

# **АРХИТЕКТУРА АЛГЕБРАИЧЕСКОЙ БАЙЕСОВСКОЙ СЕТИ ДЛЯ ИДЕНТИФИКАЦИИ АККАУНТОВ ПОЛЬЗОВАТЕЛЕЙ В СОЦИАЛЬНЫХ СЕТЯХ «ВКОНТАКТЕ» И «ОДНОКЛАССНИКИ»<sup>1</sup>**

Корепанова А. А., младший научный сотрудник СПИИРАН; студент  
кафедры информатики СПбГУ, aak@dscs.pro

Олисеенко В. Д., младший научный сотрудник СПИИРАН; студент  
кафедры информатики СПбГУ, subster3@gmail.com

Абрамов М. В., кандидат технических наук, заведующий  
лабораторией теоретических и междисциплинарных проблем  
информатики СПИИРАН; доцент кафедры информатики СПбГУ,  
mva@dscs.pr

## **Аннотация**

В работе рассматривается подход к идентификации пользователей в социальных сетях «ВКонтакте» и «Одноклассники», основанный на алгебраических байесовских сетях. Данный подход позволит получать вероятностную оценку принадлежности профилей в различных социальных сетях одному человеку. На основе этой оценки планируется принимать решение об использовании информации из дополнительных источников (профиле в другой социальной сети) для построения профиля уязвимостей пользователя. Профиль уязвимостей пользователя является ключевым элементом в оценке защищенности пользователей от социоинженерных атак.

## **Введение**

Проблема автоматизированной идентификации аккаунтов пользователей в социальных сетях «ВКонтакте» и «Одноклассники» является не решенной. Её решение в рамках тематики социоинженерных атак позволило бы аккумулировать большее количество информации, способствующей оценкам степени выраженности личностных

---

<sup>1</sup> Работа выполнена в рамках проекта по государственному заданию СПИИРАН № 0073-2019-0003 и при финансовой поддержке РФФИ (гранты №18-01-00626, № 20-07-00839). средствами.

особенностей пользователя и, опосредованно, его уязвимостей. В статье [1] приводится обоснование актуальности данной проблемы и предлагается один из подходов для её решения. В соответствии с данным подходом решается задача бинарной классификации для двух аккаунтов и выбранных для них атрибутов «имя», «фамилия», «друзья», «город проживания», «дата рождения». Однако, предложенный подход имеет свои минусы – он не учитывает все доступные атрибуты профиля, из-за сложности работы с ними при наличии пропусков информации, и имеет невысокую точность при сопоставлении аккаунтов различных пользователей с одинаковыми именами и/или фамилиями, живущими в одном городе.

В качестве альтернативы данному подходу предлагается представление структуры алгебраической байесовской сети [2, 3, 4, 5] для получения информации об оценках вероятности принадлежности двух аккаунтов одному человеку. Алгебраические байесовские сети, как и байесовские сети доверия, относятся к классу вероятностных графических моделей [6, 7, 8]. Преимущество применения алгебраических байесовских сетей над другими методами заключается в возможности с их помощью обрабатывать неполные данные. В качестве данных, поступающих на вход алгебраической байесовской сети используется результат сопоставления значений атрибутов.

### **Сопоставление значений атрибутов**

Для сопоставления значений атрибутов данные из социальных сетей необходимо привести к единой форме. Например, для атрибута «имя» – привести к единой словоформе, транслитерировать и привести к нижнему регистру, таким образом «EVGENII» будет преобразовано к «жени». Более подробно процесс преобразования указан в [1].

Соответствующие значения атрибутов сравниваются друг с другом при помощи различных метрик (например Джаро-Винклера для атрибутов «имя», «фамилия», «город»). На основе полученных результатов сравнения с помощью алгебраической байесовской сети будет строиться предположение о принадлежности аккаунтов одному пользователю. На текущий момент разработано обучение оценок алгебраической байесовской сети на основе бинарных данных [2]. Таким образом, чтобы сделать дальнейшее обучение возможным, результат сравнения каждого поля должен быть приведён к одному из следующих значений: 0 – несовпадение, 1 – совпадение или \* – отсутствие одного или обоих из сравниваемых значений. Для используемых методов сравнения, которые не дают бинарный результат, а дают число в промежутке (0,1) — применяются правила округления к ближайшему целому.

## Структура алгебраической байесовской сети

Для вычисления результата сравнения и получения информации о вероятности принадлежности двух аккаунтов одному человеку предложена следующая структура алгебраической байесовской сети, представляющей собой набор связанных фрагментов знаний, каждый из которых описывает небольшую часть информации, представляемой сетью. Математически фрагменты знаний представляют собой идеалы конъюнктов, дизъюнктов или квантов, при этом каждому из элементов сопоставляется интервал вероятности его истинности [2], что позволяет легко обрабатывать данные с пропусками, экспертные оценки и прочую неполную информацию.

Базовыми элементами, над которыми строится сеть, в частности, фрагменты знаний, являются атомарные высказывания. В контексте проводимого исследования этими элементами являются высказывания о совпадении значений атрибутов профиля, а также информация о принадлежности этих аккаунтов одному человеку ( $x$ ).

В контексте решаемой задачи была построена алгебраическая байесовская сеть со звездчатой структурой [5] с общим для всех фрагментов знаний элементом  $x$ . Каждый фрагмент знаний строился над близкими по смыслу атомарными утверждениями.

Разработанная в ходе исследования структура сети представлена на рисунке 1 (используются введенные в 1.1 обозначения). На рисунке 2 изображена более подробная структура части сети с фрагментами знаний в виде идеалов конъюнктов с высказываниями о поле, городах, образовании и карьере и совпадении адресов страницы (выделены светло-серым).

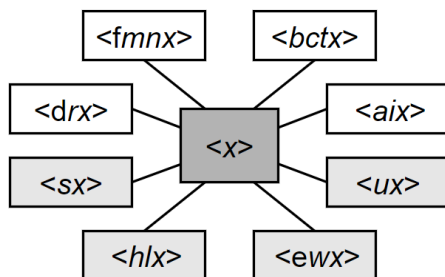


Рисунок 1. Структура алгебраической байесовской сети

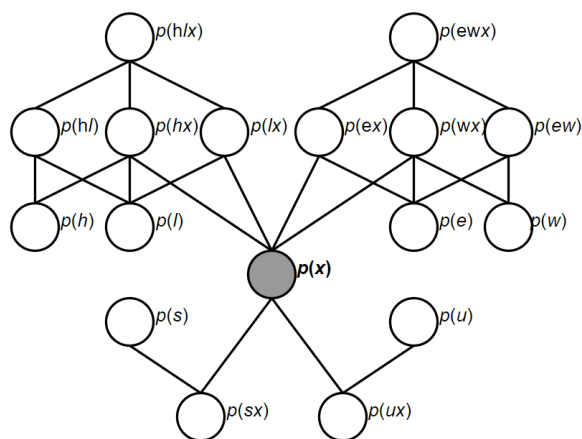


Рисунок 2. Подробная структура фрагмента алгебраической байесовской сети

После обучения сеть можно использовать для получения вероятности принадлежности двух профилей одному человеку.

На вход сети будет подаваться свидетельство, то есть набор информации о конкретных означиваниях некоторых переменных из представленных в сети. Далее на основе аппарата апостериорного вывода (2) может быть получена условная вероятность переменной  $x$ . Кроме того, произойдет уточнение оценок в сети. Например, на вход была подана информация о том, что имя, фамилия, родной город и образование в двух профилях совпадают, в то время как семейное положение и любимая музыка — нет. В таком случае формально свидетельство записывается в виде  $\langle n \wedge m \wedge h \wedge e \wedge \bar{f} \wedge \bar{t} \rangle$ , а результатом работы будет вероятность  $x$  при этих условиях:  $p(x | \langle n \wedge m \wedge h \wedge e \wedge \bar{f} \wedge \bar{t} \rangle)$ .

## Заключение

В статье предложена методика решения задачи сопоставления публичных анкет профилей в различных социальных сетях с помощью алгебраической байесовской сети. В качестве источника информации рассматриваются аккаунты социальных сетей «Одноклассники» и «ВКонтакте». Предложенная методика закладывает основу для дальнейшей работы в области автоматизации поиска, определения и слияния профилей пользователей в социальных сетях. Полученные результаты,

опосредованно, могут быть использованы с целью агрегации большого количества информации, необходимой при построении профиля уязвимостей пользователя.

Дальнейшими этапами исследования является обучение созданной структуры, то есть присвоение элементам сети интервалов вероятности их истинности на основе набора данных. Полученная таким образом сеть позволит по набору конкретных значений результатов попеременных сравнений двух профилей получать оценки вероятности их принадлежности одному пользователю.

### **Литература**

1. Корепанова А. А., Олисеенко В. Д., Абрамов М. В., Тулупьев А. Л. Применение методов машинного обучения в задаче идентификации аккаунтов пользователя в двух социальных сетях // Компьютерные инструменты в образовании. 2019. № 3. С. 29–44. doi:10.32603/2071-2340-2019-3-29-43
2. Тулупьев А.Л., Николенко С.И., Сироткин А.В. Основы теории байесовских сетей: учебник. СПб.: Изд-во С.-Петерб. ун-та, 2019. 399 с. (in Russian)
3. Zolotin A.A., Tulupyeve A.L. Sensitivity Statistical Estimates for Local A Posteriori Inference Matrix-Vector Equations in Algebraic Bayesian Networks over Quantum Propositions // Vestnik St. Petersburg University: Mathematics. Vol. 51(1), 2019. P. 42-48
4. Zolotin A.A., Tulupyeve A.L. Matrix-vector algorithms of global posteriori inference in algebraic Bayesian networks // Proceedings of 2017 20th IEEE International Conference on Soft Computing and Measurements (SCM 2017). 2017. P. 22-24
5. Kharitonov N., Malchevskaia E., Zolotin A., Abramov M. External consistency maintenance algorithm for chain and stellate structures of algebraic bayesian networks: Statistical experiments for running time analysis // Advances in Intelligent Systems and Computing. Vol. 875, 2019. P. 23-30
6. Jamali M.M., Mirzaei G. Bayesian Belief Network Based Occupancy Assessment Framework // Conference Record - Asilomar Conference on Signals, Systems and Computers. 2019. P. 792-796.
7. Baggenstoss P.M. On the Duality between Belief Networks and Feed-Forward Neural Networks // IEEE Transactions on Neural Networks and Learning Systems. Vol. 30(1). 2019. P. 190-200
8. Kou F., Du J., Yang C., Shi Y., Liang M., Xue Z., Li H. A multi-feature probabilistic graphical model for social network semantic search // Neurocomputing. Vol. 336. 2019. P. 67-78.