

# Вероятностный подход к поиску поведенческих паттернов\*

Вишневский В. В., Ветров Д. П.

valera.vishnevskiy@yandex.ru, vetrovd@yandex.ru

Москва, Московский Государственный Университет им. М. В. Ломносова

В данной работе предложен новый метод для поиска скрытых закономерностей (паттернов) в последовательностях событий, основанный на вероятностном представлении Р-Паттернов (probabilistic pattern) в дискретных последовательностях событий. Поиск производится снизу вверх: сначала находятся простые закономерности, потом, путем их соединения, образуются более сложные паттерны. Рассматривается применение данного алгоритма для анализа поведения мышей. Найденные паттерны используются для классификации животных. Проведено сравнение реализованного алгоритма с существующими аналогами, показавшее, что предложенный метод более устойчив к шуму в исходных данных.

Задача поиска закономерностей (стереотипов, паттернов, шаблонов — здесь синонимы) в поведении животных и людей крайне важна в современной нейробиологии и когнитивных науках. Выделив характерные паттерны, можно, например, делать выводы о сложности поведения различных особей, определять изменения в поведении наблюдаемых процессов, другими словами, решив задачу поиска паттернов, мы можем определенным образом измерять поведение особи, или группы особей.

В данной работе мы будем рассматривать временные паттерны в «структурном», или «эффектном» описании [3, с. 57]. Исходными данными будет размеченное поведение особи, то есть последовательность пар «момент времени», «поведенческий акт». Неформально можно сказать, что интересующие нас паттерны — это упорядоченная последовательность поведенческих актов, следующие один за другим через относительно инвариантные временные интервалы. Причем, этот паттерн должен повторяться в исходных данных достаточно часто.

Несмотря на то, что описанные выше паттерны широко распространены в описании поведения, стандартные статистические методы не подходят для их поиска: эти методы либо не учитывают всю сложность паттернов (например, периодические орбиты [4]), либо оперируют такими понятиями как циклы, волны, тренды, что невозможно напрямую использовать для поиска интересующих нас паттернов.

На сегодняшний день, для анализа таких поведенческих закономерностей наиболее широкое распространение получил метод поиска Т-Паттернов

нов (temporal patterns), предложенный в 2000-ом году Магнусом Магнуссоном в [2].

## Основные определения. Понятие Т-Паттерна

Пусть время наблюдения разбито на  $N_t$  интервалов. В каждый момент периода наблюдения  $[1, N_t]$  может произойти некоторое событие  $e$  (действие, поведенческий акт, event) <sup>1</sup> из множества допустимых событий  $\mathcal{E}$ . Соответственно, каждому типу события сопоставляется множество моментов времени  $TS(e)$ , когда это событие имело место:

$$TS(e) = \{t_1^e, \dots, t_{N_e}^e\}, \\ e \in \mathcal{E}, \quad 0 \leq t_i^e \leq N_t, \quad (i = 1, \dots, N_e),$$

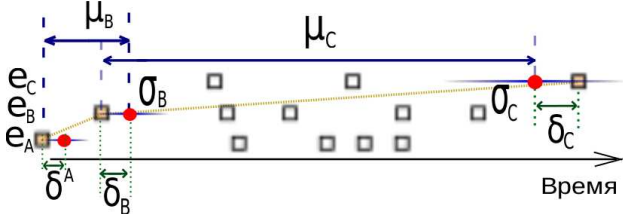
здесь  $N_e$  — количество появлений события  $e$  в данных.

**Понятие Т-Паттерна** включает в себя определение модели связи последовательных событий и способа их появления в данных (т.е. определения, что данная структура не случайна, а является закономерностью). Каждое событие паттерна определяется фиксированным временным интервалом, в течение которого это событие должно присутствовать после предыдущего события. Другими словами, расстояния между событиями моделируются равномерным распределением.

**Основной недостаток метода поиска Т-Паттернов** заключается, во-первых, в том, что само определение Т-Паттерна не позволяет ему иметь пропуски событий. По этой причине метод становится крайне чувствителен к шуму в исходных данных, из-за чего можно пропустить информативные длинные и сложные паттерны. Во-вторых, полученные Т-Паттерны сильно специфичны особи, в поведении которой они были найдены.

<sup>1</sup>Чаще всего понимается, что в этот момент времени имеет место начало действия

Работа выполнена при финансовой поддержке РФФИ (проект №08-01-00405), гранта Президента РФ (МК3827.2010.9), и федеральной целевой программы «Научные и научно-педагогические кадры инновационной России» на 2009–2013 годы (контракт №П1265). Авторы статьи выражают благодарность членам «Лаборатории Нейробиологии Памяти» при Институте Нормальной Физиологии П.К. Анохина, возглавляемой членом-корреспондентом РАН К.В. Анохиным. Отдельно благодарим Ирину Зарайскую, предоставившую нам экспериментальные данные по поведению.



**Рис. 1.** Вхождение нечеткого паттерна  $\mathbf{P} = e_A[\mu_A, \sigma_A]e_B[\mu_B, \sigma_B]e_C[\mu_C, \sigma_C]$ . в данные. Маркеры-кружки соответствуют ожидаемым позициям событиям, закрашенные маркеры-квадраты соответствуют позиции реальной позиции соответствующего события в данных.

## Вероятностная модель паттерна. Р-Паттерны

**Определение 1.** Нечетким паттерном, или Р-Паттерном  $\mathbf{P}$  длины  $N_{\mathbf{P}}$  назовем упорядоченную последовательность событий  $e_i$ , ( $i = 1, \dots, N_{\mathbf{P}}$ ), где каждое событие паттерна характеризуется смещением и разбросом от предыдущего события. Будем записывать паттерн  $\mathbf{P}$  в следующем виде:

$$\mathbf{P} = [\mu_1, \sigma_1]e_1[\mu_2, \sigma_2]e_2 \dots [\mu_{N_{\mathbf{P}}}, \sigma_{N_{\mathbf{P}}}]e_{N_{\mathbf{P}}}, \quad \mu_1 = 0.$$

Здесь  $\mu_i$  и  $\sigma_i$  — математическое ожидание и корень из дисперсии нормального распределения, моделирующего величину времени, прошедшего между событиями.

Представление Р-Паттерна иллюстрировано на рис. 1.

Далее, чтобы иметь возможность обрабатывать пропуски в Р-Паттернах, введем понятие *функции потерь*, которая определяет «штраф» за пропуск  $m$  событий в паттерне длины  $N_{\mathbf{P}}$  следующим образом:

$$f_{LOSS}(m, N) = \begin{cases} \exp\left(-\frac{\lambda m}{N_{\mathbf{P}}}\right), & m < N, \\ 0, & m = N. \end{cases}$$

Здесь  $\lambda$  является структурным параметром, определяющим уровень «нечеткости» паттернов. Если этот параметр велик, то мы, по сути, запрещаем реализациям паттерна иметь пропуски. Если выставить этот параметр слишком малым, то будут обнаруживаться паттерны, не разу полностью не встречающиеся в данных, то есть закономерности могут быть найдены даже в случайных данных.

**Определение 2.** Правдоподобие паттерна  $\mathbf{P}$  — это функция, определенная в каждый момент времени наблюдения  $\varepsilon$  ( $\varepsilon = 1, \dots, N_t$ ) следующим об-

разом:

$$L_{\mathbf{P}}(\varepsilon) = f_{LOSS}(N_{-}, N_{\mathbf{P}}) \prod_{i=1}^{N_{\mathbf{P}}} \left( \frac{1}{\sqrt{2\pi} \sigma_i} \right) \times \prod_{i \in N_{+}} \exp\left(-\frac{\delta_i^2}{2\sigma_i^2}\right), \quad (1)$$

где  $\delta_i$  — расстояния от ожидаемой позиции события в Р-Паттерне до ближайшего события в данных (более наглядно см. рис. 1). Т.е.:

$$\delta_i = \min_{x \in TS(e_i)} \left| \underbrace{\varepsilon + \sum_{j=1}^{i-1} (\mu_j + \delta_j) + \mu_i}_{\text{ожидаемая позиция события}} - x \right|,$$

здесь, если событие было пропущено, то соответствующее  $\delta_i = 0$ . Далее,  $N_{-}$  — количество пропущенных событий в паттерне, а  $N_{+}$  — множество индексов присутствующих в паттерне событий. Событие считается пропущенным, если  $\exp\left(-\frac{\delta_i^2}{2\sigma_i^2}\right) < \exp\left(-\frac{\lambda}{N_{\mathbf{P}}}\right)$ .

По сути, правдоподобие показывает насколько можно быть уверенным, что данный Р-Паттерн начинается в определенный момент времени  $\varepsilon$ .

Заметим, что правдоподобие Р-Паттерна может быть отсчитано с конца, или с  $m$ -го события паттерна.

**Утверждение 1.** Математическое ожидание функции правдоподобия (1) в момент времени  $\varepsilon$ , при условии, что в данный момент времени имеет место начало модельного Р-Паттерна  $\mathbf{P}$  вычисляется следующим образом:

$$E[L_{\mathbf{P}}(\varepsilon)] = \frac{1}{(2\sqrt{\pi})^{N_{\mathbf{P}}} \sigma_1 \dots \sigma_{N_{\mathbf{P}}}}.$$

Доказательство приведено в [1].

Здесь используется тот факт, что межточечные расстояния между событиями в модельном Р-Паттерне распределены по нормальному закону:

$$\delta_i \sim \mathcal{N}(0, \sigma_i), \quad (i = 1, \dots, N_{\mathbf{P}}).$$

Теперь мы можем считать, что Р-Паттерн  $\mathbf{P}$  имеет место быть только в следующие моменты времени:

$$t: L_{\mathbf{P}}(t) \geq \gamma E[L_{\mathbf{P}}(\varepsilon)], \quad (2)$$

где  $\gamma$  — заданная константа.

**Конструирование Р-Паттернов.** Рассмотрим пару Р-Паттернов  $\mathbf{P}_L$  (левый) и  $\mathbf{P}_R$  (правый).

Пусть  $\{\alpha_i\}_{i=1,\dots,N_L}$  и  $\{\beta_j\}_{j=1,\dots,N_R}$  — значения правдоподобия соответствующих Р-Паттернов в моменты времени, соответствующие паттерны имели вхождения (2). Важно, что правдоподобие левого Р-Паттерна отсчитывается с конца, так как мы ищем связь между концом левого паттерна и началом правого. Также пусть,  $\{t_{L,i}\}_{i=1,\dots,N_L}$  и  $\{t_{R,j}\}_{j=1,\dots,N_R}$  — моменты времени, когда эти Р-Паттерны имели место в смысле (2).  $N_L$  и  $N_R$  — количество вхождений паттернов  $\mathbf{P}_L$  и  $\mathbf{P}_R$ , соответственно. Определим множество межточечных расстояний:

$$\rho = \{t_{R,j} - t_{L,i} \mid t_{R,j} \geq t_{L,i}\}.$$

Для каждого расстояния из этого множества введем соответствующий вес  $w_l = \ln(1 + \alpha_i \beta_j)$ ,  $l = 1, \dots, M$ , где  $M = |\rho|$ .

Рассмотрим гипотезу  $H_0$ , что моменты времени и веса вхождения Р-Паттернов распределены независимо и равномерно на всем наблюдаемом промежутке. Тогда плотность распределения введенных выше межточечных расстояний  $\{t\}$  имеет следующий вид [1]:

$$p_{LR}(t) = \begin{cases} (N_t - t) \frac{2}{N_t^2}, & t \in [0, N_t], \\ 0, & x \notin [0, N_t]. \end{cases} \quad (3)$$

Введем статистическую модель связи между паттернами (проверяемые параметры связи  $\mu$  и  $\sigma$  фиксированы):

$$g_{\mu,\sigma}(t_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t_i - \mu)^2}{2\sigma^2}\right).$$

Рассмотрим следующую сумму:

$$k = \sum_{i=1}^M w_i g_{\mu,\sigma}(t_i), \quad (4)$$

где  $t_i \sim p_{LR}$ .

**Теорема 2.** При выполнении  $H_0$  статистика  $k$  распределена по нормальному закону:

$$\begin{aligned} k &\sim \mathcal{N}(\mu_*, \sigma_*^2), \text{ где} \\ \mu_* &\approx M \mathbb{E}[w] \frac{2}{N_t} \left(1 - \frac{\mu}{N_t}\right) \\ \sigma_*^2 &\approx \frac{M}{N_t} \left(1 - \frac{\mu}{N_t}\right) \left[ \left(\frac{1}{\sqrt{\pi}\sigma} - \frac{\mu}{N_t} \left(1 - \frac{\mu}{N_t}\right)\right) \times \right. \\ &\times \left. ((\mathbb{E}[w])^2 + \mathbb{D}[w]) + \frac{4\mathbb{D}[w]}{N_t} \left(1 - \frac{\mu}{N_t}\right) \right], \end{aligned} \quad (5)$$

здесь  $\mathbb{E}[w]$  и  $\mathbb{D}[w]$  — выборочное среднее и дисперсия весов, соответственно.

Доказательство приведено в [1].

**Замечание 1.** Приближенные формулы (5) дают удовлетворительный результат для задачи поиска паттернов, однако, точные формулы для значений  $\mu_*$  и  $\sigma_*$ , можно найти в [1].

Теперь для различных пар  $\mu$  и  $\sigma$  можно вычислить статистику (4), сравнить ее с  $\alpha$ -квантилью распределения (5). Если гипотеза  $H_0$  о «случайности» данных будет отвергнута односторонним критерием, то считается, что соответствующие Р-Паттерны образуют новый паттерн с параметрами  $\mu$  и  $\sigma$ . Если существует несколько пар  $\mu$  и  $\sigma$ , для которых отвергается гипотеза  $H_0$ , то для конструирования Р-Паттернов берутся непересекающиеся<sup>2</sup> параметры соответствующие максимальным значениям  $k$ .

Более подробно о процессе формирования нового Р-Паттерна можно найти в [1].

**Редукция Р-Паттернов.** Для удаления [1] паттернов-дубликатов и неполных копий анализируется коэффициент корреляции функций правдоподобия. Пусть  $\vec{L}_{P,i}$  — вектор-столбец значений функции правдоподобия, отсчитанной от  $i$ -го события во всех моментах времени наблюдения.

$$\text{cor}(\vec{L}_1, \vec{L}_2) = \frac{\vec{L}_1^\top \vec{L}_2}{\sqrt{\vec{L}_1^\top \vec{L}_1} \sqrt{\vec{L}_2^\top \vec{L}_2}} \in [0, 1]$$

— коэффициент корреляции между двумя Р-Паттернами.

Проверяются все пары паттернов, если все поведенческие акты, присутствующие в паттерне  $\mathbf{P}_L$  также присутствуют в  $\mathbf{P}_R$  с учетом порядка, и

$$\exists m: \text{cor}(\vec{L}_{P_L,1}, \vec{L}_{P_R,m}) > \nu,$$

тогда паттерн  $\mathbf{P}_L$  удаляется из множества найденных паттернов.

#### Алгоритм поиска Р-Паттернов

1. Инициализировать текущее множество Р-Паттернов событиями (паттерны длины 1).
2. Для всевозможных пар Р-Паттернов из текущего множества провести процедуру *конструирования*.
3. Провести процедуру редукции паттернов.
4. Если текущее множество паттернов изменилось, перейти к п.2.

Параметры алгоритма и способы их настройки описаны в [1].

Сложность предложенного алгоритма —  $O(n^3)$ , где  $n$  — общее количество событий во временном ряде. В [1] представлена параллельная реализация данного метода на GPU, что позволило применять алгоритм поиска Р-Паттернов на реальных данных.

<sup>2</sup>  $[\mu' - 3\sigma', \mu' + 3\sigma'] \cap [\mu'' - 3\sigma'', \mu'' + 3\sigma''] = \emptyset$

## Эксперименты на реальных данных

Описанный ниже эксперимент демонстрирует способ применения предложенного метода на реальных поведенческих данных. Целью эксперимента является анализ того, как влияет отсутствие гиппокампа на поведение. Гиппокамп — один из древнейших отделов головного мозга млекопитающих, его функции связывают с механизмами работы памяти, обучением, пространственной навигацией.

Особь были разделены на 5 групп.

- 1) Контрольная группа, содержащая разметку поведения здоровых мышей. 12 особей.
- 2) Гиппокампальная группа. Гиппокамп этих животных разрушали путем введения в эту структуру лидокаина, растворенного в искусственной спинномозговой жидкости (2 мкл. 4% раствора). 12 особей.
- 3) Шумовая группа, с параметрами частоты и продолжительности актов первой (контрольной) группы. 12 «особей».
- 4) Шумовая группа, с параметрами частоты и продолжительности актов второй (гиппокампальной) группы. 12 «особей».
- 5) Искусственные данные, содержащие один модельный Р-Паттерн. 7 «особей».

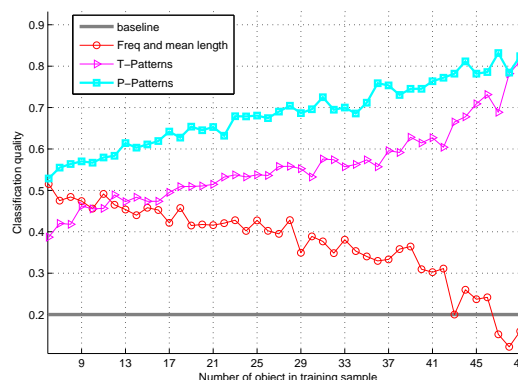
Поведение каждой особи было представлено временным рядом длины  $\sim 12$  минут, всего было 24 детектируемых поведенческих актов, более подробно условия эксперимента описаны в [1]. Требовалось решить задачу классификации особи по поведению.

Далее, данные разбивались на обучение ( $N_l$  объектов) и контроль ( $N_c$  объектов).  $N_l + N_c = 55$ . Каждая особь описывалась вектором длины  $N_l$ ,  $i$ -ое значение которого равняется количеству паттернов данной особи, также найденное в поведении  $i$ -ой особи из обучающей выборки.

На рис. 2 видно, что предложенный метод поиска Р-Паттернов дает заметно лучшее качество классификации, чем метод поиска Т-Паттернов. Наивная классификация на основе описания частот и средней продолжительности актов неприменима в данном эксперименте из-за двух шумовых классов.

**Характерные паттерны** присущие определенному классу также могут быть выделены в экспериментальных данных. Неформально, паттерн является характерным для заданного класса, если он присутствует в поведении многих особей этого класса и редко выявляется в поведении животных из других классов. Пример характерного для контрольной группы Р-паттерна (в квадратных скобках указаны смещения и дисперсии в секундах между событиями):

«Вылизывание гениталий» [2,5; 6,2] «Вылизыва-



**Рис. 2.** Качество — средняя доля правильных классификаций. По горизонтали откладывалось количество объектов в обучении (качество усреднено по ста случайным разбиениям). Классификации по Р- и Т-паттернам производилась с помощью SVM. Классификации на основе частот и длин актов производилась с помощью решающих лесов.

ние ладоней» [1,2; 7,9] «Умывание головы с ушами» [0,4; 5,7] «Вылизывание задних конечностей» [2,6; 7,9] «Умывание носа».

## Выводы

Представленный метод решает поставленные перед ним задачи и производит качественный поиск закономерностей как в синтетических временных рядах, так и в реальных поведенческих данных. В открытом доступе свободная, документированная, параллельная реализация представленного метода.

Главным преимуществом метода поиска Р-Паттернов является их вариабельность: если на то есть предпосылки, то Р-Паттерны найденные в поведении одной особи будут также найдены в поведении другой особи. Данный факт позволяет описывать поведение на основе найденных паттернов, и использовать стандартные алгоритмы машинного обучения для решения, например, задач классификации, кластеризации, или восстановления регрессии.

## Литература

- [1] В.В. Вишневский. «Параллельная реализация метода поиска закономерностей в последовательностях событий». — Дипломная работа. ВМиК МГУ, 2011.
- [2] M.S. Magnusson. Discovering hidden time patterns in behavior: T-patterns and their detection. — Behavior Research Methods, Instruments, Computers 2000.
- [3] P. Martin, P. Bateson. Measuring Behaviour: An Introductory Guide. — Cambridge University Press, second edition, 1993.
- [4] R. Stoop, B. Arthur. Periodic orbit analysis demonstrates genetic constraints, variability, and switching in Drosophila courtship behavior. — Chaos — 2008 — vol.18/2.