием объектов в изображении, коэффициент подобия между которыми сильно отличается от 2.

3 Алгоритм не сможет воспользоваться подобием объектов в изображении, угол между которыми не кратен 90° .

2.2. Прогрессивные квадраты или "Деревья"

Этот метод является развитием предыдущего способа разбиения

Также предполагаем, что $X = [0, 1]^2$

1 Все области D и R являются квадратами со сторонами, параллельными сторонам изображения.

2 При переводе ранговой области в доменную уменьшение размеров производится ровно в два, четыре, восемь раз. Это существенно упрощает как компрессор, так и декомпрессор, т.к. задача масштабирования небольших областей является нетривиальной.

3 Все доменные блоки - квадраты и имеют один из фиксированных размеров: базовый, в два раза меньший, в четыре и т.д.

4 При переводе ранговой области в доменную поворот куба возможен только на $0^\circ, 90^\circ, 180^\circ$ или 270° .
Также допускается зеркальное отражение. Общее число возможных преобразований (считая пустое) - 8.

Строится разбиение следующим образом:

1 Изображение равномерной сеткой разбивается на набор доменных блоков базового размера.

2 Для каждого доменного блока производится поиск рангового (шаг 3 алгоритма).

3 Если невозможно подобрать ранговый блок с заданной точностью, то разбиваем доменный блок на четыре квадрата меньшего размера (очевидно, что при этом не нарушается условие 1-ого шага алгоритма). Теперь каждый меньший квадрат будем считать новым доменным блоком.

4 Для новых доменных блоков производим поиск ранговых, при необходимости, снова разбивая их.

Таким образом, разбиение изображения вначале не известно и уточняется в ходе работы алгоритма. Глубина разбиения обычно ограничивается разумными пределами, например ограничением минимального размера блока.

Для одинаковой точности кодирования этот метод обычно требует меньшее число доменных блоков при разбиении и, следовательно, закодированный файл имеет меньший размер. Отрицательные стороны предложенных ограничений схожи с методом фиксированных квадратов.

2.3. Тайлинги

2.3.1. Определение

Компактное множество $T \in R^d$ называется *тайлингом* в R^d , если существует счетный набор непересекающихся множеств $\{T_1, T_2, \dots\}$, такой что каждое T_j совпадает с T и их объединение совпадает с R^d .

Тайлинг T называется *рептайлингом*, если T может быть разделено на n компактных подмножеств $\Omega_1, \dots, \Omega_n$ с непересекающимися внутренностями. Все Ω_j совпадают друг с другом с точностью до поворота и сдвига (но не отражения), и все Ω_j подобны множеству T .

Приравняв объемы множества T и суммы объемов множеств Ω_j получаем, что каждое Ω_j уменьшено относительно T в $\sqrt[n]{n}$ раз. Можно определить T следующим образом:

$$T = \bigcup_{j=1}^n \frac{1}{\sqrt[n]{n}} v_j(T), \quad (1)$$

где v_j изометрия в R^d .

Иными словами T является неподвижной точкой системы итерируемых функций $\{n^{-1/d}v_j : 1 \leq j \leq n\}$ [5].

Как уже стало ясно, рептайлинги в R^2 естественным образом подходят для разбиения изображения на доменные и ранговые блоки.

2.3.2. Использование 2-п рептайлингов

Опишем применение 2-п рептайлингов для разбиения изображения на ранговые и доменные блоки при фрактальном кодировании. Наиболее удобными 2-п рептайлингами являются Twindragon, Rectangle и Tame Twindragon (рис 1).

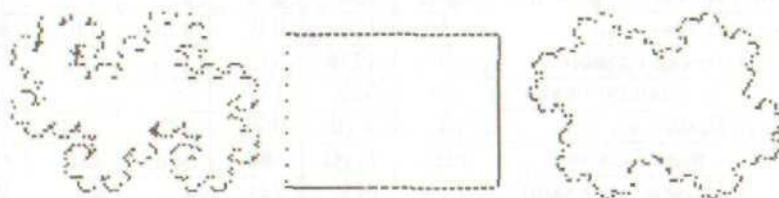


Рис. 1. Twindragon, Rectangle и Tame Twindragon

Аналогично разбиению изображения на квадратные блоки, разбиение на 2-рептайлинги может быть фиксированным и прогрессивным. Опишем прогрессивный метод.

В отличие от предыдущих методов предполагаем, что X представляет собой фигуру на плоскости - Тайлинг (один из представленных на рис 1.) Для применения метода к реальным изображениям, которые обычно представляют собой прямоугольник, достаточно поместить изображение внутрь тайлинга.

- 1 Итак, все области D и R являются тайлингами, подобными исходному изображению.
- 2 При переводе ранговой области в доменную уменьшение размеров производится ровно в $\sqrt{2}$ раз, два раза и т.д.
- 3 Все доменные блоки из фиксированных размеров: базовый, в $\sqrt{2}$ раз меньший, в два раза, в четыре и т.д.
- 4 Перевод ранговой области в доменную осуществляется одной из функций построения тайлинга. Выбирается из v_i .

Заметим, что для построения любого рангового или доменного блока необходимо последовательно применять функции v_i к исходному тайлингу. Т.к. рассматриваются 2-рептайлинги, то существуют функции v_1 и v_2 . Следовательно, любой блок, а точнее говоря, последовательность функций, необходимых для его построения, удобно кодируется двоичным числом.

Положительными сторонами данного метода являются:

1 Более "плавное" деление блоков. Не всегда для достижения нужной точности необходимо дробить блоки на 4 части, очень часто достаточно на 2-3 части.

2 Фигуры Twindragon и Tame Twindragon значительно отличаются от квадратов и, следовательно, в некоторых случаях дают преимущество.

3 Достаточно простое построение ранговых и доменных блоков, а также очень компактное хранение.

Не обойтись и без отрицательных сторон:

1 При преобразовании необходимо применять только функции построения самого тайлинга, т.е. всего два преобразования.

2 Сложности при преобразовании маленьких областей.

2.4. Численный эксперимент

Для кодирования было использовано изображение LENA 256×256 пикселей и 256 градаций серого. Рассматривались разбиения на "прогрессивные" квадраты и 2-рептайлины rectangle и tame twindragon. Для выбранной допустимой погрешности рассматривалось общее число потребовавшихся блоков, а также число блоков, которые не удалось закодировать. Глубина разбиения была ограничена для квадратов четырьмя, а для 2-рептайлингов восьмью делениями. В таблице представлены полученные результаты.

Таблица 1. Результаты эксперимента

	Погрешность	5	10	20	30	40	50
Square	Число блоков	3106	2242	1168	496	187	103
	Не закодировано	1751	833	185	29	2	1
	Время	1:16	1:15	0:24	0:06	0:02	40с
Rectangle	Число блоков	2701	1716	873	365	166	84
	Не закодировано	1293	535	137	26	2	1
	Время	3:19	1:10	0:27	0:07	0:02	44с
Tame twindragon	Число блоков	2880	1780	883	393	146	80
	Не закодировано	1901	878	224	28	1	0
	Время	3:32	1:33	0:34	0:09	0:02	45с

Как видно из таблицы, наиболее выгодно применение прямоугольников. При применении Tame twindragon также получается меньше блоков, чем при применении квадратов. Однако сказываются небольшие размеры изображения и, следовательно, большие погрешности при построении тайлингов.

Литература

- [1] Bruno Forte and Edward R. Vrscay. Theory of Generalized Fractal Transforms. Department of Applied Mathematics Faculty of Mathematics University of Waterloo, 1996.
- [2] Bruno Forte and Edward R. Vrscay. Inverse Problem Methods for Generalized Fractal Transforms. Department of Applied Mathematics Faculty of Mathematics University of Waterloo, 1996.
- [3] Yoval Fisher. Fractal Image Compression. SIGGRAPH, 1992.
- [4] Ватолин Д.С. Алгоритмы сжатия изображений. М.: Издательство МГУ, 1999.
- [5] On 2-Reptiles in the Plane. Sze-Man Ngai, Victor F. y Sirvent, J. J. P. Veerman, Yang Wang School of Mathematics, Georgia Institute of Technology Atlanta, GA 30332, USA December 16, 1999
- [6] Ходаковский В.А. Модификация IFS алгоритма кодирования изображений // Модел. и анализ информ. систем. Т.10, №1. 2003. С.31-34.

Планирование волн в кольцевых структурах ассоциаций импульсных нейронов-детекторов

Майоров В.В., Мышкин И.Ю., Мячин М.Л., Куксов А.Г.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

We propose the neural network consist of the detector neurons being described by the differential equation with delay. The network architecture is the oriented circle formed of the full-connected neural associations. We propose the rule for weight assigning so that the network has the preassigned cyclic wave mode.

В настоящее время нейрофизиологами найдены экспериментальные базовые факты и разработаны представления о возможности кодирования информации в мозге в виде незатухающих циклически повторяющихся волн нейронной активности. Такие волны создаются коллективным поведением относительно синхронно функционирующих нейронных ансамблей, причем одни нейроны могут выключаться на время из конкретного ансамбля, но их замещают нейроны другого ансамбля, генерируя импульсы в подходящие моменты. Таким образом, нейронные коалиции являются динамическими образованиями.

При волновом подходе существенно соотношение фаз колебаний, поскольку как отдельные нейроны, так и нейронные ассоциации обладают свойством рефрактерности — снижением восприимчивости к воздействию после генерации импульсов. Поэтому волновой подход к кодированию часто отождествляют с фазово-частотным. Волновой подход созвучен с идеями нелинейной динамики, согласно которым информация может кодироваться не стационарным распределением двоичной величины, а колебательным процессом динамической системы.

В работах [1, 2] предложена модель нейронной сети, генерирующей последовательности заранее заданной структуры. Ее элементы (импульсные автогенераторы) описываются уравнением с запаздыванием. Сама сеть устроена как кольцо нейронных ассоциаций. Именно по нему циклически распространяются волны импульсов. На каждом такте прохождения волны в любой ассоциации импульсы генерируют одни и те же элементы.

W-нейроны [3] не обладают авторитмичностью и являются пороговыми элементами. Тем не менее по кольцевой структуре нейронных ассоциаций (после специального выбора весов) могут распространяться циклические волны нейронной активности (волны импульсов). При этом на разных тактах прохождения волны в каждой ассоциации могут быть активны различные элементы. Следует отметить, что W-нейроны функционируют в дискретном времени.

Ниже рассматриваются сети, элементы которых описываются уравнениями с запаздыванием, и при этом являются детекторами. Достаточно сильное воздействие на нейрон (если он не рефрактерен) приводит к генерации импульса. Архитектура сети — ориентированное кольцо полностью связанных нейронных ассоциаций. Будет показано, что при определенном выборе синаптических весов по такой сети распространяется заранее запланированная циклическая волна импульсов. На разных тактах в ассоциациях импульсы генерируются, вообще говоря, различные группы нейронов.

1. Модель нейрона

Нейрон — структурная единица нервной системы — состоит из тела (центральной части) и древовидных сильно ветвящихся отростков. Разветвление единственного длинного отростка — аксона — образует на других нейронах контактные окончания — синапсы. На мембране, определяющей индивидуальные границы нейрона, наблюдается разность потенциалов — мембранный потенциал. Большую часть времени мембранный потенциал отрицателен и велик по абсолютной величине (состояние поляризации). Мембранный потенциал обусловлен относительным дефицитом ионов натрия внутри клетки. Нейрон способен генерировать электрические импульсы. Спонтанно или в результате воздействия в мембране открываются каналы и ионы натрия устремляются внутрь клетки. Абсолютное значение мембранныго потенциала снижается, а затем он меняет знак. Здесь натриевые каналы начинают закрываться. Открываются калиевые каналы, по которым ионы калия движутся из клетки, где наблюдается их избыток. Мембрана возвращается в состояние поляризации. Описанный натриево-калиевый цикл генерации импульса был открыт А. Ходжкиным и А. Хаксли. Импульсы называются потенциалами действия или спайками.

Рожденный в теле нейрона спайк, распространяясь по аксону и его разветвлениям, достигает синапсов. В пресинаптической части (непосредственном окончании аксона) происходит выброс заключенного в пузырьки специфического вещества — медиатора. Пузырьки пересекают синаптическую щель. На мембране нейрона-приемника открывается медиатор-зависимые ионные каналы. Соответствующие ионные токи могут или дополнительно поляризовать мембрану нейрона-приемника (тормозное действие), или уменьшить абсолютное значение мембранныго потенциала (возбуждающее действие). В последнем случае, если деполяризация была существенной, нейрон-приемник может генерировать импульс.

Следует отметить, что нейроны бывают двух типов — автогенераторы и детекторы. Нейроны-автогенераторы способны периодически генерировать импульсы, а медиаторное воздействие приближает или отдаляет момент генерации следующего импульса. Мембранный потенциал нейрона-детектора в отсутствие воздействия находится в состоянии равновесия (покоя, который, вообще говоря, относителен). Нейроны-детекторы генерируют импульсы в результате достаточно сильного (запорогового) воздействия.

Ниже рассматриваются сети нейронов-детекторов. Их модель базируется на предложенных в [2, 4] уравнениях. Примем уровень наибольшей поляризации за начало отсчета. Через $u(t) \geq 0$ обозначим отклонение мембранныго потенциала от этого уровня. В соответствии с идеологией работ [2, 4] примем для $u(t)$ уравнение

$$\dot{u} = \lambda[-1 - f_{Na}(u) + f_K(u(t-1))]u + \lambda\alpha f_L(u(t-1))\exp(-\lambda\sigma). \quad (1)$$

Здесь $\lambda \gg 1$ — скоростной параметр, достаточно гладкие функции $f_{Na}(u)$, $f_K(u)$, $\Theta(\cdot)$ положительны и монотонно убывают к нулю при $u \rightarrow \infty$ быстрее, чем u^{-1} . Функции $f_{Na}(u)$ и $f_K(u)$ описывают проводимости натриевых и калиевых каналов, а $f_L(u)$ — проводимость всех прочих ионных каналов (токи утечки). Калиевая проводимость запаздывает по времени (А. Ходжкин, А. Хаксли). Величина задержки (длительность восходящего участка спайка) принята за единицу времени. По нашим представлениям токи утечки также запаздывают.

Введем следующие обозначения

$$\begin{aligned}\alpha &= 1 + f_{Na}(0) - f_K(0), \\ \alpha_1 &= f_K(0) - 1, \\ \alpha_2 &= f_{Na}(0) + 1\end{aligned}$$

и будем считать, что $\sigma > 0$, $\alpha > 0$, $\alpha_1 > 0$, $f_K(0) - f_{Na}(1) - 1 > 0$, $f_L(0) = 1$, $\alpha_2 > \sigma$. При этих ограничениях уравнение (1) имеет устойчивое состояние равновесия $u_* \approx \exp(-\lambda\sigma)$ и решения импульсного типа. Начало и окончание импульса (спайка) условно связем с моментами времени, когда $u(t)$ пересекает единичное значение соответственно с положительной и отрицательной производной. Для $\beta > 0$ через S_β обозначим класс непрерывных на отрезке $s \in [-1, 0]$ функций $u(s)$, для которых $u(0) = 1$ и $0 \leq u(s) \leq \max\{\exp(\lambda\beta s), \lambda \exp(-\beta\sigma)\}$.

Асимптотический анализ при $\lambda \gg 1$, если пренебречь бесконечно малыми по времени переходными участками, показывает следующее:

$$\begin{aligned}u(t) &\approx \exp(\lambda\alpha_1 t), & t \in (0, 1), \\ u(t) &\approx \exp(\lambda(t-1)), & t \in (1, 1+\alpha_1), \\ u(t) &\approx \exp(-\lambda\alpha_2(t-1-\alpha_1)), & t \in (1+\alpha_1, 2+\alpha_1), \\ u(t) &\approx \exp(-\lambda\sigma), & t > 2+\alpha_1.\end{aligned}$$

Тем самым продолжительность спайка составляет примерно $T_1 = 1 + \alpha_1$. Импульс покоящегося нейрона можно инициализировать с помощью внешнего воздействия, например следующим способом: на промежутке времени продолжительности $\delta < 1$ положить $\sigma = 0$, а затем восстановить его ненулевое значение.

Опишем модель взаимодействия. Будем учитывать явление рефрактерности: в течение спайка и на некотором промежутке времени после него нейрон не реагирует на внешнее воздействие. Период рефрактерности T_R обычно составляет две-три продолжительности спайка ($T_R = 2T_1 \div 3T_1$). Далее, будем считать, что медиатор мгновенно появляется в начале спайка нейрона-передатчика, а в момент окончания спайка мгновенно разрушается.

Пусть $u(t)$ и $v(t)$ — мембранные потенциалы нейронов приемника и передатчика. Для $u(t)$ примем [1] уравнение

$$\dot{u} = \lambda[-1 - f_{Na}(u) + f_K(u(t-1)) + \alpha q H(u)\Theta(v(t)-1)]u + \lambda\alpha f_L(u(t-1))\exp(-\lambda\sigma). \quad (2)$$

Здесь $H(u)$ — функционал, который обеспечивает наличие рефрактерного периода. Например, достаточно следить за интегралом от $u(t)$ на промежутке времени $[t - T_R, t]$. Если он больше λ , то значение функционала равно нулю. Далее, $\Theta(\cdot)$ — функция Хевисайда, т.е. величина $\Theta(v(t) - 1)$ представляет собой индикатор присутствия медиатора. Наконец, q — синаптический вес.

Пусть в (2) выполнено неравенство $q > 1$ и $u(t) = u_*$ (состояние равновесия) при $t \leq 0$, а в нулевой момент времени начался спайк нейрона-передатчика. Асимптотический анализ уравнения (2) при $\lambda \gg 1$ показывает, что импульс нейрона-приемника начнется в момент времени $\tau \approx \sigma/\alpha(q-1)$.

Если нейрон находится под воздействием нескольких нейронов с мембранными потенциалами $v_i(t)$, то множитель $q\Theta(v - 1)$ заменяется на $\sum q_i\Theta(v_i - 1)$, где q_i — синаптические веса.

2. Архитектура и система уравнений сети

Рассмотрим N нейронных ассоциаций, каждая из которых состоит из m нейронов. Циклически пронумеруем ассоциации, т.е. будем отождествлять номера i и $i + kN$, где $i = 1, \dots, N$, а k — произвольное целое число. Будем считать, что любой нейрон каждой ассоциации может действовать на нейроны следующей по номеру ассоциации, т.е. иметь на них синаптические окончания. Кроме того, пусть в каждой ассоциации каждый нейрон действует на все остальные нейроны с одинаковым фиксированным весом $q > 0$. Любые другие синаптические связи запретим. Занумеруем нейроны парой индексов (i, j) , где i — номер ассоциации, j — номер нейрона в ассоциации. Пусть u_{ij} — мембранные потенциалы нейронов. Для произвольного (i, j) -го нейрона через q_{ij}^s ($s = 1, \dots, m$) обозначим синаптический вес воздействия на него со стороны нейрона с номером $(i-1, s)$, т.е. s -го нейрона из $(i-1)$ -ой ассоциации. Для мембранных потенциалов в соответствии с предыдущим пунктом получим систему уравнений

$$\dot{u}_{ij} = \lambda \left[-1 - f_{Na}(u_{ij}) + f_K(u_{ij}(t-1)) + \right. \\ \left. + \alpha H(u_{ij}) \left(\sum_{s=1}^m q_{ij}^s \Theta(u_{i-1,s} - 1) + q \sum_{\substack{s=1 \\ s \neq j}}^m \Theta(u_{is} - 1) \right) \right] u_{ij} + \\ + \lambda \alpha f_L(u_{ij}(t-1)) \exp(-\lambda \sigma). \quad (3)$$

Эта система имеет устойчивое состояние равновесия $u_{ij} \equiv u_*$. При надлежащем выборе весов система может иметь колебательные режимы заранее спланированной структуры.

3. Выбор синаптических весов

Скажем, что в данный момент времени нейрон находится в активном состоянии, если он генерирует импульс. Пусть $U_i(t) \in R^m$ — вектор, составленный из мембранных потенциалов i -ой ассоциации. Вектор, составленный из единиц, обозначим через e . Введем бинарный вектор $X_i(t) = \Theta(U_i(t) - e)$, где функция $\Theta(\cdot)$ вычисляется покоординатно. Вектор $X_i(t)$ ($i = 1, \dots, N$) назовем вектором активности. Если его j -ая координата равна единице, то j -ый нейрон i -ой ассоциации в данный момент времени находится в активном состоянии.

Введем наборы бинарных векторов

$$X_i^s, \quad i = 1, \dots, N, \quad s = 1, \dots, k \quad (4)$$

и поставим задачу о выборе синаптических весов q_{ij}^s и q так, чтобы система (3) имела колебательный режим следующей структуры:

- с некоторым временным рассогласованием $\tau > 0$ последовательно в сторону возрастания номеров активными становятся нейроны первой, второй и т.д. ассоциаций (вслед за N -ой идет первая ассоциация);
- вектор $X_i(t)$, соответствующий активной в текущий момент времени ассоциации, принимает последовательно следующие значения:

$$X_1^1, X_2^1, \dots, X_N^1, X_1^2, \dots, X_N^2, \dots, X_1^k, \dots, X_N^k, X_1^1, \dots$$

Ясно, что каждый из векторов последовательности наблюдается в течение времени T_1 . Тем самым по кольцевой нейронной структуре распространяется волна активности. В каждой ассоциации на первых k разных обходах волны активны разные нейроны. Далее ситуация циклически повторяется. Выбор весов основан на следующем утверждении.

Лемма (о возможности обучения). Пусть векторы $X_1, \dots, X_p, Y_1, \dots, Y_l \in R^m$ ($p + l \leq m$) линейно независимы. Тогда существует ненулевой вектор $q \in R^m$, для которого

$$\begin{aligned} (q, X_i) &= \|q\|^2, & i &= 1, \dots, p, \\ (q, Y_i) &= 0, & i &= 1, \dots, l. \end{aligned}$$

Для доказательства (нахождения вектора q) достаточно для последовательности векторов

$$X_2 - X_1, \dots, X_p - X_1, Y_1, \dots, Y_l, X_1$$

применить алгоритм ортогонализации Грамма-Шмидта. Последний ортогонализованный вектор и есть вектор q .

Пусть в последовательности (4) при любом $i = 1, \dots, N$ векторы X_i^s ($s = 1, \dots, k$) линейно независимы. Тогда в силу леммы существует матрица Q_i , для которой

$$Q_i X_{i-1}^s = \gamma X_i^s, \quad (5)$$

где константа $\gamma > 1$. Элементы q_{ij}^l ($i, j = 1, \dots, m$) матрицы Q_i есть искомые синаптические веса в системе (3).

4. Обоснование существования колебательного режима

Бинарные координаты вектора X_i^s обозначим через x_{ij}^s ($j = 1, \dots, m$). Пусть до нулевого момента времени вектор активности ($i-1$ -ой ассоциации был нулевым, а начиная с нулевого момента в течение времени T_1 (длительность спайка) стал равным X_{i-1}^s . Предположим, что до нулевого момента времени нейроны i -ой ассоциации находились в состоянии равновесия: $u_{ij} \equiv u_* \approx \exp(-\lambda\sigma)$. Тогда начиная с нулевого момента времени, пока $u_{ij}(t) < \lambda^{-1}$ уравнения (3) для i -ой ассоциации приобретают вид:

$$\dot{u}_{ij} = \lambda\alpha[-1 + \gamma x_{ij}^s + o(1)]u_{ij} + \lambda\alpha \exp(-\lambda\sigma)(1 + o(1)),$$

где $o(1)$ — слагаемые, пренебрежимые по сравнению с единицей при $\lambda \gg 1$. Ясно, что при $x_{ij}^s = 0$ имеют место равенства $u_{ij} = 0$. Если же $x_{ij}^s = 1$ (нейроны подвергаются воздействию), то $u_{ij} = \exp(\lambda|\sigma| + (\gamma - 1)t + o(1))$. Пусть $0 < \tau < T_1$ и $\gamma = \sigma/\tau + 1$. Тогда $u_{ij}(\tau + o(1)) = 1$, т.е. все нейроны, подверженные ненулевому воздействию, начинают генерировать импульсы. В дальнейшем рассуждение повторяется для ($i+1$ -ой и т.д. ассоциаций). Таким образом, по кольцевой структуре распространяется запланированная заранее волна активности нейронных ассоциаций.

Реально импульсы нейронов внутри ассоциации начинаются с некоторым временным рассогласованием. Тем более, что на практике трудно обеспечить выполнение равенства (5). Ненулевые положительные веса q взаимодействия нейронов внутри ассоциаций препятствуют десинхронизации. Существование описанных колебательных режимов можно гарантировать на достаточно длительных интервалах времени.

5. Результаты моделирования

При численном моделировании использовались функции $f_K(u)$ и $f_{Na}(u)$ следующего вида:

$$\begin{aligned} f_K(u) &= 2.0 \cdot \exp(-u^2), \\ f_{Na}(u) &= 1.2 \cdot \exp(-u^2), \\ f_L(u) &= 1.0 \cdot \exp(-u^2). \end{aligned}$$

Были также выбраны следующие значения параметров: $\lambda = 5.0$, $\sigma = 0.3$, $\alpha = 0.2$. При выбранных значениях параметров были получены следующие значения для характерных времен: $T_1 \approx 2.0$ и $T_R = 4.5$. Стационарному режиму соответствует значение $u_* \approx 0.189$.

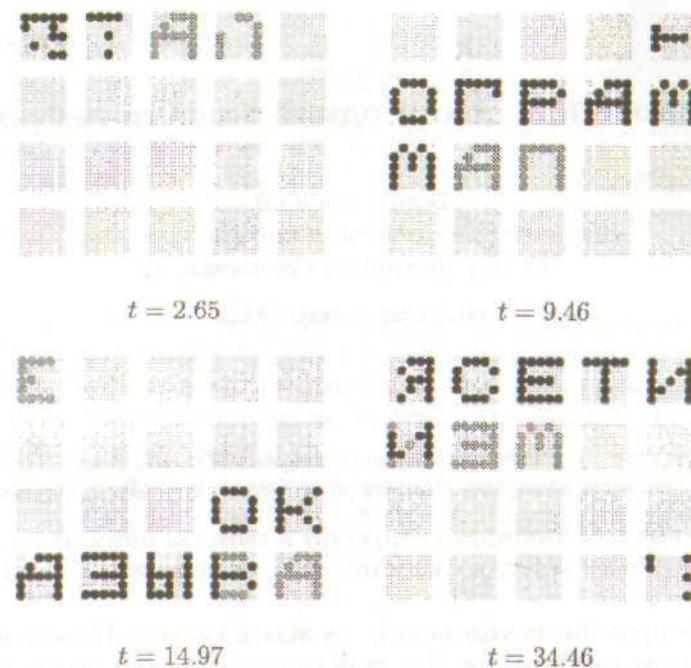


Рис. 1. Этапы эволюции состояния нейронной сети.

Моделировалась сеть, состоящая из $N = 20$ ассоциаций по $t = 25$ нейронов в каждой. Сеть была обучена на фразу, состоящую из 80 символов: «эта программа показывает возможность обучения сети из модулей нейронов на примере из 4 циклов» (пробел символом не считался). При этом на каждый нейронный модуль приходится по 4 символа. Моделирование сети производилось в дискретном времени с шагом $\Delta t = 0.01$. Активация сети производилась на первом модуле, в котором случайно задавался первый бинарный паттерн активности.

На рисунках продемонстрированы отдельные этапы эволюции состояния нейронной сети. Нейроны показаны в виде условных графических элементов, цвет которых определяется состоянием соответствующего нейрона. Светло-серый цвет соответствует состоянию ожидания; темно-серый — состоянию генерации спайка; черный — состоянию рефрактерности. Графические элементы, соответствующие нейронам из одного модуля, расположены в виде компактной группы размером 5×5 . Нейронные модули упорядочены слева направо и сверху вниз (нейроны каждого модуля действуют на нейроны модуля, расположенного справа от него; нейроны крайнего правого модуля в каждом ряду действуют на нейроны первого модуля в следующем ряду; нейроны правого нижнего модуля действуют на нейроны левого верхнего модуля).

Список литературы

1. Кащенко С.А., Майоров В.В., Мыскин И.Ю. Исследование колебаний в кольцевых нейронных системах // Доклады РАН. 1993. Т. 333, № 5. С. 594–597.
2. Кащенко С.А., Майоров В.В., Мячин М.Л. Колебания в системах уравнений с запаздыванием и разностной диффузией, моделирующих локальные нейронные сети // Доклады РАН. 1995. Т. 344, № 3. С. 1274–1279.
3. Майоров В.В., Шабаршина Г.В. Сети W -нейронов в задаче ассоциативной памяти // ЖВМиМФ. 2001. Т. 41, № 8. С. 1289–1298.
4. Майоров В.В., Мыскин И.Ю. Математическое моделирование нейронной сети на основе уравнений с запаздыванием // Математическое моделирование. 1990. Т. 2, № 11. С. 64–76.

О хаотическом поведении одной модели нейронной сети

Богомолов Ю.В.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

получена 3 июня 2003

Рассмотрена модель сети на основе формального нейрона Мак-Каллока – Питтса. Получено численное подтверждение наличия весов, для которых характерна хаотическая динамика нейронной сети.

Искусственные нейронные сети имитируют структуру и свойства мозга. В последние десятилетия интерес к нейронным сетям сильно вырос, исследуются сети различной структуры на основе различных моделей нейрона.

Одной из первых была предложена модель нейрона Мак-Каллока – Питтса, исследование динамики сети на основе которого проведено в работе. Для этой сети обнаружены такие параметры, при которых возможно хаотическое ее поведение.

Вычислительным экспериментом получено подтверждение того, что при некоторых значениях параметров сеть ведет себя хаотически. В статье приводится описание экспериментов и основные их результаты.

Статья состоит из нескольких разделов. В первом из них дается описание математической модели исследуемой нейронной сети. Одной из используемых при исследовании динамики данной сети характеристик является энтропия, которая описывается во втором разделе. В третьем разделе описываются некоторые методы анализа и отображения результатов.

На первом этапе численного эксперимента выявлены значения параметров, при которых динамика сети похожа на хаос. С целью подтверждения наличия хаотического поведения проведено исследование динамики нейронной сети при малых возмущениях начального состояния, в результате чего обнаружена чувствительность сети к такого рода возмущениям. Результаты анализа приведены в четвертом разделе.

1 Математическая модель сети

В этом разделе будет дано формальное описание используемой модели нейронной сети.

Рассмотрим формальный нейрон Мак-Каллока – Питтса [2], который определяется следующим образом.

На вход нейрона подаются n сигналов (вещественных чисел) x_1, x_2, \dots, x_n . Входные сигналы умножаются на вещественные синаптические веса w_1, w_2, \dots, w_n и суммируются. К результату добавляется величина смещения w_0 , полученное значение (состояние нейрона)

$$X = \sum_{i=1}^n w_i x_i + w_0$$

является аргументом активационной функции нейрона $f(X)$. В качестве активационной функции выбрана сигмоида следующего вида (рис. 1):

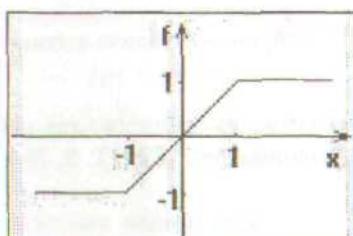


Рис. 1 Функция активации

Задается она уравнением

$$f(x) = \frac{|mx + 1| - |mx - 1|}{2}$$

при $m = 1$. Значение $f(X)$ является выходом нейрона.

Рассматриваемая в статье нейронная сеть состоит из N нейронов Мак-Каллока – Питтса, каждый из которых имеет N входов. Динамика такой сети может быть описана следующим рекуррентным соотношением:

$$X(t+1) = W \cdot F(X(t)) + I,$$

где $X(t) = (x_1(t), x_2(t), \dots, x_N(t))$ – N -мерный вектор состояния нейросети в момент времени t , представляющий собой совокупность состояний нейронов. Поведение сети определяется матрицей синаптических весов $W = \{w_{ij}\}$, состоящей из N строк и N столбцов (здесь w_{ij} – синаптический вес связи от i -го нейрона к j -му), N -мерным вектором смещения I , а также функцией активации $F(X) = (f(x_1), f(x_2), \dots, f(x_N))$ (x_i – компоненты вектора состояния сети X).

Исследуемая сеть по своей структуре аналогична рассмотренной в работе [1]. Рассмотрен случай $N = 3$. В качестве начального состояния выбирался нулевой вектор. Матрица синаптических весов имеет следующий вид:

$$W = \begin{pmatrix} 1 & -1 & 0 \\ 1 & p_0 & -1 \\ 0 & 1 & p_1 \end{pmatrix}$$

Здесь p_0 и p_1 – веса, принимающие в ходе анализа значения на отрезке $[-3, 3]$. Вектор смещения в общем случае может быть произвольным, в статье рассматривались ненулевые векторы с нормой меньше единицы.

2 Энтропия нейронной сети

Для определения хаотического поведения будет использоваться энтропия нейронной сети, которую, следуя работе [1], будем вычислять следующим образом.

Пространство состояний нейронной сети разбивается на элементарные ячейки сколь угодно малого диаметра δx (при этом считаем векторы состояния сети равными в том случае, если они попали в одну элементарную ячейку, в противном случае считаем их различными). Пусть вектор состояния нейросети за рассматриваемый период времени длины T_a (будем называть эту величину временем наблюдения) принимает N_s различных состояний s_j ($j = 0, \dots, N_s$). Частота попадания вектора состояния нейронной сети в состояние s_j равна $p_j = \frac{n_j}{T_a}$, где n_j – количество попаданий вектора состояния нейросети в данное состояние s_j за период наблюдения T_a .

Введем величину, которую будем считать оценкой энтропии:

$$H(N, T_a) = - \sum_{j=0}^{N_s} p_j \log p_j$$

Отметим, что, если вектор состояния нейросети совершает периодические колебания с длиной периода T , то $N_s = T$, а для любого состояния s_j имеем: $p_j \rightarrow \frac{1}{N_s}$, в этом случае $H(N, T_a) \rightarrow \log N_s$ (при $T_a \rightarrow \infty$). В случае хаотических колебаний, как легко убедиться, величина $H(N, T_a)$ будет неограниченно расти при увеличении параметра T_a и уменьшении δx .

По этой причине, чтобы иметь количественную характеристику поведения системы и в случае хаотической динамики, энтропия нейронной сети вводится иначе [1]: будем считать оценкой энтропии величину

$$h(N, T_a) = \frac{2^{H(N, T_a)}}{T_a}$$

Заметим, что при этом в случае периодической динамики нейронной сети $h(N, T_a) \rightarrow 0$ (при $T_a \rightarrow \infty$). В то же время для хаотической динамики характерны ненулевые значения энтропии. Легко видеть, что в случае хаотического поведения нейронной сети наибольшее значение энтропии равно 1.

В ходе численного эксперимента рассматривалось поведение не самого вектора состояния нейронной сети, а только его нормы (действительно, если норма вектора совершает хаотические колебания, то и

динамика вектора состояния сети также будет хаотической). При этом для приближенного вычисления энтропии ($H(\mathbf{N}, T_a)$, $h(\mathbf{N}, T_a)$) нейронной сети диапазон принимаемых нормой вектора состояния значений разбивался на элементарные отрезки длины $\delta x = 0.0001$. При вычислении энтропии в качестве начального состояния выбирался нулевой вектор, после чего сеть функционировала в течение времени t_0 (будем называть эту величину количеством "скрытых" итераций, при вычислении энтропии использовалось значение $t_0 = 512$). После этого в течение периода времени T_a вычислялось значение энтропии. Также в качестве показателя энтропии использовалось значение $N_s(T_a)$, то есть количество различных состояний, принимаемых нормой вектора состояний сети в течение периода наблюдения.

3 Способы отображения результатов

Для анализа нейронной сети большая часть параметров фиксируется, а некоторое количество (чаще всего один или два) параметров (синаптических весов) изменяются в некотором диапазоне значений, при этом для каждого набора параметров производятся наблюдения за вектором состояния нейросети, а также вычисляется ее показатель энтропии. Результаты такого численного анализа представимы в различных формах.

1. Диаграмма Фейгенбаума (Feigenbaum Plots, FP-диаграмма, бифуркационная диаграмма)

Для отображения результатов в таком виде один из весов выбирается в качестве свободного параметра и принимает различные значения на некотором отрезке, напротив этого свободного параметра откладываются значения евклидовой нормы вектора состояния нейросети, подсчитываемые в течение времени наблюдения T_a . При этом в случае периодических колебаний нормы вектора состояний сети напротив соответствующего параметра будем иметь некоторый набор точек, количество которых будет тем большим, чем больше будет период колебаний, в случае хаотических колебаний будут наблюдаваться полосы [1].

2. Одномерный график энтропии (1DEP's, one-dimensional entropic plots).

В этом случае напротив свободного параметра откладывается значение общего показателя энтропии (в работе напротив параметра откладывается число различных состояний нормы вектора состояния нейросети). При этом будем получать картину из линий различной высоты, высота их будет тем больше, чем больше показатель энтропии.

3. Двумерная карта энтропии (2DEM, two-dimensional entropic map).

Для отображения в таком виде уже оба параметра (p_0 и p_1) изменяются в некотором диапазоне. Каждая точка $p(x, y)$ на карте будет тем светлее, чем больше будет значение общего показателя энтропии для данных параметров ($p_0 = x$ и $p_1 = y$, соответствующими координатам (x, y) точки на карте). В работе оба параметра (p_0 и p_1) изменяются в диапазоне $-3 \leq p_0, p_1 \leq 3$. При этом точки будут тем светлее, чем выше значения энтропии сети при соответствующих значениях параметров.

Также для фиксированных параметров отображается зависимость нормы вектора состояния нейросети от текущей итерации, одна или две компоненты вектора состояния в течение периода наблюдения, а также проводится анализ последовательности значений вектора состояния и его нормы на периодичность.

4 Результаты анализа сети

Первая часть вычислительного эксперимента выявляла значения весов, при которых характерны высокие значения энтропии сети [1].

Нейронные сети, для которых получены высокие значения энтропии, были подвергнуты дополнительному анализу, в ходе которого проводилось исследование последовательности состояний нейронной сети на периодичность. Сети, периодичности в динамике которых обнаружено не было, рассматривались при малых возмущениях начального состояния, в результате чего отмечена чувствительность этих нейронных сетей к такого рода возмущениям.

4.1 Анализ двумерной карты энтропии

На рис. 2 приведено инверсное изображение двумерной карты энтропии (высоким значениям энтропии соответствуют более темные точки на диаграмме, и наоборот) для стандартных параметров (нулевой начальный вектор, вектор смещения $I = \begin{pmatrix} 0.02 \\ 0.04 \\ -0.03 \end{pmatrix}$, $t_0 = 512$, $T_a = 512$, $m = 1$; в дальнейшем по умолчанию будет использоваться этот набор параметров).

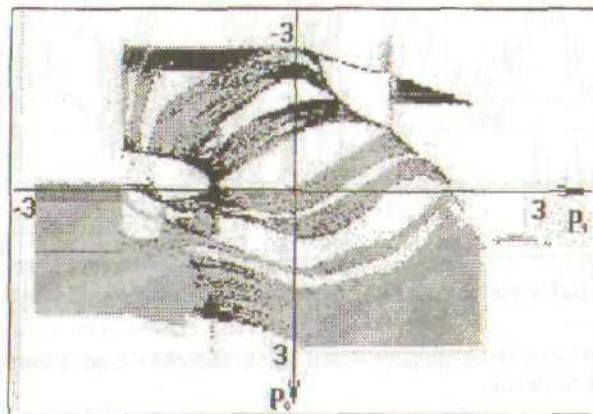


Рис. 2 Двумерная карта энтропии

На диаграмме видны области, соответствующие низким значениям показателя энтропии, а значит, короткопериодическим колебаниям (белые зоны на инверсном изображении двумерной карты энтропии), участки, соответствующие колебаниям вектора состояния нейросети с большой длиной периода (серые области), а также области, соответствующие параметрам (весам), при которых вектор состояния нейросети совершает хаотические колебания или колебания с длиной периода порядка времени наблюдения (черные участки на инверсном изображении, будем называть их зонами хаоса).

Для дополнительной проверки был проведен вычислительный эксперимент. В нем проводилось исследование двумерной карты энтропии при различных значениях параметров t_0 , m и I . Увеличение t_0 не внесло заметных изменений в двумерную карту энтропии. Предположительно, это можно связать с тем, что в случае периодической динамики сети с малой длиной периода (наиболее типичный случай) процесс стабилизации происходит достаточно быстро, а для сети с хаотической динамикой стабилизация не наступает вообще. В результате при анализе оказалось возможным ограничить количество "скрытых" итераций до $t_0 = 512$. После этого рассматривалось изменение двумерной карты энтропии при использовании функции активации с большим значением параметра m . В результате такого увеличения параметра m двумерная карта энтропии деформируется, зоны хаоса сокращаются, в основном преобладает периодическая динамика нейронной сети (причем со сравнительно небольшой длиной периода). Объяснить это можно тем, что увеличение параметра m приводит к сжатию функции активации вдоль оси абсцисс, поэтому компоненты вектора состояния нейросети чаще попадают за пределы отрезка $\left[-\frac{1}{m}, \frac{1}{m}\right]$, где обрезаются до величины 1 или -1, это ведет к тому, что значения компонент вектора состояния чаще повторяются в процессе функционирования, что способствует установлению периодических колебаний. Использование вектора смещения с большей нормой также приводит к сокращению зон хаоса. Изменение вектора начального состояния значительного влияния на вид двумерной карты энтропии не оказывает.

После дополнительной проверки было принято решение ограничиться вышеуказанными значениями основных параметров сети (t_0 , m , I).

4.2 Анализ одномерного графика энтропии и диаграммы Фейгенбаума

На графике энтропии заметны полосы из столбцов большой высоты, на двумерной карте энтропии это будет соответствовать пересечению зон хаоса прямой $p_0 = \text{const}$ ($p_1 = \text{const}$), участки с низкой высотой столбцов соответствуют прохождению такой секущей через зоны с низким значением энтропии (светлые области на инверсном изображении двумерной карты энтропии).

На диаграмме Фейгенбаума можно выделить черные вертикальные полосы, соответствующие высоким значениям энтропии (на одномерном графике энтропии таким значениям весов соответствуют полосы

из высоких столбцов), что соответствует хаотической динамике нейросети или колебаниям с периодом, сравнимым со временем наблюдения. На одномерном графике энтропии также заметны провалы между полосами из высоких столбцов.

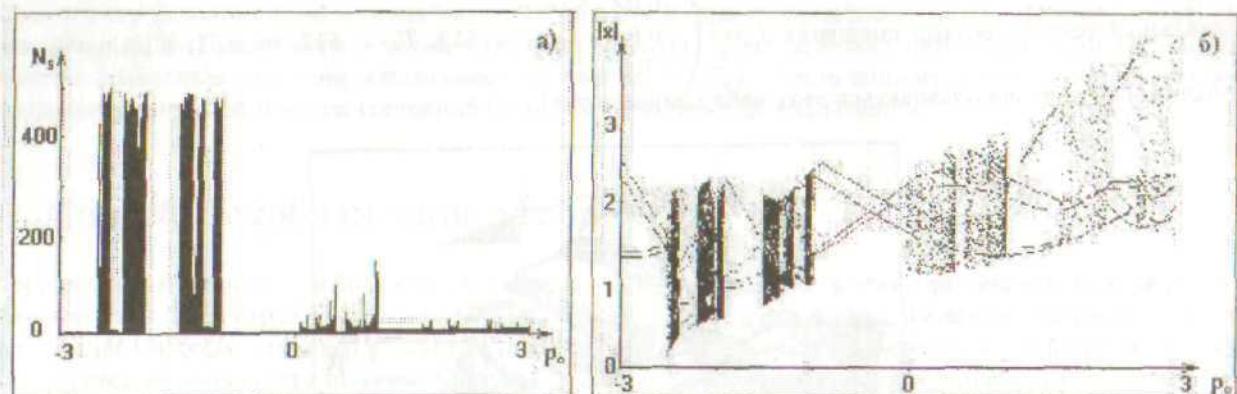


Рис. 3 Одномерный график энтропии и диаграмма Фейгенбаума

Рассмотренные выше диаграммы и графики также показывают возможность хаотического поведения сети при некоторых значениях весов.

4.3 Анализ возмущенной системы

В результате исследования двумерной карты энтропии, диаграммы Фейгенбаума и одномерного графика энтропии выявлены как значения весов, при которых сеть совершает периодические колебания, так и значения весов, при которых динамика сети похожа на хаос. Для таких значений весов проведено исследование последовательности векторов состояния сети на периодичность, в ходе которого количество "скрытых" итераций и время наблюдения увеличивалось (то есть $t_0 = 10000$, $T_a = 50000$ и более), в результате чего для некоторых значений параметров не удалось обнаружить периодичности изменения вектора состояния сети. Это позволяет предположить хаотический характер поведения вектора состояния нейронной сети при таких значениях синаптических весов.

Так как хаотическое движение характеризуется чувствительностью к изменению начальных данных и быстрым расхождением траекторий [3], то для дополнительной проверки проведен анализ поведения вектора состояния возмущенных систем при различных значениях параметров.

Опишем проведенный вычислительный эксперимент. Рассматривалась нейронная сеть при найденных значениях весов. При стандартных значениях основных параметров (нулевой вектор $X(0)$, $m = 1$,

$$I = \begin{pmatrix} 0.02 \\ 0.04 \\ -0.03 \end{pmatrix}$$

) сеть первоначально функционировала в течение периода времени $t_0 = 512$. Обозначим

вектор состояния нейросети после этого периода как X_0 ($X(t_0) = X_0$) и будем называть его начальным состоянием. Далее рассматривалось значение вектора состояния $\tilde{X}_0 = X_0 + \delta X_0$, то есть возмущенное состояние, отличающееся на величину δX_0 (вектор δX_0 будем называть возмущением исходного состояния, $|\delta X_0| \ll 1$). После этого в течение некоторого времени наблюдения отдельно рассматривается поведение одинаковых по своей структуре сетей, но с различными значениями начального состояния (X_0 и \tilde{X}_0 соответственно), векторы состояния этих нейронных сетей обозначим через $X(t)$ и $\tilde{X}(t)$. При этом исследовалось поведение разности векторов $X(t)$ и $\tilde{X}(t)$ (будем обозначать эту разность как $\delta X(t)$) в течение некоторого времени наблюдения T_a (в работе было использовано значение $T_a = 5000$). На рис.4 приведен график зависимости нормы вектора $\delta X(t)$ от текущей итерации при одном из найденных наборов весов.

Серия численных экспериментов показала, что для значений весов, при которых вектор состояния нейронной сети совершает периодические колебания, разность $X(t) - \tilde{X}(t)$ через несколько итераций сходится к нулевому вектору, далее $X(t)$ и $\tilde{X}(t)$ совпадают. При этом число итераций, за которое разность $X(t) - \tilde{X}(t)$ становилась равной нулевому вектору, в среднем было больше у нейросетей, совершающих периодические колебания с большим периодом.

Совершенно другие результаты получены при рассмотрении нейросетей с хаотической динамикой. Разность между векторами состояния $X(t)$ и $\tilde{X}(t)$ не обращалась в нуль даже при значительном увеличении времени наблюдения (до $T_a = 500000$).

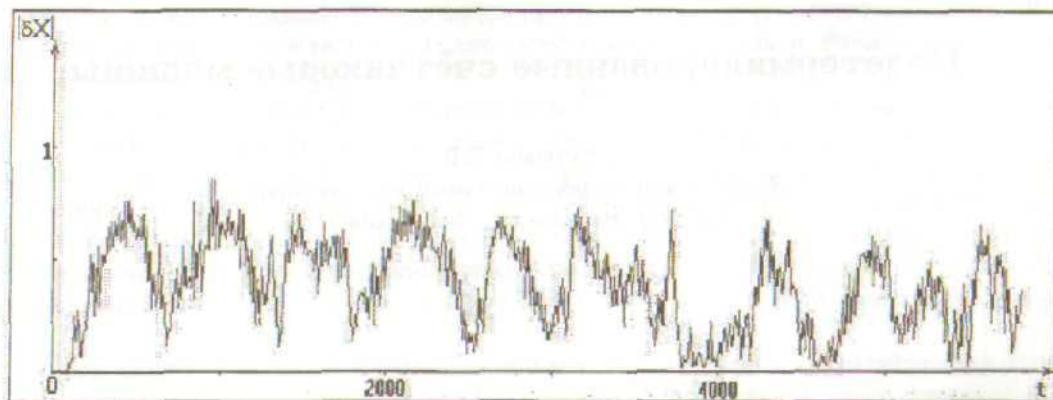


Рис. 4 График отклонения $\delta X(t)$ при $p_0 = -0.09$, $p_1 = -0.91$, $I = \begin{pmatrix} 0.001 \\ 0.001 \\ 0.001 \end{pmatrix}$

Таким образом, на основании проведенного анализа можно сделать вывод, что вектор состояния нейронной сети с периодической динамикой при малом возмущении исходного состояния сходится к вектору состояния сети с невозмущенным исходным состоянием; напротив, для сети с хаотической динамикой разница между векторами состояния при возмущенном и невозмущенном исходном состоянии не сходится к нулю.

5 Заключение

В ходе анализа данной модели нейронной сети получены следующие результаты:

- Получено экспериментальное подтверждение возможности возникновения хаотического поведения в простой модели нейронной сети.
- Рассмотрено различие в динамике нейронной сети при различных значениях основных параметров (t , I). Зоны хаоса при увеличении нормы вектора смещения и сжатии функции активации вдоль оси абсцисс сокращаются. При увеличении параметра t активационной функции двумерная карта энтропии деформируется.
- Вектор состояния сети с периодической динамикой устойчив к малым возмущениям исходного состояния, что позволяет предположить прямую зависимость времени стабилизации возмущенной системы от периода колебаний вектора состояния сети.
- Вектор состояния сети при значениях весов, для которых значения показателя энтропии высоки, чувствителен к малым изменениям вектора исходного состояния, что дает еще одно численное подтверждение хаотического поведения нейронной сети при таких значениях весов.

Анализ данной модели показал, что нейросеть при значениях параметров, для которых показатель энтропии высок и динамика нейронной сети непериодическая, можно использовать на практике (к примеру, в качестве генератора хаоса). Также затронута проблема чувствительности моделей нейронных сетей к входным данным. Подобного рода исследования можно провести и для нейронных сетей с другой структурой, в том числе и на основе других моделей нейрона. Как несложно убедиться, в этом направлении открыты широкие возможности для дальнейших исследований.

Литература

- [1] Radu Dogaru, A.T.Murgan, Daniel Ioan. Robust Oscillations and Bifurcations in Cellular Neural Networks//Proceedings of IEEE Int. Workshop on Cellular Neural Networks and Their Applications, (CNNA'94), pp.297-302, Rome, 1994,
- [2] Уоссермен Ф. Нейрокомпьютерная техника: теория и практика. М.: Мир, 1992. 240 с.
- [3] Мелик-Гайказян И.В., Мелик-Гайказян Н.В., Тарасенко В.Ф. Методология моделирования нелинейной динамики сложных систем. М.: ФИЗМАТЛИТ, 2001. 272 с.

Недетерминированные счётчиковые машины

Кузьмин Е.В.

Ярославский государственный университет

150 000, Ярославль, Советская, 14

получена 11 июля 2003

В работе рассматриваются недетерминированные счётчиковые машины, использующиеся как общее средство для демонстрации неразрешимости ряда проблем для систем, способных моделировать эти машины.

1. Введение

В настоящее время большое внимание уделяется моделированию, анализу поведенческих свойств и верификации программных, аппаратных, технологических систем и процессов. Особый интерес в этом направлении представляют распределённые системы, характеризующиеся отсутствием централизованного управления работой системы. Такие системы обычно моделируются системами переходов с конечным числом управляющих состояний и с различными видами переменных и структур данных, таких как счётчики, очереди и т.д. Многие системы переходов можно рассматривать как машины Тьюринга (системы равномощные машинам Тьюринга) с ослаблениями, такими как наличие недетерминизма в правилах переходов из одной конфигурации в другую или же включения отношения потери. Чтобы оценить уровень ослабления для некоторых классов систем, вводятся абстрактные машины, с помощью которых определяется круг задач, которые по-прежнему остаются неразрешимыми. К таким абстрактным машинам можно отнести счётчиковые машины с потерями [6], которые были введены для систем с потерями, использующихся для моделирования передачи данных через ненадёжные каналы связи (например, FIFO-канальные системы с потерями).

В данной работе рассматриваются системы переходов, где каждый переход определяется недетерминированно в соответствии с управляющими состояниями и независимо от манипулируемых данных. Вводятся недетерминированные счётчиковые машины в качестве общего средства для демонстрации неразрешимости ряда проблем для систем, которые могут моделировать эти машины, в частности, для взаимодействующих раскрашивающих автоматов [5], использующихся для моделирования перемещения данных различного типа между компонентами распределённой системы. Показывается неразрешимость для недетерминированных счётчиковых машин проблем ограниченности, достижимости, эквивалентности и включения. Из чего можно сделать вывод, что даже при таком ослаблении (недетерминизм переходов) могут строиться довольно выразительные системы, такие как взаимодействующие раскрашивающие автоматы.

2. Счётчиковые машины

МашинаМинского или счётчиковая машина M – это набор $(\{q_0, \dots, q_n\}, \{x_1, \dots, x_m\}, \{\delta_0, \dots, \delta_{n-1}\})$, где x_i – счётчик, q_i – состояние, q_0 – начальное состояние, q_n – финальное (заключительное) состояние, δ_i – правило переходов для q_i ($0 \leq i \leq n-1$).

Состояния q_i , $0 \leq i \leq n-1$, подразделяются на два типа. Состояния первого типа имеют правила переходов вида ($1 \leq j \leq m$, $0 \leq k \leq n$):

$$\delta : x := x + 1; \text{ goto } q .$$

Для состояний второго типа имеем ($0 \leq k' \leq n$):

$$\delta : \text{if } x > 0 \text{ then } (x := x - 1; \text{ goto } q) \text{ else goto } q' .$$

Конфигурация машины Минского – это набор (q_i, c_1, \dots, c_m) , где q_i – состояние машины, c_1, \dots, c_m – натуральные числа, являющиеся значениями соответствующих счётчиков. Размер конфигурации определяется как $\text{size}((q, c_1, \dots, c_m)) = \sum_{i=1}^m c_i$.

Исполнение машины – это последовательность конфигураций $s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow \dots$, начинающаяся в начальной конфигурации s_0 и индуктивно определяемая в соответствии с правилами переходов. Необходимо отметить, что исполнение детерминировано, т.к. каждое состояние имеет не более одного правила переходов. Машина Минского останавливается, если исполнение содержит конфигурацию с состоянием q_n , т.е. достигает финального состояния. Поскольку машина Минского уже всего с двумя счётчиками может моделировать машину Тьюринга, проблема (останова) достижения финального состояния из некоторой начальной конфигурации (q_0, c_1, c_2) для машины Минского с двумя счётчиками является неразрешимой [8].

Говорят, что машина Минского ограничена при некоторой начальной конфигурации тогда и только тогда, когда существует натуральное число c' такое, что на протяжении всего исполнения машины из начальной конфигурации выполняется условие $x_1 + \dots + x_m \leq c'$, т.е. в любой момент времени сумма значений счётчиков не превосходит c' .

Утверждение 1. Для 3-х счётчиковой машины Минского проблема ограниченности неразрешима.

Доказательство. Доказательство осуществляется сведением проблемы останова для машины Минского с двумя счётчиками к проблеме ограниченности машины Минского с тремя счётчиками. Рассмотрим некоторую машину Минского $2cM = (\{q_0, \dots, q_n\}, \{x_1, x_2\}, \{\delta_0, \dots, \delta_{n-1}\})$. Преобразуем $2cM$ в машину Минского $3cM$ с тремя счётчиками. Добавим новый счётчик x_3 . И для каждого состояния $q_i, 0 \leq i \leq n-1$, добавим состояния q'_i, q''_i , правила переходов δ'_i, δ''_i и изменим δ_i следующим образом. Для состояния первого типа имеем: $(\delta) q : x := x + 1; \text{ goto } q'$; $(\delta') q' : x := x + 1; \text{ goto } q$. Для состояния второго типа: $(\delta) q : \text{if } x > 0 \text{ then } (x := x - 1; \text{ goto } q') \text{ else goto } q''$; $(\delta') q' : x := x + 1; \text{ goto } q$; $(\delta'') q'' : x := x + 1; \text{ goto } q'$.

Машина Минского $2cM$ имеет конечное исполнение тогда и только тогда, когда машина Минского $3cM$ ограничена, т.к. счётчик x_3 , использующийся для подсчёта количества шагов исполнения машины $2cM$, ограничен только в случае конечного исполнения. \square

Говорят, что машина Минского M *тотально ограничена*, если и только если она ограничена для всех возможных начальных конфигураций.

Утверждение 2. Для 3-х счётчиковой машины Минского тотальная ограниченность неразрешима.

Доказательство. Доказательство проводится сведением рассмотренной выше проблемы ограниченности машины Минского с тремя счётчиками и начальной конфигурацией $(q_0, 0, 0, 0)$ к проблеме тотальной ограниченности. Рассмотрим некоторую машину Минского $M = (\{q_-, \dots, q_+\}, \{x_-, x_+, x_0\}, \{\delta_-, \dots, \delta_+\})$. Преобразуем M в машину Минского M' следующим образом. Добавим новые состояния q'_0, q'_1, q'_2 и правила переходов $\delta', \delta'', \delta'$:

$$\begin{aligned} (\delta') q' : \text{if } x_- > 0 \text{ then } (x_- := x_- - 1; \text{ goto } q') \text{ else goto } q'; \\ (\delta') q' : \text{if } x_+ > 0 \text{ then } (x_+ := x_+ - 1; \text{ goto } q') \text{ else goto } q'; \\ (\delta') q' : \text{if } x_0 > 0 \text{ then } (x_0 := x_0 - 1; \text{ goto } q') \text{ else goto } q . \end{aligned}$$

Машина Минского M ограничена тогда и только тогда, когда машина $M' = (\{q'_-, q'_+, q'_0, q_-, \dots, q_+\}, \{x_-, x_+, x_0\}, \{\delta', \delta'', \delta', \delta_-, \dots, \delta_+\})$, где q'_0 – начальное состояние, totally ограничена, т.к. правила переходов $\delta', \delta'', \delta'$ обнуляют любые начальные значения счётчиков и переводят M' в состояние q_0 . \square

Определим *отношение потери* [6], которое описывает самопроизвольное изменение конфигураций машины Минского. Обозначим \rightarrow_s (от ‘sum’) отношение на конфигурациях n -счётчиковой машины, которое определяется следующим образом:

$$(q, c_1, \dots, c_m) \rightarrow_s (q', c'_1, \dots, c'_m) \iff q = q' \wedge \left(\forall i : c_i = c'_i \vee \sum_{i=1}^m c_i > \sum_{i=1}^m c'_i \right).$$

Это отношение означает, что либо никаких изменений не происходит, либо сумма значений всех счётчиков уменьшается. Пусть id – отношение тождества. Отношение \rightarrow_l (от ‘lossy’) является *отношением потери* тогда и только тогда, когда $id \subseteq \rightarrow_l \subseteq \rightarrow_s$.

Счётчиковая машина с потерями (*LCM* – lossy counter machine) [6] определяется на основе обычной счётчиковой машины M добавлением отношения потери \rightarrow_l . Пусть \rightarrow – обычное отношение переходов на множестве конфигураций машины M . Отношение переходов с потерей \Rightarrow для *LCM* определяется как

$$s_1 \Rightarrow s_2 \iff \exists s'_1, s'_2 : s_1 \rightarrow_l s'_1 \rightarrow s'_2 \rightarrow_l s_2.$$

Примером счётчиковой машины с потерями может служить *счётчиковая машина с обнулениями* (*RCM* – reset counter machine). Предполагается, что в состоянии второго типа машины Минского при проверке на ноль счётчик может внезапно обнуляться. И тогда счётчиковая машина с обнулениями *rcM* для состояний q_i второго типа имеет следующее правило переходов:

$$\delta : \text{if } x > 0 \text{ then } (x := x - 1; \text{goto } q) \sqcup x := 0; \text{goto } q,$$

где \sqcup – оператор недетерминированного выбора.

Определение 1. Недетерминированная счётчиковая машина (*ncM*) *ncM* – это набор из четырёх элементов $(Q, \{x_1, \dots, x_m\}, \rightarrow, T)$, где x_i – счётчик, $Q = \{q_0, \dots, q_n\}$ – конечное множество состояний, q_0 – начальное состояние, q_n – финальное (заключительное) состояние, T – конечное множество меток действий, соответствующих выражению над счётчиками, $\rightarrow \subseteq Q \setminus \{q_n\} \times T \times Q$ – отношение переходов, запись $q \xrightarrow{t} q'$ используется для обозначения (правила) перехода $(q, t, q') \in \rightarrow$.

Переход из одной конфигурации в другую может быть одного из следующих видов (при $q \xrightarrow{t} q'$):
 $(q, c_1, \dots, c_i, \dots, c_m) \xrightarrow{t} (q', c_1, \dots, c_i + 1, \dots, c_m)$, если t соответствует выражению $x := x + 1$;
 $(q, c_1, \dots, c_i, \dots, c_m) \xrightarrow{t} (q', c_1, \dots, c'_i, \dots, c_m)$, где $c'_i = 0$ при $c_i = 0$ и $c'_i = c_i - 1$ при $c_i > 0$, если метка t соответствует выражению $x := x \ominus 1$;
 $(q, c_1, \dots, c_i, \dots, c_m) \xrightarrow{t} (q', c_1, \dots, 0, \dots, c_m)$, если t соответствует выражению $x := 0$;
 $(q, c_1, \dots, c_i, \dots, c_m) \xrightarrow{t} (q', c_1, \dots, c'_j, \dots, c_m)$, где $c'_j = c_i + 1$ при $c_j \geq 1$ и $c'_j = c_i$ при $c_j = 0$, если метка t соответствует выражению $x := x + \min(x, 1)$.

Для некоторой счётчиковой машины M и начальной конфигурации s_0 обозначим $R(M)$ множество достижимых из s_0 векторов значений счётчиков. Тогда $R_m(M)$ – множество проекций векторов из $R(M)$ на их первые m координат, т.е. $R_m(M) = \{(c_1, \dots, c_m) \mid (c_1, \dots, c_m, \dots, c_n) \in R(M)\}$.

Утверждение 3. Для любой m -счётчиковой машины с обнулениями *rcM* можно построить недетерминированную $(m+1)$ -счётчиковую машину *ncM* такую, что для любого $\bar{c} \in R(\text{rcM})$ существует $\bar{c}' \in R_m(\text{ncM})$ такой, что $\bar{c} \leq \bar{c}'$.

Доказательство. Рассмотрим произвольную m -счётчиковую машину с обнулениями *rcM*. Предположим, что для состояний второго типа q_i в правиле переходов $\delta : \text{if } x > 0 \text{ then } (x := x - 1; \text{goto } q) \sqcup x := 0; \text{goto } q$ при проверке счётчика $x_i = 0$ на нуль может произойти ошибка и первым сработает переход, соответствующий левой части правила. Добавим “исправляющие команды” таким образом, что если произошла ошибка, то в дальнейшем при любом переходе значения счётчиков не будут увеличиваться.

Заменим в машине *rcM* все выражения вида $x := x + 1$ (соответствующие правилам переходов для состояний первого типа) на выражения $x := x + \min(z, 1)$, где z – новый счётчик, и для каждого состояния q_i второго типа заменим переход $\text{if } x > 0 \text{ then } (x := x - 1; \text{goto } q)$ на $(q : z := z + \min(x, 1); \text{goto } q'); (q' : z := z \ominus 1; \text{goto } q''); (q'' : x := x \ominus 1; \text{goto } q)$, где q', q'' – новые состояния. Полагаем, что начальное значение счётчика $z = 1$.

Таким образом, для машины *rcM* с начальной конфигурацией (q_0, c_1, \dots, c_m) мы построили недетерминированную машину *ncM* с начальной конфигурацией $(q_0, c_1, \dots, c_m, 1)$ такую, что множество проекций достижимых машиной *ncM* векторов значений счётчиков при отсутствии ошибок $R_m^{\text{true}}(\text{ncM})$ равно $R(\text{rcM})$, а в случае ошибок для любого $\bar{c} \in R_m^{\text{false}}(\text{ncM})$ существует $\bar{c}' \in R_m^{\text{true}}(\text{ncM})$ такой, что $\bar{c} \leq \bar{c}'$. \square

3. Проблема ограниченности

Утверждение 4. Для машины Минского *3cM* с тремя счётчиками и нулевой начальной конфигурацией, т.е. $(q_0, 0, 0, 0)$, можно построить недетерминированную 4-х счётчиковую машину *n4cM* такую, что *n4cM* будет ограничена тогда и только тогда, когда машина *3cM* ограничена.

Доказательство. Преобразуем машину Минского $3cM$ в машину Минского $4cM$ с четырьмя счётчиками (и с несколько изменённым набором правил переходов) следующим образом. Добавим новый “ёмкостный” счётчик K . Заменим в каждом правиле перехода команду $x := x - 1$ на команды $x := x - 1; K := K + 1$ и $x := x + 1$ на “управляемое приращение” ($i = 1, 2, 3$):

`if $K > 0$ then ($K := K - 1; x := x + 1$) else goto q_+ .`

Наконец, добавим четыре новых состояния $q_{n+1}, q_{n+2}, q_{n+3}, q_{n+4}$ и соответствующие им правила переходов $\delta_+, \delta_+, \delta_+, \delta_+$:

`q_+ : if $x > 0$ then ($x := x - 1; K := K + 1; \text{goto } q_+$) else goto q_+ ,`

`q_+ : if $x > 0$ then ($x := x - 1; K := K + 1; \text{goto } q_+$) else goto q_+ ,`

`q_+ : if $x > 0$ then ($x := x - 1; K := K + 1; \text{goto } q_+$) else goto q_+ ,`

`q_+ : $K := K + 1; \text{goto } q_+$.`

В машине $4cM$ только правило перехода δ_+ изменяет (увеличивает на единицу) сумму счётчиков $x + x + x + K$.

Машина $4cM$ с начальной конфигурацией $(q_0, 0, 0, 0, 0)$ ограничена тогда и только тогда, когда ограничена машина Минского $3cM$ с начальной конфигурацией $(q_0, 0, 0, 0)$ (наличие у машины Минского $3cM$ бесконечного неограниченного исполнения означает наличие бесконечного неограниченного исполнения и у машины $4cM$, т.к. значение счётчика K (изначально $K = 0$) бесконечно возрастает из-за прохождения бесконечного числа раз через состояние q_{n+4} , где правило перехода δ_+ увеличивает значение K на 1).

Разрешим ещё одну команду $x := 0$ для правил переходов. И добавим новый оператор недетерминированного выбора \sqcup . В $4cM$ заменим каждое правило перехода `if $x > 0$ then comm else comm`, где $x \in \{x_1, x_2, x_3, K\}$, на правило `comm \sqcup $x := 0$; comm`. В `comm` заменим каждую последовательность команд $x := x - 1; K := K + 1$ на

$$K := K + \min(x, 1); K := K \ominus 1; K := K + \min(x, 1); x := x \ominus 1;$$

и каждую последовательность команд $K := K - 1; x := x + 1$ на

$$x := x + \min(K, 1); x := x \ominus 1; x := x + \min(K, 1); K := K \ominus 1.$$

Таким образом, мы получили недетерминированную машину $n4cM$. И в машине $n4cM$ только δ_+ увеличивает сумму счётчиков $x + x + x + K$ в то время, как все остальные правила переходов уменьшают или оставляют без изменения эту сумму.

Машина $n4cM$ ограничена, если все её исполнения ограничены. Машина $n4cM$ с начальной конфигурацией $(q_0, 0, 0, 0, 0)$ ограничена тогда и только тогда, когда ограничена машина $4cM$ с начальной конфигурацией $(q_0, 0, 0, 0, 0)$. Действительно, если существует некоторое неограниченное исполнение машины $n4cM$, то оно бесконечно часто проходит через состояние q_{n+4} , также как и через состояния $q_{n+3}, q_{n+2}, q_{n+1}, q_0$. И, поскольку исполнение неограниченное, то существует такое неограниченное исполнение машины $n4cM$, что в нём все команды $x := 0$ выполняются только тогда, когда $x = 0$, что означает существование неограниченного исполнения машины $4cM$.

Легко видеть, что машина $n4cM$ представляет собой недетерминированную 4-х счётчиковую машину. \square

Примечание. Доказательство данного утверждения останется верным, если команду $x := 0$ заменить командой $x := x \ominus 1$.

Следствие 1. Проблема ограниченности для недетерминированной счётчиковой машины неразрешима.

Доказательство. Из утверждений 1 и 4 следует неразрешимость задачи ограниченности для недетерминированной счётчиковой машины с начальной конфигурацией $(q_0, 0, \dots, 0)$.

Очевидно, что для любой недетерминированной счётчиковой машины ncM с начальной конфигурацией $(q_0, 0, \dots, 0)$ может быть построена недетерминированная машина ncM' с начальной конфигурацией (q'_0, c_1, \dots, c_m) ($m \geq 4$) такая, что ncM' будет ограничена тогда и только тогда, когда ограничена ncM . Действительно, это может быть достигнуто добавлением к ncM' новых состояний (в том числе и нового начального q'_0) и переходов с выражениями $x := x \ominus c_1, \dots, x := x \ominus c_j$, ведущих от q'_i к q_j , $x := x \ominus c_j$ означает применение команды $x := x \ominus 1$ c_j раз. \square

Утверждение 5. Проблема тотальной ограниченности для недетерминированной 4-х счётчиковой машины неразрешима.

Доказательство. Сведём проблему ограниченности для недетерминированной 4-х счётчиковой машины $n4cM$ и начальной конфигурацией $(q_0, 0, 0, 0, 0)$ к проблеме тотальной ограниченности.

Построим недетерминированную машину $n4cM'$ добавлением к $n4cM$ новых состояний (в том числе и нового начального q'_0) и переходов с выражениями $x := 0, \dots, x := 0$, ведущих от q'_0 к состоянию q_0 .

Очевидно, что машина $n4cM'$ totally ограничена тогда и только тогда, когда ограничена машина $n4cM$ с начальной конфигурацией $(q_0, 0, 0, 0, 0)$.

Запретим на переходах выражения вида $x := 0$. Тогда, используя построение из доказательства утверждения 4 для любой машины Минского, можно построить недетерминированную машину, которая будет totally ограничена тогда и только тогда, когда totally ограничена машина Минского. \square

4. Проблемы включения и эквивалентности

Пусть задан класс недетерминированных счётчиковых машин, которые имеют одно и то же множество счётчиков. Проблема включения состоит в определении существования алгоритма, устанавливающего для любых двух недетерминированных счётчиковых машин ncM_1 и ncM_2 из этого класса, имеет ли место соотношение $R(ncM_1) \subseteq R(ncM_2)$, где $R(ncM_i)$ – множество достижимых векторов значений счётчиков из начальной конфигурации для машины ncM_i ($i = 1, 2$). В случае проблемы эквивалентности соотношение имеет вид $R(ncM_1) = R(ncM_2)$.

Пусть дан полином P от n переменных с целыми коэффициентами; существует ли такой вектор целых (c_1, c_2, \dots, c_n) , что $P(c_1, c_2, \dots, c_n) = 0$? Уравнение $P(x_1, x_2, \dots, x_n) = 0$ называется диофантовым. В общем оно представляет собой сумму членов: $P(x_1, \dots, x_n) = \sum_i P_i(x_1, \dots, x_n)$, где $P_i(x_1, x_2, \dots, x_n) = a_i \cdot x_{s_1} \cdot x_{s_2} \cdots x_{s_k}$.

Граф $G(P)$ диофантова полинома $P(x_1, \dots, x_n)$ с неотрицательными коэффициентами – это множество $G(P) = \{(x_1, \dots, x_n, y) \mid y \leq P(x_1, \dots, x_n) \wedge 0 \leq x_1, \dots, x_n, y\}$.

Задача включения графов полиномов заключается в определении для двух диофантовых полиномов A и B , выполняется ли $G(A) \subseteq G(B)$.

Утверждение 6. Задача включения графов полиномов неразрешима.

Доказательство этого утверждения проводится сведением десятой проблемы Гильберта к задаче включения графов полиномов (например, см. [7]).

4.1. Слабое вычисление

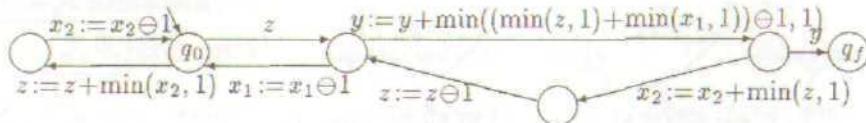
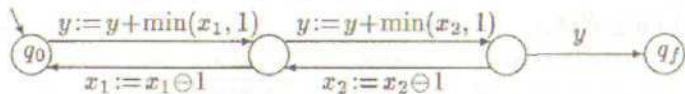
Необходимо показать, что недетерминированные счётчиковые машины могут (в определённом смысле) вычислять значение полинома $P(x_1, x_2, \dots, x_n)$. Ограничим полином P до неотрицательных значений полинома, неотрицательных коэффициентов и неотрицательных значений переменных. Это позволяет нам представить значения переменных и значение полинома значениями счётчиков недетерминированной машины. Входные значения переменных x_1, x_2, \dots, x_n задаются начальной конфигурацией машины. Выходная переменная y будет использоваться для хранения значения полинома, где $y \leq P(x_1, x_2, \dots, x_n)$. Вычисление полинома начинается в начальной конфигурации $(q_0, x_1, \dots, x_n, 0, \dots, 0)$ недетерминированной машины и заканчивается в финальной конфигурации $(q^f, x_1, \dots, x_n, y, \dots, c)$, где $y \leq P(x_1, x_2, \dots, x_n)$, т.е. машина попадает в состояние, из которого нет переходов.

Недетерминированная счётчиковая машина будет слабо вычислять значение $P(x_1, \dots, x_n)$. Слабое вычисление означает, что вычисленное значение не будет превышать $P(x_1, \dots, x_n)$, но может быть любым (неотрицательным) значением, меньшим $P(x_1, \dots, x_n)$.

Сейчас мы хотим показать, что можно построить недетерминированную машину, слабо вычисляющую функцию умножения (двух чисел). На её основе мы можем построить составную машину, которая слабо вычисляет значение каждого члена путём последовательной композиции машин умножения и затем суммирует результаты.

Недетерминированная машина умножения приводится на рисунке 1. В самом лучшем случае эта машина вычисляет значение $y = x_1 \cdot x_2$, в остальных $y < x_1 \cdot x_2$. На рис. 1 переход с выражением $y := y + \min((\min(z, 1) + \min(x'_1, 1)) \ominus 1, 1)$ соответствует набору переходов $k := k + \min(x_1, 1)$; $k := k + \min(z, 1)$; $k := k \ominus 1$; $y := y + \min(k, 1)$; $k := k \ominus 1$.

Машине сложения представлена на рис. 2. Недетерминированная счётчиковая машина, слабо вычисляющая некоторый одночлен $R_i = a_i \cdot x_{s_1} \cdot x_{s_2} \cdots x_{s_k}$ диофантова уравнения, показана на рис. 3 и представляет собой комбинацию машин умножения.

Рис. 1. Недетерминированная счётчиковая машина, слабо вычисляющая $y = x_1 \cdot x_2$ Рис. 2. Недетерминированная счётчиковая машина, слабо вычисляющая $y = x_1 + x_2$

Лемма 1. Для любого полинома $P(x_1, \dots, x_n)$ с неотрицательными целыми коэффициентами можно построить недетерминированную счётчиковую машину M такую, что $G(P) = R_n(M)$, где $G(P)$ – граф полинома P , $R_n(M)$ – множество проекций векторов из множества достижимости $R(M)$ на их первые $n+1$ координаты. (Множество проекций $R_n(M) = \{(c_1, \dots, c_n, c_{n+1}) \mid (c_1, \dots, c_n, c_{n+1}, \dots, c_k) \in R(M)\}$).

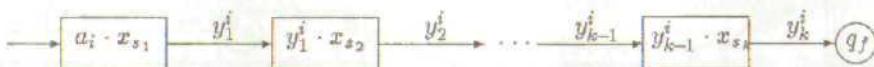
Доказательство. Основой для построения недетерминированной счётчиковой машины M (см. рис. 4) служит недетерминированная машина M_P , слабо вычисляющая полином P . К этой машине добавляются новые состояния (в том числе и новое начальное), переходы и счётчики x_1, \dots, x_n . Для каждого счётчика x_i создаются следующие правила переходов: $(q_1^i, x_i := x_i + 1, q_2^i), (q_2^i, x_i^1 := x_i^1 + 1, q_3^i), \dots, (q_{s+1}^i, x_i^s := x_i^s + 1, q_1^i), (q_1^i, x_i, q_1^{i+1})$, где x_i^s – счётчик недетерминированной машины M_{P_s} , слабо вычисляющей одночлен P_s , который соответствует переменной x_i . Для нового начального состояния q_0 построим переход $q_0 \rightarrow q_1^1$, а для старого начального состояния q_0' переход $q_s^n \rightarrow q_0'$. Эти переходы позволяют задавать произвольные значения для счётчиков x_1, \dots, x_n (и для счётчиков x_i^s каждой машины M_{P_s} , причём $x_i^s = x_i$) и передают управление машине M_P . Зададим на множестве счётчиков недетерминированной машины M порядок так, чтобы первыми $n+1$ счётчиками были x_1, \dots, x_n, y , где y – счётчик, в который помещается результат от слабого вычисления полинома P .

Для начальной конфигурации s_0 машины M установим значения счётчиков равными нулю, за исключением счётчиков a , использующихся для хранения констант a_i одночленов P_i , которые положим равными соответствующим константам a_i . Тогда для недетерминированной счётчиковой машины M с начальной конфигурацией s_0 может быть получено множество $R_n(M)$ проекций достижимых векторов значений счётчиков на счётчики x_1, \dots, x_n, y , которое будет равно графу $G(P)$ полинома P . \square

4.2. Проблема включения

Теорема 1. Проблема включения неразрешима.

Доказательство. Для двух произвольных полиномов P и S с одинаковым числом переменных n можно построить недетерминированные счётчиковые машины M_P и M_S такие, что $G(P) \subseteq G(S) \iff R(M_P) \subseteq R(M_S)$. Пусть M_P и M'_S – недетерминированные счётчиковые машины, слабо вычисляющие полиномы P

Рис. 3. Недетерминированная счётчиковая машина, слабо вычисляющая одночлен диофанта уравнения $P_i = a_i \cdot x_{s_1} \cdot x_{s_2} \cdot \dots \cdot x_{s_k}$. Каждый блок представляет собой автомат умножения

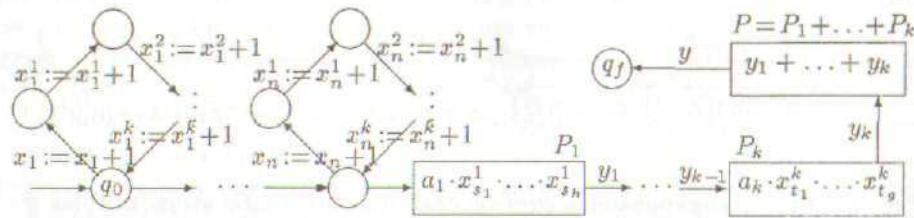


Рис. 4. Недетерминированная счётчиковая машина M , порождающая граф $G(P)$ полинома P , где $G(P) = \{(x_1, \dots, x_n, y) \mid y \leq P(x_1, \dots, x_n)\}$

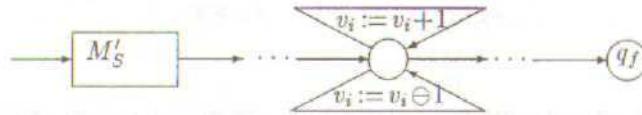


Рис. 5. Недетерминированная счётчиковая машина M_S , где $v_i \in V \setminus \{x_1, \dots, x_n, y\}$ ($1 \leq i \leq r - n - 1$)

и соответственно S , причём $G(P) = R_n(M_P)$ и $G(S) = R_n(M'_S)$ (см. лемму 1). Если первая машина имеет k счётчиков, а вторая l , то число счётчиков в обеих машинах можно уравнять, добавив $|k - l|$ счётчиков для машины с меньшим их количеством. Эти счётчики не входят в выражения каких-либо переходов, начальное значение каждого из них полагается равным 0.

Недетерминированная машина M_S строится на основе машины M'_S так, как показано на рис. 5. К машине M'_S добавляются новые состояния и переходы. Для каждой машины различаем счётчики x_1, \dots, x_n, y . Остальные счётчики, которые будем называть внутренними, переименуем как $v_1, v_2, \dots, v_{r-n-1}$, где r – мощность множества счётчиков V . Для каждого внутреннего счётчика $v_i \in V \setminus \{x_1, \dots, x_n, y\}$ недетерминированной машины M'_S построим новые переходы: $(q'_i, v_i := v_i + 1, q'_i)$ и $(q'_i, v_i := v_i \oplus 1, q'_i)$, где $1 \leq i \leq r - n - 1$. После слабого вычисления полинома S машиной M'_S и помещения результата вычисления в счётчик y , эти новые переходы позволяют установить произвольное значение для каждого внутреннего счётчика.

Таким образом, $R(M_P) \subseteq R(M_S) \iff R(M_P) \subseteq R_n(M'_S) \times N^{r-n-1}$, где N^{r-n-1} – множество векторов размерности $r - n - 1$ из целых неотрицательных чисел. Правая часть эквивалентна отношению $R_n(M_P) \subseteq R_n(M'_S)$. Следовательно, $R(M_P) \subseteq R(M_S) \iff R_n(M_P) \subseteq R_n(M'_S)$.

Это означает, что для решения проблемы включения графов двух произвольных полиномов P и S достаточно построить недетерминированные счётчиковые машины M_P, M'_S, M_S , убедиться, имеет ли место включение $R(M_P) \subseteq R(M_S)$, на основании этого сделать заключение о том, имеет ли место $R_n(M_P) \subseteq R_n(M'_S)$, и на основании леммы 1 установить, содержит ли $G(P)$ в $G(S)$. Таким образом, если предположить разрешимость проблемы включения для множеств достижимых векторов значений счётчиков недетерминированных машин, получаем разрешимость проблемы включения графов полиномов с неотрицательными коэффициентами, что не верно. \square

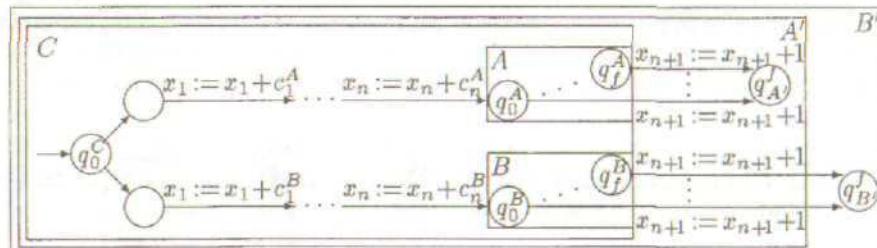


Рис. 6. Переход с выражением $x := x + c$ соответствует последовательности из c переходов $x := x + 1$

4.3. Проблема эквивалентности

Теорема 2. Проблема эквивалентности неразрешима.

Доказательство. Для доказательства неразрешимости методом сведения достаточно показать, что проблема включения сводится к проблеме эквивалентности. Это означает, что нужно указать способ, с помощью которого для любых двух недетерминированных счётчиковых машин A и B можно построить машины A' и B' такие, что $R(B) \subseteq R(A) \iff R(A') = R(B')$.

Обе машины A' и B' строятся добавлением нового счётчика и разных переходов к одной и той же машине C , которая, в свою очередь, строится из машин A и B .

Способ конструирования машин C , A' , B' показан на рис. 6. Предполагается, что машины A и B имеют одно и то же множество счётчиков (в противном случае добавляются новые счётчики, как в доказательстве теоремы 1). Пусть A и B имеют начальные конфигурации $(q_0^A, c_1^A, \dots, c_n^A)$ и соответственно $(q_0^B, c_1^B, \dots, c_n^B)$. Машина C строится таким образом, что из начального состояния q_0^C существуют две последовательности переходов в начальные состояния q_0^A машины A и соответственно q_0^B машины B такие, что из начальной конфигурации машины C , которая полагается равной $(q_0^C, 0, \dots, 0)$, в состоянии q_0^A мы достигаем конфигурацию $(q_0^A, c_1^A, \dots, c_n^A)$, а в состоянии q_0^B конфигурацию $(q_0^B, c_1^B, \dots, c_n^B)$. Таким образом, в зависимости от того, какое направление было выбрано, машина C может функционировать далее как машина A или машина B .

Машина A' строится из машины C добавлением нового счётчика x_{n+1} и добавлением для каждого состояния q_i^A машины A перехода $(q_i^A, x_{n+1} := x_{n+1} + 1, q_{A'}^f)$, где $q_{A'}^f$ – финальное состояние машины A' . Машина B' строится из A' добавлением нового перехода для каждого состояния q_i^B машины B : $(q_i^B, x_{n+1} := x_{n+1} + 1, q_{B'}^f)$. Легко видеть, что $R(C) = R(A) \cup R(B) \cup R'(C)$, где $R'(C)$ – множество достижимых векторов значений счётчиков машины C до состояний q_0^A и q_0^B ,

$$R'(C) = \{(c'_1, 0, \dots, 0) \mid c'_1 \leq c_1^A\} \cup \{(c_1^A, c'_2, \dots, 0) \mid c'_2 \leq c_2^A\} \cup \dots \cup \{(c_1^A, c_2^A, \dots, c'_n) \mid c'_n \leq c_n^A\} \cup \{(c'_1, 0, \dots, 0) \mid c'_1 \leq c_1^B\} \cup \{(c_1^B, c'_2, \dots, 0) \mid c'_2 \leq c_2^B\} \cup \dots \cup \{(c_1^B, c_2^B, \dots, c'_n) \mid c'_n \leq c_n^B\}.$$

$$R(A') = R(C) \times \{0\} \cup R(A) \times \{1\}; \quad R(B') = R(A') \cup R(B) \times \{1\} = R(C) \times \{0\} \cup (R(A) \cup R(B)) \times \{1\};$$

Так как из машины C не достижима никакая конфигурация, где $x_{n+1} = 1$, то из этого факта, что $R(B) \subseteq R(A) \iff R(A) = R(A) \cup R(B)$, следует $R(A') = R(B') \iff R(B) \subseteq R(A)$. \square

Несмотря на то, что задача эквивалентности является неразрешимой, всё же можно в некоторых случаях уменьшить число управляющих состояний недетерминированной счётчиковой машины так, чтобы множество достижимости осталось прежним. Для этого необходимо рассмотреть недетерминированную счётчиковую машину как конечный автомат $(Q, T, \rightarrow, q_0, F)$, где Q – конечное множество управляющих состояний, T – множество входных символов, q_0 – начальное состояние, \rightarrow – функция переходов, а $F = Q$ – множество заключительных или допускающих состояний, и провести процедуру минимизации [9].

5. Проблема достижимости

Проблема достижимости состоит в обнаружении алгоритма, с помощью которого для любой недетерминированной счётчиковой машины M и для любого вектора значений счётчиков $\bar{c} = (c_1, \dots, c_m)$ можно выяснить выполнимость условия $\bar{c} \in R(M)$.

Теорема 3. Проблема достижимости для недетерминированной 4-х счётчиковой машины и для вектора значений счётчиков $(c_1, c_2, c_3, c_4) > (0, 0, 0, 0)$ является неразрешимой.

Доказательство. Рассмотрим машину Минского с двумя счётчиками $2cM$. В этом доказательстве предполагается, что если машина $2cM$ останавливается, то её счётчики x_1 и x_2 равны соответственно c_1 и c_2 , где $c_1, c_2 \geq 0$. Для машины $2cM$ построим недетерминированную счётчиковую машину $n4cM$ с четырьмя счётчиками x_1, x_2, x_3, x_4 такую, что $(c_1, c_2, c_3, c_4) \in R(n4cM)$ тогда и только тогда, когда $2cM$ останавливается.

Построим для машины $2cM$ её “слабую” модель $n2cM$ следующим образом. Добавим оператор недетерминированного выбора \sqcap . И в $2cM$ заменим каждое правило перехода $\text{if } x > 0 \text{ then } \text{comm} \text{ else } \text{comm}$ на правило $\text{comm} \sqcap \text{comm}$. В выражении $x_j := x_j - 1$ оператор ‘ $-$ ’ заменим оператором вычитания до нуля ‘ \ominus ’. Таким образом, мы получили недетерминированную машину $n2cM$.

Преобразуем машину $n2cM$ в недетерминированную машину $n4cM$ следующим образом. В каждом правиле переходов $\text{comm} \sqcap \text{comm}$ заменим каждый переход $(q, x_j := x_j \ominus 1, q')$ на последовательность из

переходов $(q, x_3 := x_3 + \min(x_j, 1), q'_1)$ $(q'_1, x_j := x_j \ominus 1, q'_2)$ $(q'_2, x_3 := x_3 \ominus 1, q')$, а переход $(q, x_j := x_j, q'')$ на последовательность $(q, x_j := x_j, q''')$ $(q''', x_4 := x_4 + \min(x_j, 1), q'')$.

Рассмотрим машину $n4cM$ с начальной конфигурацией $s_0 = (q_0, 0, 0, c_3, c_4 - 1)$, где $c_3, c_4 > 0$. Добавим к $n4cM$ новое финальное q'_f состояние и переход $(q_f, x_4 := x_4 + 1, q'_f)$. Если машина $n4cM$ из начальной конфигурации s_0 достигает конфигурацию $(q'_f, c_1, c_2, c_3, c_4)$, то машина Минского $2cM$ останавливается, т.к. в машине $n2cM$ из $(q_0, 0, 0)$ в (q_f, c_1, c_2) существует путь, равный пути машины $2cM$. Действительно, если бы такого пути не было, то “неправильный” переход изменил бы значение одного из счётчиков x_3 или x_4 , т.е. либо уменьшил значение счётчика x_3 , либо увеличил счётчик x_4 . Значения этих счётчиков могут остаться равными соответственно c_3 и c_4 , если $n2cM$ имеет конечный путь в состояние q_f , который будет равен конечному пути в состояние q_f машины Минского $2cM$. Обратно, очевидно, что если $2cM$ из начальной конфигурации $(q_0, 0, 0)$ переходит в финальную (q_0, c_1, c_2) , то машина $n4cM$ из $(q_0, 0, 0, c_3, c_4 - 1)$ достигает $(q'_f, c_1, c_2, c_3, c_4 - 1)$, а затем и $(q'_f, c_1, c_2, c_3, c_4)$. \square

Необходимо отметить, что поскольку, очевидно, недетерминированные счётчиковые машины являются вполне структурированными системами переходов с совместимостью по возрастанию и убыванию, задача достижимости конфигурации $(q, 0, \dots, 0)$ или же просто вектора $(0, \dots, 0)$ разрешима [4].

Литература

- [1] Araki T., Kasami T. Some decision problems related to the reachability problem for Petri nets // Theor. Comp. Sci. 1977. 3(1). P. 85-104.
- [2] Dufourd C., Jancar P., Schnoebelen Ph. Boundedness of Reset P/T nets // Proc. ICALP'99, volume 1644 of Lecture Notes in Computer Science. Springer. 1999. P. 301-310.
- [3] Dufourd C., Finkel A., Schnoebelen Ph. Reset nets between decidability and undecidability // Proc. ICALP'98. volume 1443 of Lecture Notes in Computer Science. Springer. 1998. P. 103-115.
- [4] Finkel A., Schnoebelen Ph. Well-structured transition systems everywhere! // Theoretical Computer Science. 2001. 256(1-2). P. 63-92.
- [5] Kouzmin E., Sokolov V. Communicating Colouring Automata // Proc. International Workshop on Program Understanding, 2003.
- [6] Mayr R. Lossy counter machines // Tech. Report TUM-I9830, Institut für Informatik, TUM, Munich, Germany, October 1998.
- [7] Peterson J. L. Petri Net Theory and the Modeling of Systems. Prentice-Hall Int., 1981.
- [8] Минский М. Вычисления и автоматы. М.: Мир, 1971.
- [9] Хопкрофт Дж., Мотвани Р., Ульман Дж. Введение в теорию автоматов, языков и вычислений. 2-е изд.: Пер. с англ. М.: Вильямс, 2002. 528 с.

УДК 519.22

Об обобщении критерия однородности А.Ю. Левина

Янкевич А.П.

Ярославский государственный университет
150 000, Ярославль, Советская, 14

получена 15 сентября 2003

В статье рассматривается классическая проблема однородности для многомерного случая. Изучается одно из возможных обобщений многомерного непараметрического критерия однородности А.Ю. Левина, связанное с "расширением кругозора" статистики. В определение статистики вводится натуральный параметр k , варьирование которого позволяет усилить устойчивость критерия к небольшим отклонениям от исходной статистической модели (свойство робастности). Приведенное исследование показывает, что обобщенный случай критерия состоятелен против любых альтернатив.

1. Введение. Основной результат

Подтверждение того, что два распределения различны, или же демонстрация того, что они однородны, это задача, которая очень часто возникает в различных областях науки.

Для скалярного случая имеются полностью состоятельные классические непараметрические критерии однородности (Колмогорова-Смирнова, Крамера-Мизеса-Смирнова, Вальда-Волфовича, Рены и др. см., напр., [1, 2, 3]). Существуют многомерные критерии, ориентированные на конкретные типы альтернатив, например, классический критерий χ^2 и многомерные аналоги критерия Вилкоксона, а также критерии, предполагающие нормальность распределений.

В работах А.Ю. Левина [4, 5] был предложен новый мощный многомерный критерий однородности, основанный на отождествлении "быть ближайшим". Там же показаны положительные свойства критерия, такие как полная состоятельность, непараметричность, нечувствительность к размерности, гибкость при получении инвариантных версий, робастность и др.

Целью настоящей работы является изучение одного из возможных обобщений γ -критерия А.Ю. Левина, сохраняющего все положительные свойства оригинальной идеи, но при этом дающего большую устойчивость статистики к грубым ошибкам в статистических данных (свойство робастности).

Итак, в \mathbf{R}^N рассматриваются две независимые выборки

$$S_i = \{x_{i1}, x_{i2}, \dots, x_{in_i}\}, \quad i = 1, 2,$$

с непрерывными плотностями p_1, p_2 соответственно. Проверке подлежит гипотеза однородности

$$H_0 : p_1(x) = p_2(x) \text{ всюду в } \mathbf{R}^N. \quad (1)$$

Обозначим через S объединенную выборку объема $n = n_1 + n_2$. Как обычно предполагается, что выборки равномощны, то есть с ростом n величины $n_1 n^{-1}, n_2 n^{-1}$ стремятся к каким-либо ненулевым пределам:

$$\frac{n_i}{n} \rightarrow s_i (\neq 0), \quad i = 1, 2 \quad (n \rightarrow \infty). \quad (2)$$

Зададим натуральный параметр k . Пусть θ_i ($0 \leq \theta_i \leq n_i$) — число элементов выборки S_i , для которых из остальных $n - 1$ элементов объединенной выборки S среди $2k - 1$ ближайших k или более принадлежат той же выборке S_i ($i = 1, 2$).

Теорема 1. Критерий, отклоняющий гипотезу H_0 при больших значениях статистики

$$\gamma_n^{2k-1} = \sqrt{n} \left(\frac{\theta_1}{n_1} + \frac{\theta_2}{n_2} - 1 \right), \quad (3)$$

состоителен против любых альтернатив.

А именно, имеют место два утверждения:

а) если гипотеза H_0 верна, то для любого $\varepsilon > 0$ можно указать постоянную $C(\varepsilon)$, зависящую лишь от ε , такую что при всех n

$$P(\gamma_n^{2k-1} < C(\varepsilon)) > 1 - \varepsilon; \quad (4)$$

б) если H_0 не верна (то есть плотности p_1 и p_2 различны в некоторой области), то при любом C

$$P(\gamma_n^{2k-1} < C) \rightarrow 0 \quad (n \rightarrow \infty). \quad (5)$$

Критерий, основанный на этих фактах, будем называть обобщенным γ -критерием А.Ю. Левина.

Отличие от классического γ -критерия состоит в том, что для каждого элемента x_{ij} рассматривается не один, а $2k-1$ ближайших элементов из S . Идея направлена на то, чтобы усилить устойчивость критерия к небольшим отклонениям от исходной статистической модели, например, в случаях, когда нарушается строгая независимость результатов близких по очертанию экспериментов.

Заметим, что при $n_1, n_2 \geq k$ событие $A_{l,n}(z)$, состоящее в том, что среди $2k-1$ ближайших к $z \in \mathbb{R}^N$ элементов объединенной выборки S k или более принадлежат выборке S_l , равносильно событию $B_{l,n}(z)$ такому, что

$$\min^{(k)}\{|z - x_{li}| | i = 1, \dots, n_l, j = 1, \dots, n_m, m \neq l\}, \quad (6)$$

где $\min^{(k)}\{x_1, \dots, x_n\} = x_k$, если $x_1 \leq x_2 \leq \dots \leq x_n$.

Таким образом можно говорить, что θ_l ($l = 1, 2$) — это число элементов выборки S_l , для которых расстояние до k -го ближайшего элемента из той же выборки S_l меньше, чем расстояние до k -го ближайшего элемента из другой выборки S_m ($m \neq l$). Этот факт полезно иметь в виду при программной реализации критерия.

Далее отдельно изучается поведение $E\gamma_n^{2k-1}$ и $D\gamma_n^{2k-1}$ при выполнении и невыполнении гипотезы H_0 , а затем на основе полученных оценок доказывается полная состоятельность обобщенного γ -критерия. Структура доказательства в целом повторяет работу А.Ю. Левина [5].

Некоторые части рассуждений без изменений переносятся на обобщенный случай и по сему будут изложены конспективно со ссылками на источник. Моменты же, представляющие принципиальную значимость, будут сопровождены строгими доказательствами.

2. $E\gamma_n^{2k-1}$ в случае однородности

В предположении, что гипотеза H_0 выполнена, $E\gamma_n^{2k-1}$ легко вычисляется. Среди $2k-1$ ближайших к элементу x_{li} выборки S_l могут равновероятно оказаться любые $2k-1$ из $n-1$ элементов объединенной выборки S . Пусть I_{li} — индикатор того события, что среди этих $2k-1$ ближайших k или более принадлежат выборке S_l . Тогда

$$P(I_{li} = 1) = \frac{\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j}}{\binom{n-1}{2k-1}}; \quad P(I_{2i} = 1) = \frac{\sum_{j=k}^{2k-1} \binom{n_2-1}{j} \binom{n_1}{2k-1-j}}{\binom{n-1}{2k-1}}. \quad (7)$$

Поскольку θ_l есть сумма n_l индикаторов I_{li} , то $E\theta_l = n_l P(I_{li} = 1)$, и, следовательно,

$$E\gamma_n^{2k-1} = \sqrt{n} (P(I_{li} = 1) + P(I_{2i} = 1) - 1).$$

Воспользовавшись свойствами сочетаний $\binom{n}{m} = \binom{n}{n-m}$, $\binom{n+1}{m+1} = \binom{n}{m} + \binom{n}{m+1}$, $\sum_{j=0}^k \binom{n}{j} \binom{m}{k-j} = \binom{n+m}{k}$, имеем

$$\begin{aligned} E\gamma_n^{2k-1} &= \frac{\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j}}{\binom{n-1}{2k-1}} + \frac{\sum_{j=0}^{k-1} \binom{n_2-1}{j} \binom{n_1}{2k-1-j}}{\binom{n-1}{2k-1}} - 1 = \\ &= \frac{\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j} + \sum_{j=1}^{k-1} (\binom{n_1-1}{j} + \binom{n_1-1}{j-1}) \binom{n_2-1}{2k-1-j} + \binom{n_2-1}{2k-1} - \binom{n-1}{2k-1}}{\binom{n-1}{2k-1}} = \end{aligned}$$

$$\begin{aligned}
&= \frac{\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j} + \sum_{j=0}^{k-1} \binom{n_1-1}{j} \binom{n_2-1}{2k-1-j} + \sum_{j=1}^{k-1} \binom{n_1-1}{j-1} \binom{n_2-1}{2k-1-j} - \binom{n-1}{2k-1}}{\binom{n-1}{2k-1}} = \\
&= \frac{\sum_{j=0}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j} - \sum_{j=0}^{k-1} \binom{n_1-1}{j} \binom{n_2-1}{2k-2-j} + \sum_{j=1}^{k-1} \binom{n_1-1}{j-1} \binom{n_2-1}{2k-1-j} - \binom{n-1}{2k-1}}{\binom{n-1}{2k-1}} = \\
&= \frac{-\sum_{j=0}^{k-1} \binom{n_1-1}{j} \binom{n_2-1}{2k-2-j} + \sum_{i=0}^{k-2} \binom{n_1-1}{i} \binom{n_2-1}{2k-2-i}}{\binom{n-1}{2k-1}}.
\end{aligned}$$

Таким образом

$$\mathbf{E}\gamma_n^{2k-1} = -\sqrt{n} \frac{\binom{n_1-1}{k-1} \binom{n_2-1}{k-1}}{\binom{n-1}{2k-1}}. \quad (8)$$

Итак, величина $\mathbf{E}\gamma_n^{2k-1}$ отрицательна и стремится к 0 при $n \rightarrow \infty$.

3. Асимптотика $\mathbf{E}\gamma_n^{2k-1}$ в общем случае

В этом параграфе будет показано, что в случае неоднородности предел $\mathbf{E}\gamma_n^{2k-1}/\sqrt{n}$ ($n \rightarrow \infty$), в отличие от однородного случая, положителен.

Через \mathbf{R}_+^N здесь и далее обозначается множество таких точек для функций p_1, p_2 , в которых по крайней мере одна из них отлична от нуля. Пусть $A_{l,n}(z) = A_l(z, n)$ — событие, состоящее в том, что среди $2k-1$ ближайших к z ($\in \mathbf{R}_+^N$) точек из S k или более принадлежат выборке S_l ($l = 1, 2$).

Теорема 2. Для любой точки $z \in \mathbf{R}_+^N$

$$\mathbf{P}(A_{1,n}(z)) \rightarrow g_1(z) = \frac{\sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1(z))^j (s_2 p_2(z))^{2k-1-j}}{(s_1 p_1(z) + s_2 p_2(z))^{2k-1}}, \quad (n \rightarrow \infty); \quad (9)$$

$$\mathbf{P}(A_{2,n}(z)) \rightarrow g_2(z) = \frac{\sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1(z))^j (s_2 p_2(z))^{2k-1-j}}{(s_1 p_1(z) + s_2 p_2(z))^{2k-1}}, \quad (n \rightarrow \infty). \quad (10)$$

Для доказательства теоремы 2 рассмотрим сначала частный случай.

Лемма 1. Пусть p_1 и p_2 постоянны в некоторой окрестности $U = U(z)$ точки z . Тогда имеют место соотношения (9), (10).

Пусть $A_n(z)$ — событие, состоящее в том, что в U попало α_n (α_{in}) элементов выборки S (S_i) ($i = 1, 2$). Ясно, что $\alpha_n = \alpha_{1n} + \alpha_{2n}$.

Пусть $\alpha_n > 2k-1$. Все α_n элементов, которые попадут в U , имеют одинаковое условное распределение, а именно, равномерное в U , в силу условия леммы. Так как они, очевидно, независимы, то любой из них равновероятно с остальными может оказаться в числе $2k-1$ ближайших к z . Поэтому

$$\mathbf{P}(A_{1,n}(z)|A_n) = \sum_{j=k}^{2k-1} \frac{\binom{\alpha_{1n}}{j} \binom{\alpha_{2n}}{2k-1-j}}{\binom{\alpha_n}{2k-1}}, \quad (\alpha_n > 2k-1); \quad (11)$$

$$\mathbf{P}(A_{2,n}(z)|A_n) = \sum_{j=0}^{k-1} \frac{\binom{\alpha_{1n}}{j} \binom{\alpha_{2n}}{2k-1-j}}{\binom{\alpha_n}{2k-1}}, \quad (\alpha_n > 2k-1); \quad (12)$$

при $\alpha_n \leq 2k-1$ правая часть заменяется величиной $\mathbf{P}(A_{l,n}(z)|\alpha_n \leq 2k-1)$ ($l = 1, 2$), несущественной для дальнейшего.

Пусть $|U|$ — объем U . В силу (2) и усиленного закона больших чисел при любых i, l ($1 \leq i, l \leq 2$)

$$\frac{\alpha_{in}}{n_i} \xrightarrow{n \rightarrow \infty} |U| p_i(z), \quad (n \rightarrow \infty),$$

$$\begin{aligned} \frac{\alpha_{ln}}{\alpha_n} &= \frac{n_l}{n} \frac{\alpha_{ln}}{n_l} \left(\frac{\alpha_n}{n} \right)^{-1} = \frac{n_l}{n} \frac{\alpha_{ln}}{n_l} \left(\frac{n_1}{n} \frac{\alpha_{1n}}{n_1} + \frac{n_2}{n} \frac{\alpha_{2n}}{n_2} \right)^{-1} \xrightarrow{n \rightarrow \infty} \frac{s_l |U| p_l(z)}{s_1 |U| p_1(z) + s_2 |U| p_2(z)} = \\ &= \frac{s_l p_l(z)}{s_1 p_1(z) + s_2 p_2(z)} \quad (n \rightarrow \infty), \quad \sum_{j=k}^{2k-1} \frac{\binom{\alpha_{1n}}{j} \binom{\alpha_{2n}}{2k-1-j}}{\binom{\alpha_n}{2k-1}} \xrightarrow{j \rightarrow \infty} \sum_{j=k}^{2k-1} \frac{\alpha_{1n}^j}{j!} \frac{\alpha_{2n}^{2k-1-j}}{(2k-1-j)!} / \frac{\alpha_n^{2k-1}}{(2k-1)!} = \\ &= \sum_{j=k}^{2k-1} \binom{2k-1}{j} \left(\frac{\alpha_{1n}}{\alpha_n} \right)^j \left(\frac{\alpha_{2n}}{\alpha_n} \right)^{2k-1-j} \xrightarrow{n \rightarrow \infty} g_1(z) \quad (n \rightarrow \infty). \end{aligned}$$

Аналогично для второй выборки.

Поскольку $z \in \mathbf{R}_+^N$, то $P(\alpha_n \leq 2k-1) \rightarrow 0$ при $n \rightarrow \infty$, так что случаем $\alpha_n \leq 2k-1$ можно пренебречь. Итак,

$$P(A_{l,n}(z)|A_n) \xrightarrow{n \rightarrow \infty} g_l(z) \quad (n \rightarrow \infty), \quad l = 1, 2.$$

Случайные величины в левой части равномерно ограничены, поэтому тот же предел должны иметь их математические ожидания:

$$E P(A_{l,n}(z)) = E P(A_{l,n}(z)|A_n) \rightarrow g_l(z) \quad (n \rightarrow \infty), \quad l = 1, 2.$$

Лемма I доказана.

Ход дальнейшего доказательства заключается в том, что произвольная плотность $p_l(x)$ заключается в некоторой окрестности точки z между двумя вспомогательными плотностями, постоянными в этой окрестности. Ко вспомогательным плотностям, применима лемма 1. Вероятность события $A_{l,n}(z)$ также оказывается заключенной между вероятностями соответствующих событий для вспомогательных плотностей. В итоге, устремив мажорирующую и минорирующую плотность друг к другу, из того, что предельные значения $g_l(z)$ для них совпадают и равны значениям, указанным в (9) или (10) (в зависимости от l), следует, что к тому же пределу стремится и $P(A_{l,n}(z))$.

Эти рассуждения, подробно изложенные в [5], переносятся на обобщенный случай с одним-единственным замечанием.

В доказательстве А.Ю. Левина использован тот факт, что для двух наборов независимых одинаково распределенных случайных величин $\xi_{11}, \dots, \xi_{1n}, \xi_{21}, \dots, \xi_{2n}$ с непрерывными функциями распределения $F_1(t), F_2(t)$ соответственно из того, что $F_1(t) \leq F_2(t)$ ($-\infty < t < \infty$), следует, что для функций $\Phi_{l,n}(t)$ распределения минимума из $\xi_{l1}, \dots, \xi_{ln}$ ($l = 1, 2$) также всюду $\Phi_{1,n}(t) \leq \Phi_{2,n}(t)$.

Для функций $\Phi_{l,n}^{(k)}(t) = \Phi_n^{(k)}(F_l, t)$ распределения k -го минимума то же свойство получается из следующих соображений. Поскольку

$$\begin{aligned} \Phi_n^{(k)}(F, t) &= 1 - \sum_{j=0}^{k-1} \binom{n}{j} F(t)^j (1-F(t))^{n-j}, \\ \frac{\partial \Phi_n^{(k)}(F, t)}{\partial F} &= - \sum_{j=0}^{k-1} \binom{n}{j} (jF(t)^{j-1}(1-F(t))^{n-j} - (n-j)F(t)^j(1-F(t))^{n-j-1}) = \\ &= \sum_{j=0}^{k-1} \binom{n}{j} (n-j)F(t)^j(1-F(t))^{n-j-1} - \sum_{i=0}^{k-2} \binom{n}{i+1} (i+1)F(t)^i(1-F(t))^{n-i-1} = \\ &= \binom{n}{k} (n-k+1)F(t)^{k-1}(1-F(t))^{n-k} \geq 0. \end{aligned}$$

Таким образом функция $\Phi_n^{(k)}(F, t)$ монотонна по F , и поэтому из того, что всюду $F_1(t) \leq F_2(t)$, следует, что $\Phi_n^{(k)}(F_1, t) \leq \Phi_n^{(k)}(F_2, t)$ ($-\infty < t < \infty$).

Сформулируем ряд следствий из теоремы 2.

Пусть $B_{l,jn}$ — событие, состоящее в том, что в числе $2k-1$ ближайших к x_{lj} среди $n-1$ остальных элементов S k или более окажутся из выборки S_l .

Следствие 1. При любых t, j и любом $x \in \mathbf{R}_+^N$

$$P(B_{l,jn}|x_{lj} = x) \rightarrow g_l(x) \quad (n \rightarrow \infty). \quad (13)$$

Единственное различие по сравнению с теоремой 2 заключается в том, что здесь из выборки S_l исключается элемент x_{lj} , и ее объем должен быть теперь заменен на $n_l - 1$. Однако эта деталь не влияет на асимптотику левой части (13) при $n \rightarrow \infty$.

Далее, откажемся от фиксации элемента x_{lj} и вернемся к рассмотрению его как случайной величины.

Следствие 2. При любых l, j

$$\mathbf{P}(B_{ljn}) \rightarrow \int_{\mathbf{R}_+^N} g_l(x)p_l(x)dx \quad (n \rightarrow \infty). \quad (14)$$

По формуле полной вероятности

$$\mathbf{P}(B_{ljn}) = \int_{\mathbf{R}_+^N} \mathbf{P}(B_{ljn} | x_{lj} = x)p_l(x)dx, \quad (j = 1, \dots, n_l). \quad (15)$$

Интегрирование по \mathbf{R}_+^N вместо \mathbf{R}^N оправдано, поскольку оно исключает лишь интегрирование по множеству, где $p_l(x) = 0$ ($l = 1, 2$). Подынтегральные функции в (15) мажорируются интегрируемой функцией $p_l(x)$, поэтому возможен предельный переход под знаком интеграла, и (15) с учетом (13) дает (14).

Вернемся к статистике γ_n^{2k-1} . Пусть I_{lm} — индикатор события $B_{lm} = B_{lmn}$. Тогда при $n \rightarrow \infty$

$$\begin{aligned} \mathbf{E}\theta_l &= \mathbf{E}\left(\sum_{m=1}^{n_l} I_{lm}\right) = \sum_{m=1}^{n_l} \mathbf{P}(B_{lm}) = n_l \mathbf{P}(B_{l1}) \rightarrow n_l \int_{\mathbf{R}_+^N} g_l(x)p_l(x)dx, \\ \mathbf{E}\left(\frac{\gamma_n^{2k-1}}{\sqrt{n}}\right) &= \frac{\mathbf{E}\theta_1}{n_1} + \frac{\mathbf{E}\theta_2}{n_2} - 1 \rightarrow \int_{\mathbf{R}_+^N} \left[p_1(x) \sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1(x))^j (s_2 p_2(x))^{2k-1-j} + \right. \\ &\quad \left. p_2(x) \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1(x))^j (s_2 p_2(x))^{2k-1-j} \right] \frac{dx}{(s_1 p_1(x) + s_2 p_2(x))^{2k-1}} - 1. \end{aligned} \quad (16)$$

Следствие 3. При невыполнении гипотезы H_0

$$\lim_{n \rightarrow \infty} \mathbf{E}\left(\frac{\gamma_n^{2k-1}}{\sqrt{n}}\right) > 0. \quad (17)$$

Рассмотрим отдельно числитель подынтегральной функции в (16). Поскольку $s_1 + s_2 = 1$, имеем

$$\begin{aligned} &(s_1 + s_2)p_1 \sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} + (s_1 + s_2)p_2 \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} = \\ &= (s_1 p_1 + s_2 p_2)^{2k} + \sum_{j=0}^{k-1} \left(\binom{2k-1}{j} - \binom{2k}{j} \right) (s_1 p_1)^j (s_2 p_2)^{2k-j} + \sum_{i=k+1}^{2k} \left(\binom{2k-1}{i-1} - \binom{2k}{i} \right) (s_1 p_1)^i (s_2 p_2)^{2k-i} - \\ &- \binom{2k}{k} (s_1 p_1)^k (s_2 p_2)^k + s_2 p_1 \sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} + s_1 p_2 \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} = \\ &= (s_1 p_1 + s_2 p_2)^{2k} - \sum_{j=1}^{k-1} \binom{2k-1}{j-1} (s_1 p_1)^j (s_2 p_2)^{2k-j} - \\ &- \sum_{j=k+1}^{2k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-j} - \left(\binom{2k-1}{k} + \binom{2k-1}{k-1} \right) (s_1 p_1)^k (s_2 p_2)^k + \\ &+ s_2 p_1 \sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} + s_1 p_2 \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} = \end{aligned}$$

$$\begin{aligned}
&= (s_1 p_1 + s_2 p_2)^{2k} + (p_1 - p_2) \left[s_2 \sum_{j=k}^{2k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} - s_1 \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} \right] = \\
&= (s_1 p_1 + s_2 p_2)^{2k} + (p_1 - p_2) \left[s_2 (s_1 p_1 + s_2 p_2)^{2k-1} - \sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1)^j (s_2 p_2)^{2k-1-j} \right]. \quad (18)
\end{aligned}$$

Вернемся к выражению для математического ожидания. Из (16), (18) имеем

$$\begin{aligned}
&\lim_{n \rightarrow \infty} \mathbf{E} \left(\frac{\gamma_n^{2k-1}}{\sqrt{n}} \right) = \int_{\mathbf{R}_+^N} (s_1 p_1(x) + s_2 p_2(x)) dx - 1 + \\
&+ \int_{\mathbf{R}_+^N} (p_1(x) - p_2(x)) dx + \int_{\mathbf{R}_+^N} (p_2(x) - p_1(x)) F(x) dx = \int_{\mathbf{R}_+^N} (p_2(x) - p_1(x)) F(x) dx, \quad (19)
\end{aligned}$$

где

$$F(x) = \frac{\sum_{j=0}^{k-1} \binom{2k-1}{j} (s_1 p_1(x))^j (s_2 p_2(x))^{2k-1-j}}{(s_1 p_1(x) + s_2 p_2(x))^{2k-1}}. \quad (20)$$

Далее, нам понадобится следующее утверждение.

Лемма 2. Для любых непрерывных плотностей $p_1(x), p_2(x)$ и любой строго монотонно убывающей положительной функции $G(t)$ при $t \geq 0$

$$I = \int_{\mathbf{R}^N} (p_2(x) - p_1(x)) G \left(\frac{p_1(x)}{p_2(x)} \right) dx \geq 0,$$

причем равенство возможно только тогда, когда $p_1(x) \equiv p_2(x)$.

Случай $p_1(x) \equiv p_2(x)$ тривиален. Предположим обратное. Разобьем пространство \mathbf{R}^N на два множества: $R_1 = \{x : p_1(x) \leq p_2(x)\}$, $R_2 = \{x : p_1(x) > p_2(x)\}$. Очевидно, что существует $\varepsilon > 0$, такое, что на некотором невырожденном интервале $X \subset R_2$ $p_1(x)/p_2(x) \geq 1 + \varepsilon$. Тогда

$$\begin{aligned}
I &= \int_{R_1} (p_2(x) - p_1(x)) G \left(\frac{p_1(x)}{p_2(x)} \right) dx - \int_{R_2 \setminus X} (p_1(x) - p_2(x)) G \left(\frac{p_1(x)}{p_2(x)} \right) dx - \\
&- \int_X (p_1(x) - p_2(x)) G \left(\frac{p_1(x)}{p_2(x)} \right) dx \geq \int_{R_1} (p_2(x) - p_1(x)) \inf_{R_1} G \left(\frac{p_1(x)}{p_2(x)} \right) dx - \\
&- \int_{R_2 \setminus X} (p_1(x) - p_2(x)) \sup_{R_2 \setminus X} G \left(\frac{p_1(x)}{p_2(x)} \right) dx - \int_X (p_1(x) - p_2(x)) \sup_X G \left(\frac{p_1(x)}{p_2(x)} \right) dx = \\
&= G(1) \int_{R_1} (p_2(x) - p_1(x)) dx - G(1) \int_{R_2 \setminus X} (p_1(x) - p_2(x)) dx - G(1 + \varepsilon) \int_X (p_1(x) - p_2(x)) dx > \\
&> G(1) \int_{R_1} (p_2(x) - p_1(x)) dx - G(1) \int_{R_2} (p_1(x) - p_2(x)) dx = 0.
\end{aligned}$$

Лемма доказана.

Покажем, что функция $F(t)$, заданная в (20), строго монотонно убывает.

$$\begin{aligned}
F(t) &= \frac{\sum_{j=0}^{k-1} \binom{2k-1}{j} t^j}{(1+t)^{2k-1}}, \quad F'(t) = \frac{\sum_{j=1}^{k-1} \binom{2k-1}{j} j t^{j-1} (1+t) - \sum_{j=0}^{k-1} \binom{2k-1}{j} t^j (2k-1)}{(1+t)^{2k}} = \\
&= \frac{\sum_{i=0}^{k-2} \binom{2k-1}{i+1} (i+1) t^i - \sum_{j=0}^{k-1} \binom{2k-1}{j} (2k-1-j) t^j}{(1+t)^{2k}} = - \frac{\binom{2k-1}{k-1} k t^{k-1}}{(1+t)^{2k}}.
\end{aligned}$$

Применение к (19), (20) леммы 2 завершает доказательство следствия 3.

4. Оценка $D\gamma_n^{2k-1}$ при гипотезе H_0

Теорема 3. Если справедлива гипотеза H_0 , то при любом $n (\geq 4k)$ существует константа C , не зависящая от n , такая что

$$D\gamma_n^{2k-1} \leq C \quad (n \geq 4k). \quad (21)$$

Из определения γ_n^{2k-1} непосредственно следует, что

$$D\gamma_n^{2k-1} = \sum_{l=1}^2 \frac{1}{n_l^2} D\theta_l + \sum_{r \neq s} \frac{1}{n_r n_s} cov(\theta_r, \theta_s). \quad (22)$$

Итак, пусть выполнена гипотеза H_0 , т.е. все n элементов различных выборок S_i имеют одну и ту же плотность $p(x)$. Как и в работе [5], введем ряд полезных обозначений. Избавимся от двойной индексации элементов x_{ij} и обозначим их просто как x_1, x_2, \dots, x_n . При этом множество индексов $H_1 = \{1, \dots, n_1\}$ отвечает элементам выборки S_1 , а множество $H_2 = \{n_1 + 1, \dots, n\}$ — элементам выборки S_2 . Для элемента x_i через $I_{i i_1 \dots i_{2k-1}} = I_{i A_i}$ обозначим индикатор того события, что среди остальных $n - 1$ элементов S в числе $2k - 1$ ближайших к x_i окажутся элементы $x_{i_1}, \dots, x_{i_{2k-1}}$. A_i , таким образом, — множество индексов $2k - 1$ ближайших к x_i элементов S . (Отметим, что эти индикаторы, вообще говоря, зависимы).

Очевидно, если $i \in A_i$ или какие-то $i_r = i_s$ ($i_r, i_s \in A_i$), то $I_{i A_i} = 0$; и $E I_{i A_i} = \binom{n-1}{2k-1}^{-1}$ ($0 \leq i, i_j \leq n$) в иных случаях.

Поскольку все x_i одинаково распределены и независимы, их совместное распределение перестановочно (т.е. симметрично). Поэтому положим

$$\begin{aligned} a_m &= E(I_{iA_i} I_{jA_j}), \quad i \notin A_j, \quad j \notin A_i, \quad |A_i \cap A_j| = m, \quad (0 \leq m \leq 2k-1); \\ b_m &= E(I_{iA_i} I_{jA_j}), \quad i \notin A_j, \quad j \in A_i, \quad |A_i \cap A_j| = m, \quad (0 \leq m \leq 2k-2); \\ c_m &= E(I_{iA_i} I_{jA_j}), \quad i \in A_j, \quad j \in A_i, \quad |A_i \cap A_j| = m, \quad (0 \leq m \leq 2k-2). \end{aligned} \quad (23)$$

Величины (23) в общем случае зависят от плотности p . Однако при этом они удовлетворяют определенному соотношению.

В самом деле, индикаторы I_{1A_1} , где множество индексов A_1 пробегает по всем сочетаниям из оставшихся $n - 1$ по $2k - 1$ элементов, отвечают событиям, образующим разбиение.

При $k = 1$ (так как $E(I_{2\{3\}} I_{1\{2\}}) = b_0$) имеем следующие равенства:

$$(n-1)^{-1} = E I_{2\{1\}} = E(I_{2\{1\}} I_{1\{2\}}) + \dots + E(I_{2\{1\}} I_{1\{n\}}) = c_0 + (n-2)b_0,$$

$$(n-1)^{-1} = E I_{2\{3\}} = E(I_{2\{3\}} I_{1\{2\}}) + E(I_{2\{3\}} I_{1\{3\}}) + \dots + E(I_{2\{3\}} I_{1\{n\}}) = b_0 + a_1 + (n-3)a_0.$$

При $k = 2$, аналогично

$$\frac{1}{\binom{n-1}{3}} = c_2 + 2(n-4)c_1 + \binom{n-4}{2}c_0 + (n-4)b_2 + 2\binom{n-4}{2}b_1 + \binom{n-4}{3}b_0,$$

$$\frac{1}{\binom{n-1}{3}} = 3b_2 + 3(n-5)b_1 + \binom{n-5}{2}b_0 + a_3 + 3(n-5)a_2 + 3\binom{n-5}{2}a_1 + \binom{n-5}{3}a_0.$$

При произвольных k эти соотношения выглядят следующим образом:

$$\frac{1}{\binom{n-1}{2k-1}} = \sum_{j=1}^{2k-2} (2k-1-j) \binom{n-2k}{2k-2-j} c_j + \binom{n-2k}{2k-2} c_0 + \sum_{j=1}^{2k-2} (2k-1-j) \binom{n-2k}{2k-1-j} b_j + \binom{n-2k}{2k-1} b_0,$$

$$\frac{1}{\binom{n-1}{2k-1}} = \sum_{j=0}^{2k-2} \binom{2k-1}{j} \binom{n-2k-1}{2k-2-j} b_j + \sum_{j=0}^{2k-1} \binom{2k-1}{j} \binom{n-2k-1}{2k-1-j} a_j.$$

Отсюда следует, что при $n \rightarrow \infty$

$$a_j \leq \frac{1}{\binom{n-1}{2k-1} \binom{2k-1}{j} \binom{n-2k-1}{2k-1-j}} = O(n^{-(4k-2-j)}), \quad (j = 0, \dots, 2k-1);$$

$$\begin{aligned}
b_0 &\leq \frac{1}{\binom{n-1}{2k-1} \binom{n-2k}{2k-1}} = O(n^{-(4k-2)}); \\
b_j &\leq \frac{1}{(2k-1-j) \binom{n-1}{2k-1} \binom{n-2k}{2k-1-j}} = O(n^{-(4k-2-j)}), \quad (j=1, \dots, 2k-2); \\
c_0 &\leq \frac{1}{\binom{n-1}{2k-1} \binom{n-2k}{2k-2}} = O(n^{-(4k-3)}); \\
c_j &\leq \frac{1}{(2k-1-j) \binom{n-1}{2k-1} \binom{n-2k}{2k-2-j}} = O(n^{-(4k-3-j)}), \quad (j=1, \dots, 2k-2).
\end{aligned}$$

Как уже отмечалось в параграфе 2, при любых k

$$\theta_1 = \sum I_{i\{i_j\}}, \quad \mathbf{E}\theta_1 = \sum \mathbf{E}I_{i\{i_j\}} = n_1 \frac{\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j}}{\binom{n-1}{2k-1}}, \quad (24)$$

где суммирование производится по всем $i \in H_1$, $A_i : |A_i \cap H_1| > |A_i \cap H_2|$,

$$\theta_2 = \sum I_{i\{i_j\}}, \quad \mathbf{E}\theta_2 = \sum \mathbf{E}I_{i\{i_j\}} = n_2 \frac{\sum_{j=0}^{k-1} \binom{n_1}{j} \binom{n_2-1}{2k-1-j}}{\binom{n-1}{2k-1}}, \quad (25)$$

суммирование производится по всем $i \in H_2$, $A_i : |A_i \cap H_1| < |A_i \cap H_2|$.

Переходим к дисперсиям и ковариации величин θ_1 . Для $\mathbf{D}\theta_1$, посчитав число ненулевых слагаемых $\mathbf{E}(I_{iA}, I_{jA})$ типов a_m, b_m, c_m , получаем

$$\begin{aligned}
\mathbf{D}\theta_1 &= \mathbf{E}\theta_1 + n_1(n_1-1) \left[a_0 \sum_{r=k}^{2k-1} \binom{n_1-2}{r} \binom{n_2}{2k-1-r} \sum_{s=k}^{2k-1} \binom{n_1-2-r}{s} \binom{n_2-2k+1+r}{2k-1-s} + \right. \\
&\quad + 2b_0 \sum_{r=k-1}^{2k-2} \binom{n_1-2}{r} \binom{n_2}{2k-2-r} \sum_{s=k}^{2k-1} \binom{n_1-2-r}{s} \binom{n_2-2k+2+r}{2k-1-s} + \\
&\quad \left. + c_0 \sum_{r=k-1}^{2k-2} \binom{n_1-2}{r} \binom{n_2}{2k-2-r} \sum_{s=k-1}^{2k-2} \binom{n_1-2-r}{s} \binom{n_2-2k+2+r}{2k-2-s} + \dots \right] - (\mathbf{E}\theta_1)^2.
\end{aligned} \quad (26)$$

С учетом асимптотических оценок величин $a_m, b_m, c_m, \mathbf{E}\theta_1$ получаем

$$\begin{aligned}
\mathbf{D}\theta_1 &= n_1(n_1-1)a_0 \sum_{r=k}^{2k-1} \binom{n_1-2}{r} \binom{n_2}{2k-1-r} \sum_{s=k}^{2k-1} \binom{n_1-2-r}{s} \binom{n_2-2k+1+r}{2k-1-s} + O(n) - \\
&\quad - n_1^2 \frac{\left(\sum_{j=k}^{2k-1} \binom{n_1-1}{j} \binom{n_2}{2k-1-j} \right)^2}{\binom{n-1}{2k-1}^2} \leq n_1^2 \frac{n^{4k-2} \sum_{r=k}^{2k-1} \frac{s_1^r}{r!} \frac{s_2^{2k-1-r}}{(2k-1-r)!} \sum_{s=k}^{2k-1} \frac{s_1^s}{s!} \frac{s_2^{2k-1-s}}{(2k-1-s)!} + O(n^{4k-3})}{\frac{n^{4k-2}}{(2k-1)!^2}} + O(n) - \\
&\quad - n_1^2 \frac{n^{4k-2} \left(\sum_{j=k}^{2k-1} \frac{s_1^j}{j!} \frac{s_2^{2k-1-j}}{(2k-1-j)!} \right)^2}{\frac{n^{4k-2}}{(2k-1)!^2}} + O(n^{4k-3}) = O(n).
\end{aligned}$$

Аналогично получаем, что $\mathbf{D}\theta_2 = O(n)$ и $\text{cov}(\theta_1, \theta_2) = O(n)$. Следовательно,

$$\mathbf{D}\gamma_n^{2k-1} = n \left[\frac{1}{n_1^2} \mathbf{D}\theta_1 + \frac{1}{n_2^2} \mathbf{D}\theta_2 + \frac{2}{n_1 n_2} \text{cov}(\theta_1, \theta_2) \right] = O(1).$$

Поэтому, очевидно, существует константа $C > 0$ такая, что при любом n $\mathbf{D}\gamma_n^{2k-1} < C$. Утверждение теоремы 3 доказано.

Для получения конкретной численной оценки дисперсии при $n \rightarrow \infty$ в (26) необходимо детально рассмотреть все слагаемые порядка $O(n)$. В свете изложенного задача эта не представляет особой сложности, однако приводит к громоздким выкладкам. По сему оставим её рассмотрение за рамками настоящей статьи, ограничившись лишь описанием результатов.

Удаётся показать, что при произвольных значениях параметра k , как и в случае $k = 1$, справедлива асимптотическая оценка

$$\mathbf{D}\gamma_n^{2k-1} = C_0 + O(n^{-1}), \quad C_0 \leq 2 \quad (n \rightarrow \infty).$$

5. Асимптотическая оценка сверху $\mathbf{D}\gamma_n^{2k-1}$ в общем случае

При рассмотрении $\mathbf{D}\gamma_n^{2k-1}$ без предположения об однородности достаточно ограничиться следующим утверждением.

Теорема 4. *Имеет место оценка*

$$\mathbf{D}\gamma_n^{2k-1} = o(n) \quad (n \rightarrow \infty). \quad (27)$$

Пусть снова I_{li} — индикатор того события, среди $2k-1$ ближайших к x_{li} элементов выборки S большая их часть также принадлежит S_l . В этих обозначениях

$$\gamma_n^{2k-1} = \sqrt{n} \left(\sum_{l=1}^2 \sum_{i=1}^{n_l} \frac{1}{n_l} I_{li} - 1 \right).$$

Утверждение теоремы 4 будет доказано, если мы проверим, что

$$\begin{aligned} \frac{1}{n} \mathbf{D}\gamma_n^{2k-1} &= \sum_{l=1}^2 \sum_{i=1}^{n_l} \frac{1}{(n_l)^2} \mathbf{D}I_{li} + \sum_{l=1}^2 \sum_{i,j=1, i \neq j}^{n_l} \frac{1}{(n_l)^2} \text{cov}(I_{li}, I_{lj}) + 2 \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \frac{1}{n_1 n_2} \text{cov}(I_{1i}, I_{2j}) = \\ &= \sum_1 + \sum_2 + \sum_3 = o(1) \quad (n \rightarrow \infty). \end{aligned}$$

Поскольку $\sum_1 = O(n^{-1})$ ($n \rightarrow \infty$), дальнейший интерес представляют лишь слагаемые \sum_2, \sum_3 .

Пусть $c_{qr} = \text{cov}(I_{qi}, I_{rj})$ при всех i, j таких, что $i = 1, \dots, n_q$, $j = 1, \dots, n_r$ и $i \neq j$, если $q = r$. Случаю $q = r$ в сумме \sum_2 отвечает $n_r(n_r - 1)$ одинаковых слагаемых, равных $n_r^{-2} c_{rr}$, а случаю $q \neq r$ в сумме \sum_3 — $2n_q n_r$ одинаковых слагаемых, равных $n_q^{-1} n_r^{-1} c_{qr}$. Поэтому

$$\{\sum_2 + \sum_3\} \leq \sum_{q,r=1}^2 |c_{qr}|.$$

При любом n в правую часть в качестве слагаемых входит фиксированное число (а именно, 4) ковариаций. Поэтому бесконечная малость \sum_2 и \sum_3 при $n \rightarrow \infty$ вытекает из того факта, что при всех q, r $c_{qr} \rightarrow 0$ ($n \rightarrow \infty$; $1 \leq q, r \leq 2$). Суть этого утверждения заключается в том, что с ростом n зависимость между различными индикаторами сходит на нет. Подробное доказательство изложено в работе [5]. Оно без изменений применимо и к обобщенному случаю γ -критерия. По сему оставим эту часть рассуждений за рамками настоящей статьи и перейдем к доказательству основного результата.

6. Доказательство теоремы 1

В случае однородности, согласно результатам параграфов 2 и 4 (теорема 3), имеем

$$\mathbf{E}\gamma_n^{2k-1} = o(1) \quad (< 0), \quad \mathbf{D}\gamma_n^{2k-1} \leq C.$$

Утверждение а) следует из неравенства Чебышева, гласящего, что для любой случайной величины X с конечными мат. ожиданием и дисперсией

$$\mathbf{P}(X - \mathbf{E}X \geq \varepsilon) \leq \frac{\mathbf{D}X}{\varepsilon^2}.$$

Поэтому, с учетом того, что $E\gamma_n^{2k-1} < 0$, в иных обозначениях имеем

$$P(\gamma_n^{2k-1} < C(\varepsilon)) > 1 - P(\gamma_n^{2k-1} - E\gamma_n^{2k-1} \geq C(\varepsilon)) \geq 1 - \frac{D\gamma_n^{2k-1}}{C(\varepsilon)^2} \geq 1 - \frac{C}{C(\varepsilon)^2} \geq 1 - \varepsilon.$$

Отсюда в качестве $C(\varepsilon)$ можно выбрать величины $\sqrt{C\varepsilon^{-1}}$ (они могут быть уточнены).

В случае неоднородности, согласно (17), (27),

$$E\left(\frac{\gamma_n^{2k-1}}{\sqrt{n}}\right) \rightarrow \alpha, \quad D\left(\frac{\gamma_n^{2k-1}}{\sqrt{n}}\right) \rightarrow 0 \quad (n \rightarrow \infty),$$

где α — некоторая положительная константа. Итак, γ_n^{2k-1} п.н. сходится к $\alpha (> 0)$. Отсюда сразу следует утверждение б).

Теорема доказана.

7. Свойства обобщенного γ -критерия

Результат "расширения кругозора" статистики будет виден в тех случаях, когда необходимые предположения (о независимости и одинаковой распределенности элементов внутри выборок) несколько нарушаются. Например, если имеется зависимость между результатами близких по очереди экспериментов, что часто возникает, когда изучаемые данные являются результатом дискретизации некоторого непрерывного процесса. В этом случае элементы внутри выборок S_1, S_2 будут либо "тяготеть", либо "отталкиваться" друг от друга. При $k = 1$ это приведет к неоправданному росту, либо уменьшению значения статистики. Рассмотрение же более одного ближайшего элемента позволяет уменьшить влияние этой ошибки и получить более адекватные результаты.

На практике оказывается полезным наблюдать изменения значений статистики с ростом параметра k . Эти изменения должны оказываться незначительными. В противном случае следует отнести к полученным результатам очень настороженно, и прежде чем принимать решение об отклонении или принятии гипотезы H_0 , стоит постараться найти причины наблюдаемой тенденции.

В заключение отметим, что проводилось множество численных экспериментов, которые свидетельствуют о практической эффективности γ -критерия. Так, в работе [6] его с успехом применили в исследованиях активности мозга.

Литература

- [1] Кендалл М.Дж., Стьюарт А. Многомерный статистический анализ и временные ряды. М.:Наука, 1976. 736 с.
- [2] Леман Э. Проверка статистических гипотез. М.:Наука, 1979. 408 с.
- [3] Боровков А.А. Математическая статистика. Дополнительные главы. М.:Наука, 1984. 144 с.
- [4] Левин А.Ю. О состоятельном многомерном непараметрическом критерии однородности // Усп. матем. наук. 1993. Т. 48, № 6. С. 155–156.
- [5] Левин А.Ю. О состоятельном многомерном непараметрическом критерии однородности // Модел. и анализ информ. систем. 2002. Т. 9, № 2. С. 32–44.
- [6] Майоров В.В., Мышкин И.Ю. Критерий количественной оценки различий вызванных потенциалов мозга // Депонир.. ИИОН РАН, № 50504, 15.6, 1995. 21 с.