## Московский Государственный Университет

## Факультет Вычислительной Математики и Кибернетики Кафедра Математических Методов Прогнозирования

### Курсовая Работа

Поиск скрытых поведенческих паттернов

Вишневский Валерий группа 317

научный руководитель: Ветров Дмитрий Петрович

#### Аннотация

В работе рассматривается алгоритм поиска скрытых поведенческих закономерностей (hidden T-Patterns), которые сложно обнаружить визуально, или с помощью стандартных статистических методов. Проводится анализ параметров алгоритма. Предложенный метод поиска паттернов основывается на определении взаимосвязи между парами событий, названной критическим отношением(critical interval relation). Поиск производится снизу вверх: алгоритм сначала находит простые закономерности, потом, путем соединения простых, образуются более сложные паттерны. На каждом шаге проводится отбор самых существенных и полных паттернов. Данный алгоритм был реализован на языке Си в виде консольного приложения и модуля Matlab.

# Содержание

1	Teo	ретическое обоснование. Определения	3
	1.1	Наблюдаемый временной ряд	3
	1.2	Понятие критического интервала	3
	1.3	Определение Т-Паттерна	4
2	Алі	соритм поиска Т-Паттернов	4
	2.1	Формальное описание	5
	2.2	О параметрах $\alpha$ и $N_{min}$	6
	2.3	Случайные паттерны	6
3	Ана	ализ результатов работы алгоритма	7
	3.1	Стратегия выбора критического интервала	7
	3.2	Запрет паттернов с одинаковыми событиями	
4	Дон	кументация модуля к Matlab	9
	4.1	Описание функций	9

## 1 Теоретическое обоснование. Определения

#### 1.1 Наблюдаемый временной ряд

Пусть время наблюдения разбито на  $N_t$  интервалов. В каждый момент nepuoda наблюдения  $(observation\ period)$   $[1, N_t]$  может произойти некоторое событие (deŭcmbue, event) из множества допустимых событий  $\mathcal{E}$   $(event\ types)$ . Соответственно, каждому типу событий сопоставляется множество моментов времени  $TS(\mathbf{A})$ 

$$TS(\mathbf{A}) = \{T_{\mathbf{A},1}, \dots, T_{\mathbf{A},N_{\mathbf{A}}}\}, \quad \mathbf{A} \in \mathcal{E}, \quad 0 \leqslant T_{\mathbf{A},i} \leqslant N_t \quad (i = 1, \dots, N_{\mathbf{A}})$$

.

#### 1.2 Понятие критического интервала

Во время поиска закономерностей в данных, нас интересуют отношения между распределениями отдельных событий. Если предположить, что в исходных не существует никаких закономерностей, то каждое событие должно появляться независимо, от других; то есть распределение компонентов должно быть независимым. Но паттерн характеризуется появлением своих компонентов в одинаковом порядке, более того, временные интервалы, разделяющие компоненты, должны быть примерно одинаковыми.

Будем говорить, что события **A** и **B** связаны отношением *критического интервала (Critical Interval, CI)*, если, после появления события **A** в момент времени t, существует интервал  $[t+d_1,t+d_2], (d2 \geqslant d1 \geqslant 0)$ , содержащий **B**, чаще, чем это ожидается из предположения о независимости событий. Данную взаимосвязь будем обозначать, как  $\mathbf{A}[d_1,d_2]\mathbf{B}$ , или, короче,  $(\mathbf{A}\mathbf{B})$ .

Далее раскроем понятие «чаще чем это ожидается». Пусть  $N_{\bf A}$  и  $N_{\bf B}$  — количество возникновений  ${\bf A}$  и  ${\bf B}$ , соответственно, в течение  $[1,N_t]$ .  $P({\bf A})=N_{\bf A}/N_t$  — вероятность появления события  ${\bf A}$  в некоторый момент времени;  $P(\neg {\bf A})=1-P({\bf A})$ .  $P(\neg {\bf A})^d$  — вероятность, что  ${\bf A}$  не появится в течение какого-либо интервала  $[d_1,d_2],\ (d=d_2-d_1+1),$  длины d. Вероятность наблюдать  ${\bf A}$  на интервале длины d один или более раз, равна  $1-P(\neg {\bf A})^d$ .

Зафиксируем события **A**, **B**, и длину интервала d. Приняв гипотезу о независимом распределении событий, событие **B** содержится<sup>2</sup> в интервале длины d, после события **A**,  $N_{\mathbf{A}} * (1 - P(\neg \mathbf{B})^d)$  раз.

$$\rho = P(\geqslant N_{\mathbf{AB}}) = 1 - P(< N_{\mathbf{AB}})$$

Чаще всего понимается, что в этот момент времени имеет место начало действия

<sup>&</sup>lt;sup>2</sup>Один или более раз.

— априорная вероятность того, что  $N_{\mathbf{AB}}$  из  $N_{\mathbf{A}}$  интервалов содержат вхождения  $\mathbf{B}$ . Очевидно, что  $P(< N_{\mathbf{AB}})$  распределено по биномиальному закону, где  $N_{\mathbf{A}}$  — количество «испытаний»,  $1 - P(\neg \mathbf{B})^d$  — вероятность «успеха». Следовательно,

$$\rho = P(\geqslant N_{\mathbf{AB}}) = 1 - \sum_{i=0}^{N_{\mathbf{AB}}-1} C_{N_{\mathbf{A}}}^{i} (1 - P(\neg \mathbf{B})^{d})^{i} P(\neg \mathbf{B})^{N_{\mathbf{A}}-i}$$

Полученная вероятность  $\rho$  сравнивается с пороговым значением  $\alpha$ , являющимся структурным параметром поиска: если  $\rho \leqslant \alpha$ , то заданный интервал признается критическим. Заметим, что  $\rho$  зависит от  $N_{\bf A}, N_{\bf B}, N_t, N_{\bf AB}, d$ :

$$\rho = \rho(N_{\mathbf{A}}, N_{\mathbf{B}}, N_t, N_{\mathbf{AB}}, d)$$

## 1.3 Определение Т-Паттерна

Дадим рекурсивное определение Т-Паттерна. Договоримся называть каждое допустимое событие nceedonammephom. Тогда Т-Паттерн  $\mathbf{Q}$  можно определить как:

$$\mathbf{Q} = \mathbf{X_1}[dL_1, dR_1]\mathbf{X_2}[dL_2, dR_2] \dots \mathbf{X_i}[dL_i, dR_i]\mathbf{X_{i+1}} \dots \mathbf{X_m},$$
 
$$Events(\mathbf{Q}) = \begin{cases} \{Events(\mathbf{X_1}), \dots Events(\mathbf{X_m})\}, & \text{если } \mathbf{Q} - \text{Т-Паттерн} \\ Q, & \text{если } \mathbf{Q} - \text{событие (псевдопаттерн)} \end{cases}$$

где  $\mathbf{X_i}$ ,  $(i=1\dots m)$  — Т-Паттерн, или псевдопаттерн. Для паттернов отношение критического интервала  $\mathbf{Q_L}[dL,dR]\mathbf{Q_R}$  вводится, как и для событий, с учетом того, что интервал [dL,dR] отсчитывается от последнего элемента  $\mathbf{Q_L}$ , и вхождение  $\mathbf{Q_R}$  определяется его первым элементом. Будем называть  $\partial soйными$  сериями(double series, DS) паттерна  $\mathbf{Q}$ , множество  $\{\{Left_i, Right_i\}_{i=1\dots N_Q}\}$ , где  $N_{\mathbf{Q}}$  — количество появлений паттерна  $\mathbf{Q}$ ,  $Left_i, Right_i$  — индексы начального и конечного событий i-го появления паттерна  $\mathbf{Q}$ .

## 2 Алгоритм поиска Т-Паттернов

Алгоритм поиска Т-Паттернов, описанный ниже, заключается в итеративном повторении двух стадий: конструирование новых паттернов и удаление неполных паттернов. На выходе алгоритм выдает множество Т-Паттернов с их сопутствующими характеристиками: критическим интервалом, уровнем значимости, частотой встречаемости.

Пусть  $\mathcal{D}_i$  — множество паттернов, обнаруженных к i-ой итерации. Фактически, множество

$$\mathcal{D}_m \setminus \mathcal{E}$$
,

где m — номер последней итерации, и будет результатом работы алгоритма.

#### Стадия 1: Поиск и конструирование:

На данном шаге, для любой упорядоченной пары

$$(\mathbf{Q}',\mathbf{Q}''):\mathbf{Q}',\mathbf{Q}''\in\mathcal{D}_i,$$

проверяется существование критической связи. Если критическая связь [dL, dR] была найдена, и паттерн  $(\mathbf{Q'Q''})$  встречается чаще, чем  $N_{min}$  раз, то в множество  $\mathcal{D}_{i+1}$  добавляется новый паттерн  $\mathbf{Q'}[dL, dR]\mathbf{Q''}$ .

Однако такие действия приводят к тому, что один и тот же паттерн **ABCD** может быть сконструирован разными способами (например как (**A**(**BCD**)) и ((**AB**)(**CD**))), что ведет к заполнению множества  $\mathcal{D}$  лишними паттернами, и как следствие, к замедлению работы программы. Для избежания данной проблемы предлагается следующее решение: паттерн (**Q**'**Q**") добавляется в множество  $\mathcal{D}$  тогда и только тогда, когда не существует паттерна **P** из  $\mathcal{D}$  такого, что  $Events(\mathbf{P}) = Events((\mathbf{Q}'\mathbf{Q}"))$  и  $DS(\mathbf{P}) = DS((\mathbf{Q}'\mathbf{Q}"))$ . Другими словами, найденный паттерн, перед добавлением, сравнивается с уже существующими паттернами и проверяется, не является ли он дубликатом.

#### Стадия 2: Удаление неполных паттернов:

На данной стадии алгоритм стремится удалить найденные паттерны, являющиеся меньшими частями, или неполными версиями других обнаруженных паттернов. Для индикации таких паттернов можно применять разные эвристики. Ниже опишем условие, взятое из статьи [?].

Итак, паттерн  $\mathbf{Q_x}$  считается менее полным, чем  $\mathbf{Q_y}$ , если  $\mathbf{Q_x}$  и  $\mathbf{Q_y}$  появляются одинаково часто, и все события возникающие в  $\mathbf{Q_x}$ , также возникают в  $\mathbf{Q_y}$ .

### 2.1 Формальное описание

**Require:**  $\mathcal{E}$  — допустимые события,

 $N_t$  — продолжительность наблюдения,

 $\alpha$  — минимальный уровень значимости,

 $N_{min}$  — минимальное количество вхождений паттерна.

- 1:  $\mathcal{D}_{-1} = \emptyset$
- 2:  $\mathcal{D}_0 = \mathcal{E}$

⊳ Инициализируем множество паттернов

- 3: t = 0
- 4: while  $\mathcal{D}_t \neq \mathcal{D}_{t-1}$  do
- 5: t = t + 1
- 6:  $\mathcal{D}_t = \mathcal{D}_{t-1}$
- 7:  $\mathbf{for} \ P_L \in \mathcal{D}_{t-1} \ \mathbf{do}$   $\triangleright$  Стадия 1
- 8: for  $P_R \in \mathcal{D}_{t-1}$  do

```
9:
                 for all d_L, d_R \in [1, N_t] таких, что d_L \leqslant d_R do
                      if \rho(N_{P_L}, N_{P_R}, N_t, N_{P_L P_R}, d_R - d_L + 1) < \alpha then
10:
                          if isUnique(P_L \cup P_R, \mathcal{D}_t) then
11:
                               \mathcal{D}_t = \mathcal{D}_t \cup \{P_L \cup P_R\}
12:
         for P_L \in \mathcal{D}_t \setminus \mathcal{E} do
                                                                                                  ⊳ Стадия 2
13:
             for P_R \in \mathcal{D}_t do
14:
                 if P_L \subset P_R^3 and |DS(P_L)| = |DS(P_R)| and isIntersect(P_L, P_R) then
15:
16:
                      удалить P_L из \mathcal{D}_t
17: function ISUNIQUE(Q, \mathcal{D})
         for P \in \mathcal{D} do
18:
             if DS(Q)=DS(P) and Events(Q)=Events(P) then
19:
20:
                 return false
21:
        return true
22: function ISINTERSECT(P_L, P_R)
         for i = 1 ... |DS(P_L)| do
23:
             if DS_{i,Left}(P_L) > DS_{i,Right}(P_R) or DS_{i,Right}(P_L) < DS_{i,Left}(P_R) then
24:
                 return false
25:
26:
             return true
27:
         return true
```

### 2.2 О параметрах $\alpha$ и $N_{min}$

Так как для задачи поиска Т-Паттернов не вводится определение функционала качества найденных паттернов, то алгоритм требует ручной настройки параметров  $\alpha$  и  $N_{min}$ . Выбор значений параметров должен основываться на специфике наблюдаемых процессов и ожидаемых результатов. Однако выбор значений  $\alpha=0.005$  и  $N_{min}=3$  обычно удовлетворителен [?, с. 99]. Для более тонкой настройки поиска, можно использовать разные  $\alpha$  и  $N_{min}$  для паттернов разной длины. Например, при значениях  $\alpha=0.00001,\ N_{min}=7$  для паттернов длины 2, и  $\alpha=0.005,\ N_{min}=3$  для паттернов длины 3 и более, алгоритм найдет только несколько самых ярко выраженных патернов длины 2, и уже потом будет конструировать множество более длиных паттернов.

## 2.3 Случайные паттерны

Когда исследования проводятся на больших наборах данных, то Т-Паттерны могут возникать даже, если выборка была сгенерированна случайно. Поэтому, для найденного множества паттернов было бы полезно оценить, являются ли они случайными,

<sup>&</sup>lt;sup>3</sup>Имеется в виду упорядоченная вложенность множеств. Т.е.  $abd \subset abcd$ , но  $bdc \not\subset abcd$ 

или «структурными». Одним из подходов к решению данной задачи, является анализ рандомизированных данных:

По множеству  $\mathcal{E}$  строится множество  $\mathcal{E}'$ :

$$\mathcal{E}'=\{rand(\mathbf{A})|\mathbf{A}\in\mathcal{E}\}$$
, где $rand(\mathbf{A})=\mathbf{A}':TS(\mathbf{A}')=\{\xi_1,\ldots,\xi_{N_\mathbf{A}}\}$   $\xi_1,\ldots,\xi_{N_\mathbf{A}}\sim\mathbf{U}[1,N_t].$ 

Проще говоря, создаются данные с такой же протяженностью периода наблюдения и мощностью множества допустимых событий, но для каждого допустимого события  $\mathbf{A}$ , моменты времени его появления генерируются случайно с вероятностью появления  $N_{\mathbf{A}}/N_t$ .

Для полученного множества  $\mathcal{E}'$  применяется процесс поиска паттернов с *теми* энсе параметрами, которые применялись на исходных данных. Описанные действия исполняются несколько раз, после чего, результаты поиска на рандомизированных данных сравниваются с исходными.

Считается, что поиск Т-Паттернов прошел успешно, если в исходных данных было выявлено значительно больше Т-Паттернов, чем в рандомизированных, или они оказались длиннее.

## 3 Анализ результатов работы алгоритма

## 3.1 Стратегия выбора критического интервала

Во время поиска связи критического интервала  $\mathbf{Q_L}[dL,dR]\mathbf{Q_R}$  между двумя паттернами, вообще говоря, проверяются все возможные интервалы [dL,dR]. На практике, не является редкостью случай, когда для двух паттернов существует несколько пар dL и dR, удовлетворяющих отношению критического интервала. Для выбора конкретных значений dL и dR, предлагается использовать одну из нижеописанных стратегий. Для каждой стратегии представлен результат работы алгоритма на тестовых данных. Данные содержат паттерн длины 7, встречающийся 8 раз.

Выбор кратчайшего критического интервала: При использовании данной стратегии, на каждом шаге выбирается критический интервал, имеющий наименьшую длину d (d = dR - dL + 1). Такой подход позволяет уменьшить длину критических связей, тем самым выявляя более выраженные и «стройные» паттерны. Одним из недостатков данного метода является эффект «расщепления». Поясним данный эффект на примере: пусть в исходных данных существует критическая связь  $\mathbf{A}[4,20]\mathbf{B}$ , наблюдаемая 15 раз. Данная связь настолько ярко

выражена, что алгоритму не требуется подбирать границы критического интервала, чтобы  ${\bf B}$  появлялось после  ${\bf A}$  все 15 раз. Алгоритм, скорее, выделит два(или даже больше) критических интервала:  ${\bf A}[4,12]{\bf B}$ , наблюдаемый 6 раз, и  ${\bf A}[13,18]{\bf B}$ , наблюдаемый 7 раз. Таким образом, один ярко выраженный паттерн будет распадаться на несколько более редких, что в свою очередь, может помешать дальнейшему выявлению закономерностей. Ниже описанный метод позволяет избавиться от такой проблемы.

Выбор самого длинного критического интервала: Как следует из названия, при использовании этой стратегии, среди всех значимых [dL, dR], выбирается интервал, соответствующий наибольшему значению d. В результате некоторые паттерны «загрубляются», и алгоритм стремится найти максимальное количество вхождений каждого критического интервала.

**Выбор самого значимого критического интервала:** При этом подходе выбирается критический интервал, имеющий максимальный уровень значимости  $\rho$ .

В общем случае, нету каких-либо рекомендаций по выбору стратегии поиска критического интервала. Выбор должен основываться на априорных сведениях о наблюдаемом процессе и ожидаемых результатов эксперимента.

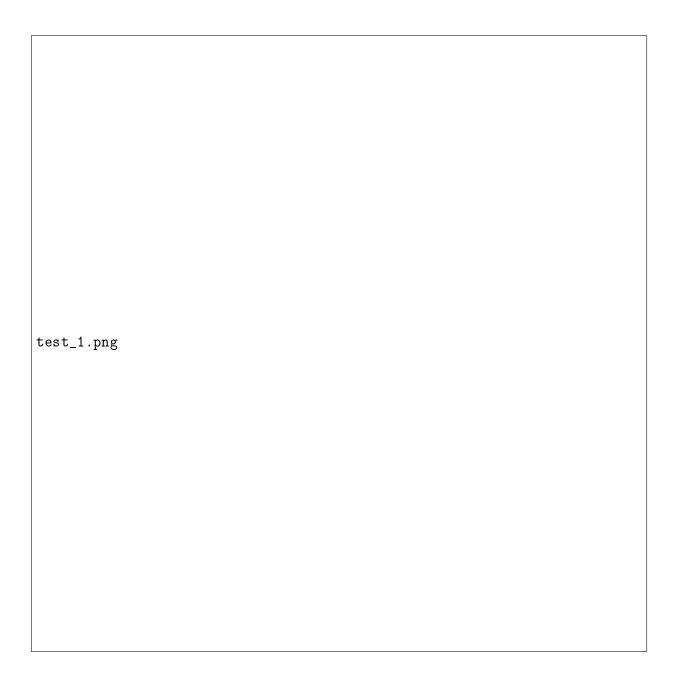


Рис. 1: Визуализация исходного паттерна abcdef.

Стратегия	Найдено паттернов	Найдено вхождений целевого паттерна
Короткий	40	Найдено 3 представления целевого паттерна.
		Для каждого представления 4 вхождения.
Длинный	5	8
Значимый	7	7

#### 3.2 Запрет паттернов с одинаковыми событиями

Если при поиске поведенческих закономерностей известно, что в один и тот же паттерн не могут входить два одинаковых события, то существует возможность сообщить это алгоритму. Таким образом, если,  $Events(\mathbf{Q_L}) \cap Events(\mathbf{Q_R}) \neq \varnothing$ , то паттерн  $(\mathbf{Q_LQ_R})$ , просто не будет создан. Это условие позволяет отбросить множество лишних паттернов, и найти более длинные и ценные закономерности.

## 4 Документация модуля к Matlab

Matlab-реализация алгоритма состоит из следующих файлов:

mexPattern.mex\*: откомпилированный mex-файл, реализующий алгоритм поиска паттернов.

mexPattern.m: объявление mex-функции.

T DRAW PATTERNS.m: графический вывод найденных паттернов.

T\_GENERATE\_PATTERNS.m: создание искусственных паттернов во временных рядах.

**T\_LOAD\_FILE.m:** загрузка временного ряда из файла для дальнейшей работы с ним.

Т STAT VALIDATE.m: процедура статистической валидации.

test.m: пример использования модуля.

### 4.1 Описание функций

Создает временной ряд, содержащий один искусственный паттерн.

Параметр	Описание
Вход:	
pat_sym	Матрица $\mathtt{1x}N_p$ типа $\mathtt{char},$ определяющая паттерн. Каждый символ — событие.
noise_sym	Матрица 1xN типа char. Определяет события которые будут генерироваться случайно.

CIs	Матрица ( $N_p$ -1)х2 типа int. В $i$ -й строке которой, записан со-
	ответствующий критический интервал.
Npat	Количество паттернов, которые требуется сгенерировать.
dist_b_patterns	Максимальное расстояние между двумя появлениями паттерна.
P_noise1	Частота встречаемости шумовых символов.
P_noise2	Вероятность того, что символ из паттерна будет зашумлен.
Выход:	
events	Массив структур 1xN, где N — количество событий. Каждая
	структура состоит из двух полей: event_name — строка назва-
	ния события, и indexes — матрица 1xN типа int. Определяет
	времеа появления событий.
Nt	Продолжительность получившегося периода наблюдений.
ts	Символьная матрица 1xNt.

function patterns = mexPattern(events, Nt, levels, allow\_same\_events, ci\_strategy);
Реализует поиск паттернов во временных рядах.

длину паттернов, к которым должны применяться следую параметры; минимальный уровень значимости $\alpha$ ; минимал количество вхождений паттерна $N_{min}$ .  аllow_same_events 1 если, разрешается появление одинаковых событий в патте 0 иначе.  Стратегия выбора критического интервала. 1 — стратегия бора длиннейшего интервала, 2 — кратчайшего, 3 — самого чимого.  Выход:  раtterns Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие	Параметр	Описание
levels Vатрица Nx3 уровней значимости. Каждая строка содерод длину паттернов, к которым должны применяться следую параметры; минимальный уровень значимости $\alpha$ ; минимальный уровень $\alpha$ ; минимальный уровень $\alpha$	Вход:	
длину паттернов, к которым должны применяться следую параметры; минимальный уровень значимости $\alpha$ ; минимал количество вхождений паттерна $N_{min}$ .  аllow_same_events 1 если, разрешается появление одинаковых событий в патте 0 иначе.  Сі_strategy Стратегия выбора критического интервала. 1 — стратегия бора длиннейшего интервала, 2 — кратчайшего, 3 — самого чимого.  Выход:  раtterns Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. СІѕ — интервалы, соответствующие	textttevents, Nt	См. определение T_GENERATE_PATTERN.
allow_same_events       1 если, разрешается появление одинаковых событий в патте о иначе.         ci_strategy       Стратегия выбора критического интервала. 1 — стратегия бора длиннейшего интервала, 2 — кратчайшего, 3 — самого чимого.         Выход:       раtterns         Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие	levels	Vатрица Nx3 уровней значимости. Каждая строка содержит: длину паттернов, к которым должны применяться следующие параметры; минимальный уровень значимости $\alpha$ ; минимальное
о иначе.  Стратегия выбора критического интервала. 1 — стратегия бора длиннейшего интервала, 2 — кратчайшего, 3 — самого чимого.  Выход:  раtterns  Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие		
бора длиннейшего интервала, 2 — кратчайшего, 3 — самого чимого.  Выход: раtterns Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие	allow_same_events	
чимого.         Выход:         patterns       Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие	ci_strategy	Стратегия выбора критического интервала. 1 — стратегия вы-
Выход:           patterns         Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие		бора длиннейшего интервала, $2$ — кратчайшего, $3$ — самого зна-
раtterns Массив структур. Каждая структура описывает найден паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие		чимого.
паттерн. Поля структуры: Events — индексы событий, кото составляют паттерн. CIs — интервалы, соответствующие	Выход:	
паттерна. DS — двойные серии $(\mathit{double\ series})$ паттерна. Stri	patterns	Массив структур. Каждая структура описывает найденный паттерн. Поля структуры: Events — индексы событий, которые составляют паттерн. CIs — интервалы, соответствующие критическим связям между событиями в паттерне. Sign — уровень значимости найденного паттерна. Nab — количество появлений паттерна. DS — двойные серии(double series) паттерна. String —
строка, описывающая паттерн в следующем формате: $Event_1[dL_1,dR_1]Event_2\dots Event_m<$ уровень значимости $>\{N_{ab}\}:(DS_1)(DS_2)\dots$		

#### function T\_DRAW\_PATTERNS(patterns, events, Nt, np)

Строит диаграмму найденных паттернов.

Параметр	Описание
Bxod:	
patterns, events,	См. предыдущие определения.
Nt	
np	Номер паттерна, который нужно представить. Или -1 для по-
	следовательного вывода всех паттернов.

function [events, Nt] = T\_LOAD\_FILE(fname)

Загружает временной ряд из фалйа для последующей работы с ним. Формат входного файла:

Time Event
0 :
time event
...
time event
time &

function [p] = T\_STAT\_VALIDATE( Nt, events, levels, nvalidations )

Процедура статистической валидации.

Параметр	Описание
Вход:	
Nt, events,	См. предыдущие определения.
levels	
nvalidations	Количество повторений процедуры рандомизации.
Bыход:	
p	Целочисленная матрица 1 х nvalidations. В каждой ячейке —
	количество паттернов, найденных в рандомизированных дан-
	ных.

# Список литературы

- [1] Magnus Magnusson
- [2] Кнут Д. Всё про Т<br/>EX. Протвино, RDT<br/>EX, 1993.