

Апробация модели интенсивности поведения со скрытыми переменными на данных респондентов о последних публикациях в сети Инстаграм

А. В. Торопова

Санкт-Петербургский государственный университет
alexandra.toropova@gmail.com

Т. В. Тулупьева

Санкт-Петербургский государственный университет;
Санкт-Петербургский федеральный исследовательский
центр РАН;
Северо-Западный институт управления РАНХиГС
tvt@dscs.pro

Аннотация. В разных областях науки существует множество задач, связанных с исследованием поведения человека. Одной из важнейших характеристик поведения является его интенсивность. Однако часто получить прямую оценку интенсивности поведения невозможно из-за ограничений в ресурсах и возможностях, поэтому разработка методов косвенной оценки характеристик поведения является актуальной. Ранее была предложена модель на основе байесовской сети доверия, оценивающая интенсивность поведения по данным о последних эпизодах поведения, учитывающая тот факт, что при ответах респонденты могут ошибаться. Был разработан опросник, с помощью которого были собраны ответы респондентов о последних эпизодах публикаций постов в социальной сети Инстаграм. Публикацию постов в социальной сети Инстаграм можно рассмотреть в качестве изучаемого поведения. Благодаря тому, что время публикации постов фиксируется, появляется возможность узнать действительную интенсивность этого поведения за исследуемый период. Таким образом, появилась возможность проверить модель интенсивности со скрытыми переменными на реальных данных. Цель данной статьи – апробация модели интенсивности поведения со скрытыми переменными на реальных данных. На основе полученных результатов можно оценить, насколько эффективным будет использование модели интенсивности поведения со скрытыми переменными в задачах, связанных с другими видами поведения, интенсивность которого нельзя узнать с помощью прямых методов.

Ключевые слова: байесовские сети доверия; оценка интенсивности поведения; интенсивность поведения; эпизоды поведения

I. ВВЕДЕНИЕ

В разных областях науки существует множество задач, связанных с исследованием поведения человека. Одной из важнейших характеристик поведения является его интенсивность. В ряде работ, зная интенсивность, исследователи оценивают определенные характеристики и параметры, связанные с поведением. В [1–3] по

Работа выполнена в рамках проекта по государственному заданию СПб ФИЦ АН № 0073-2019-0003, при финансовой поддержке РФФИ, проекты №19-37-90120, № 20-07-00839

интенсивности взаимодействия пользователей оценивается вероятность успеха социоинженерных атак. В [4] показано, что по интенсивности запросов Google пользователей об изоляции можно сделать выводы о распространении COVID-19.

Безусловно самым надежным способом оценить интенсивность поведения является прямое наблюдение, однако часто этот метод не может быть осуществим [5–6], в связи с этим требуются методы и подходы, позволяющие определить этот параметр косвенными методами, например, на основе данных респондентов, полученных в результате опроса.

В данной статье рассматривается функционирование модели интенсивности поведения со скрытыми переменными на данных, собранных с помощью специально разработанного опросника об эпизодах публикации постов в Инстаграм.

II. ОПИСАНИЕ МОДЕЛИ

Ранее [7] была предложена модель на основе байесовской сети доверия, оценивающая интенсивность поведения по данным о последних эпизодах поведения, учитывающая тот факт, что при ответах респонденты могут ошибаться. Модель представлена на рис. 1.

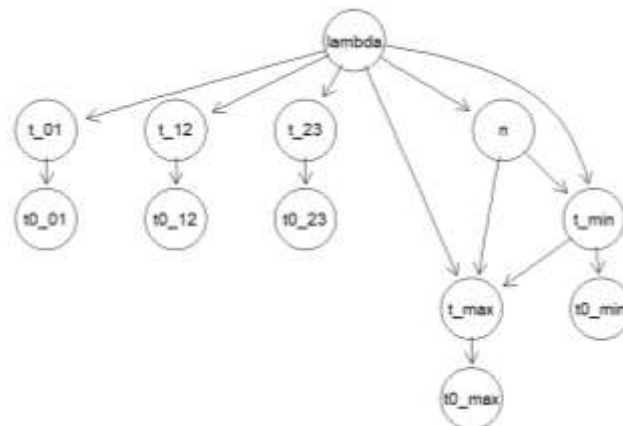


Рис. 1. Модель интенсивности поведения со скрытыми переменными

Вершина λ представляет собой оценку интенсивности поведения, эта оценка может быть получена на основе следующих данных, полученных от респондентов: временной интервал между моментом сбора данных и последним эпизодом исследуемого поведения (t_{0_01}), временной интервал между последним и предпоследним эпизодами поведения (t_{0_12}), временной интервал между предпоследним и предпредпоследним эпизодами поведения (t_{0_23}), минимальный временной интервал между эпизодами поведения за исследуемый период (t_{0_min}) и максимальный временной интервал между эпизодами поведения за исследуемый период (t_{0_max}). Так как ответы респондентов и реальные значения указанных временных интервалов могут отличаться (одна из наиболее вероятных причин состоит в том, что респонденты дают ответы по памяти и поэтому могут ошибиться) модель содержит скрытые переменные t_{01} , t_{12} , t_{23} , t_{min} , t_{max} , соответствующие реальным значениям временных интервалов. Вершина p характеризует количество эпизодов поведения за исследуемый период.

III. ОПИСАНИЕ ДАННЫХ

Для апробации модели было решено использовать данные о публикации постов в Инстаграм. Инстаграм является одной из самых популярных социальных сетей в России [8]. Публикация постов в Инстаграм может рассматриваться как вид поведения, а благодаря тому, что время постов в Инстаграм фиксируется, появляется возможность сравнить оценку интенсивности, предсказанную моделью с настоящим значением интенсивности этого поведения.

Для сбора данных респондентов был разработан опросник на основе инструментов Google Сайты, Google Формы и Google Таблицы. Респондент заходит на сайт опросника [9], отвечает на вопросы, после этого введенные сведения через дополнительную форму помещаются в таблицу.

Опросник содержит следующие вопросы: имя пользователя Инстаграм, сведения об интервалах между последними тремя эпизодами публикации постов и сведения о минимальном и максимальном интервалах между публикациями за год.

Ответ об интервалах между последними тремя эпизодами публикации постов может быть дан тремя способами:

- ввести дату и время (время необязательно, рис. 2);
- ввести интервал в выбранных единицах времени (рис. 3);
- ввести ответ в текстовое поле в свободной форме.

Ответ о минимальном и максимальном интервалах может быть дан в свободной форме или с помощью ввода интервала в выбранных единицах времени (рис. 3).

С помощью данного опросника были собраны данные о 90 пользователях Инстаграм.

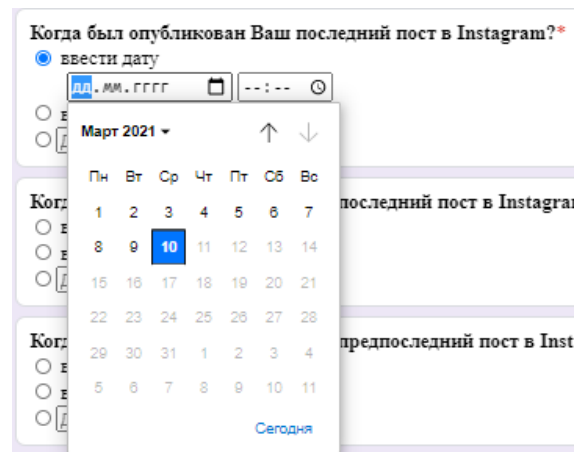


Рис. 2. Ввод даты в опросник

Для того, чтобы привести эти данные к единому виду, а также для того, чтобы определить значение интенсивности для каждого пользователя на языке C# была написана специальная программа.

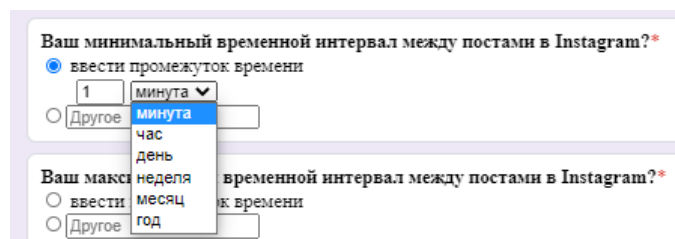


Рис. 3. Ввод временного интервала в опросник

Так как набранный датасет оказался довольно небольшим, для обучения модели были использованы синтетические данные: по гамма распределению были сгенерированы 500 значений интенсивностей, для каждого из этих значений было создано по 20 «респондентов», итоговый обучающий датасет содержит 10000 записей (более подробная информация о синтезе данных для обучения модели интенсивности со скрытыми переменными в [7]).

IV. ЭКСПЕРИМЕНТ

Для обучения и тестирования модели были использованы язык R [10] и пакет bnlearn [11] для работы с байесовскими сетями доверия.

Для работы с байесовскими сетями доверия требуется дискретизация непрерывных величин. Значения возможных значений интенсивности поведения (λ , измеряем как отношение количества эпизодов к числу дней в исследуемом периоде, то есть к 365) были разбиты на интервалы $\lambda_1=[0, 0.002]$, $\lambda_2=[0.002, 0.01]$, $\lambda_3=[0.01, 0.03]$, $\lambda_4=[0.03, 0.06]$, $\lambda_5=[0.06, 0.1]$, $\lambda_6=[0.1, 0.2]$, $\lambda_7=[0.2, 0.5]$, $\lambda_8=[0.5, \infty)$. Отметим, что первый интервал (λ_1) подобран таким образом, чтобы в него входили те пользователи Инстаграм, у которых нет публикаций за год. Значение временных интервалов (измеряем в днях) были разбиты на следующие интервалы: $[0, 0.1]$, $[0.1, 0.5]$, $[0.5, 1]$, $[1, 7]$, $[7, 14]$, $[14, 30]$, $[30, 180]$, $[180, 365]$, $[365, \infty)$.

На рис. 4 представлено распределение значений интенсивностей для пользователей Инстаграм, заполнивших опросник.

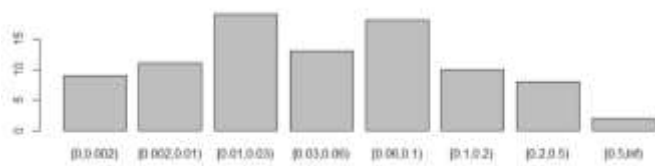


Рис. 4. Распределение интенсивностей для пользователей Инстаграм, участвовавших в опросе

После того как модель была обучена, на основе ответов 90 респондентов она сделала предсказания об интенсивности их публикаций в Инстаграм. Полученные предсказания можно сравнить с действительными значениями интенсивностей.

Табл. I представляет собой матрицу смежности, в которой строки обозначают действительные значения, а столбцы — значения интенсивностей, предсказанных моделью.

ТАБЛИЦА I МАТРИЦА СМЕЖНОСТИ

	λ_1	λ_2	λ_3	λ_4	λ_5	λ_6	λ_7	λ_8
λ_1	1	1	0	0	0	0	0	0
λ_2	1	5	1	0	1	0	0	0
λ_3	0	7	3	2	2	1	2	0
λ_4	0	2	3	4	0	3	0	0
λ_5	0	1	3	4	5	3	0	0
λ_6	0	0	1	2	2	3	1	0
λ_7	0	0	0	1	1	0	2	0
λ_8	0	0	0	0	0	0	0	0

Точность (accuracy) равна 0.338, средняя точность (average accuracy) равна 0.835, мера каппа равна 0.208. В табл. II представлены точность (precision), полнота (recall) и F-1, основные метрики качества по классам.

ТАБЛИЦА II МЕТРИКИ КАЧЕСТВА ПО КЛАССАМ

	Точность	Полнота	F-1
λ_1	0.5	0.5	0.5
λ_2	0.313	0.625	0.417
λ_3	0.273	0.176	0.214
λ_4	0.308	0.333	0.32
λ_5	0.455	0.313	0.37
λ_6	0.3	0.333	0.316
λ_7	0.4	0.5	0.444
λ_8	NaN	0	NaN

V. ЗАКЛЮЧЕНИЕ

Довольно высокий показатель средней точности говорит о том, что данная модель может быть применена

для оценки интенсивности в тех областях, где невозможно оценить этот параметр прямыми методами, а данные о поведении могут быть получены только с помощью опроса респондентов.

Для обучения модели были использованы синтетические данные, соответственно можно использовать эту модель, не имея большого набора данных для обучения модели.

Таким образом, на данных ответов респондентов об эпизодах их публикаций постов в Инстаграм была проведена апробация модели интенсивности со скрытыми переменными.

На основе полученных результатов можно сделать вывод, что описанная модель имеет потенциал для использования ее в областях науки, занимающимися исследованиями поведения человека, таких как социология, психология, эпидемиология и др.

СПИСОК ЛИТЕРАТУРЫ

- [1] Abramov M.V., Tulupyev A.L. Soft estimates of user protection from social engineering attacks: fuzzy combination of user vulnerabilities and malefactor competencies in the attacking impact success prediction // Artificial Intelligence and Natural Language. 2019. P. 47–58. DOI: 10.1007/978-3-030-34518-1_4.
- [2] Khlobystova A.O., Abramov M.V., Tulupyeva T.V. Application of the Alternatives Method Probabilities in Construction of Intensity of User Communications Estimates // 2020 XXIII International Conference on Soft Computing and Measurements (SCM), St. Petersburg, Russia, 2020, pp. 37–40, DOI: 10.1109/SCM50615.2020.9198751.
- [3] Khlobystova A.O., Abramov M.V., Tulupyev A.L. Soft Estimates for Social Engineering Attack Propagation Probabilities Depending on Interaction Rates Among Instagram Users // International Symposium on Intelligent and Distributed Computing. – Springer, Cham, 2019. P. 272–277. DOI: 10.1007/978-3-030-32258-8_32.
- [4] Bari A., Khubchandani A., Wang J. et al. COVID-19 early-alert signals using human behavior alternative data. Soc. Netw. Anal. Min. 11, 18, 2021., DOI: 10.1007/s13278-021-00723-5.
- [5] Mayer G.R., Sulzer-Azaroff B., Wallace, M. Behavior analysis for lasting change. Cornwall-on-Hudson, NY: Sloan Publishing. 2018.
- [6] Rehfeldt R.A. Clarifying the nature and purpose of behavioral assessment: A response to Newsome et al. // Journal of Contextual Behavioral Science. 2019. Vol. 14. P. 37–39.
- [7] Toropova A.V., Tulupyeva T.V. Synthesis and learning of socially significant behavior model with hidden variables // Advances in Intelligent Systems and Computing. 2019. T. 875. pp. 76–84.
- [8] SimilarWeb. URL: <https://www.similarweb.com/fr/top-websites/russian-federation> (дата обращения: 11.03.2021).
- [9] Опросник. Последние публикации в Instagram. URL: <https://sites.google.com/view/instagrampostsquestionary/> (дата обращения: 11.03.2021).
- [10] R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. URL: <http://www.R-project.org> (дата обращения: 11.03.2021).
- [11] bnlearn - an R package for Bayesian network learning and inference. URL: <https://www.bnlearn.com/> (дата обращения: 11.03.2021).