

# **ВЫДЕЛЕНИЕ ГРУППЫ РАДИКАЛЬНЫХ ФУТБОЛЬНЫХ БОЛЕЛЬЩИКОВ В СОЦИАЛЬНЫХ МЕДИА ПО ИХ ИНТЕРЕСАМ И ПОВЕДЕНИЮ НА ПРИМЕРЕ СЕТИ VK.COM**

Сергей Сергеевич Дмитриев, магистрант Университета ИТМО  
seregas.dmitrias@gmail.com

## **Аннотация**

В данной работе представлен подход к определению принадлежности пользователей к группе радикальных футбольных болельщиков в социальной сети «ВКонтакте».

## **Введение**

Последнее время социальные медиа набрали огромную популярность. Такие сайты как *Facebook.com*<sup>1</sup>, *Vk.com*<sup>2</sup>, *Twitter.com*<sup>3</sup> обладают аудиторией из многих миллионов пользователей. Они создают огромные массы контента, состоящего из их мнений и точек зрения. Однако содержание этой информации в основном остается не использованным. Тогда как оно может быть крайне важным.

Такие данные могут быть использованы для определения интересов, предпочтений и иных личных свойств пользователя. В большинстве социальных медиа профиль пользователя содержит информацию о его поле, возрасте, местоположении, интересах. Однако, зачастую, такие данные могут быть неполными, а иногда и неверными. А некоторые признаки, такие как, например, вероисповедание или политические взгляды обычно опускаются. Из-за этой неполноты возникает задача восстановления недостающих данных. Получение подобной информации может, к примеру, позволить отлавливать на ранней стадии, бандитские группировки, или прогнозировать конфликты между группами людей различных взглядов [1, 2].

Существуют исследования, показывающие возможность восстанавливать информацию, явно неуказанную в профилях пользователей [3, 4, 9], исследования, показывающие влияние социальных связей на поведение пользователей [5, 6]. В исследовании [3] информация восстанавливалась на основе социальных связей. В работе вычислялись

---

1 <http://www.facebook.com>

2 <http://www.vk.com>

3 <http://www.twitter.com>

политические предпочтения пользователей. Исследователи предложили статистическую модель, в которой, вероятность связи между двумя пользователями зависит от меры расстояния между ними в плоскости идеологии.

Для получения точек на плоскости политических предпочтений использовались «корневые» пользователи. Они обладали, как минимум, 10 подписками на политические аккаунты, для которых уже имелась информация о их предпочтениях. Этим пользователям разметили их политические предпочтения, исходя из подписок, и разместили на идеологической плоскости. Далее подписчиков корневых пользователей разместили по аналогичному признаку.

В данной работе предложены подходы для определения принадлежности пользователя к группе футбольных фанатов, а так же выделение подгруппы радикальных футбольных фанатов. Предполагается, что описанные методы могут позволить выделять возможных радикалов, а так же в последствии предсказывать динамику изменения взглядов выбранного пользователя. Аналогичный подход можно обобщить и применять для выделения других групп пользователей.

## **Описание набора данных**

Сайт *Vk.com* поддерживает открытое API. Через него были собраны данные. Использовались следующие методы: *groups.getMembers*, *board.getTopics*, *board.getComments*, *likes.isLiked*. Первоначально было вручную выбрано 60 публичных страниц на тему футбола. Для каждого из них вручную была определена характеристика их радикальности. В статье [8] показано, что радикалы используют типичные для себя фразы. Так был создан словарь потенциально радикальных фраз, считались вхождения этих фраз в постах групп. Так же многие выбранные группы были заблокированы по причине своего радикального содержания. Поэтому при определении радикальности использовались две характеристики: число вхождений радикальных фраз и блокировка группы в России. Данные подобраны так, что радикальных публичных страниц было столько же сколько и нерадикальных. Далее для каждой группы выкачивались все её посты. Так же сохранялась информация о всех лайках и репостах постов групп, информация обо всех подписчиках, друзьях подписчиков и подписок подписчиков.

## Преобразование данных

В данном разделе описаны преобразования собранных данных, а так же генерация новых.

### *Преобразование текстовых данных*

Так как зачастую в текстах групп в социальных сетях присутствуют спец. символы, символы, относящиеся к другому языку, полученные данные были очищены от символов не принадлежащих к латинице и кириллице. Так же для стемминга слов использовался Mystem<sup>4</sup>.

### *Преобразование данных в векторное пространство*

Набор пользователей и групп можно представить как циклический ориентированный граф. Пользователи и группы будут являться узлами графами. Ребра между ними соответствуют отношению нахождению в подписках. Собранные группы будут делиться на 2 типа: группы, в которых вручную размечен уровень радикальности, и новые группы собранные автоматически. Для первых групп вектор радикальности:

$$V_{group1} = (R, \frac{K}{N}) ,$$

где  $R$  вручную выставленная мера радикальности группы, принимает значение 0 или 1,  $K$  число фраз из сформированного словаря радикальных выражений входящих в посты, а  $N$  число постов группы. У подписчиков этих групп, назовем их первым поколением пользователей, вектор радикальности

$$V_{user1} = \left( \frac{\sum_1^n \sqrt{R_i^2 + \left(\frac{K_i}{N_i}\right)^2}}{n}, \frac{\sum_1^p (likes_i + reposts_i)}{p} \right) ,$$

где  $n$  число первых групп, а  $p$  число постов в этих группах.

## Вычисление радикальности для новых пользователей

Воспользуемся предположением, что обладая большим числом связей с определенной группой людей, одобряя их действия, вероятность

4 <https://tech.yandex.ru/mystem/>

принадлежать к этой к группе у человека возрастает. Таким образом, меры радикальности новых групп должны расти тем больше, чем больше на них подписано радикалов. Предположим, что наличие агрессивных фраз так же должно влиять на усиление радикальности. Таким образом для новых групп считаем вектор радикальности, основанный на описанных признаках

$$V_{group_{next}} = \left( \frac{\sum_1^u V_{user1}}{n}, Count(aggressive\ phrases) \right),$$

где n число всех пользователей первого поколения, для пользователей второго поколения вектор радикальности считался следующим образом.

$$V_{user_{next}} = \left( \frac{\sum_1^u V_{user1}}{n}, \frac{\sum_1^p (likes_i + reposts_i)}{p} \right)$$

Полученные признаки образуют двумерное пространство.

## Результат

Так как достоверных сведений о реальных радикальных футбольных фанатах в социальных сетях нет, для оценки результатов используется общедоступная статистика. В статье [7] показана связь геолокации пользователя и его социальных связей. Предлагается воспользоваться этой связью и показать, что люди из регионов где совершается большее число радикальных действий, будут иметь более радикальные признаки.

Так например, известно, что в Санкт-Петербурге ежегодно регистрируются случаи нарушения закона радикальными футбольными болельщиками<sup>5</sup>, а в городе Краснодар подобных происшествий не зарегистрировано. Таким образом можно предположить, что если модель верна, то средние меры радикальности у пользователей из двух городов будут отличаться.

---

5 <http://vz.ru/news/2014/5/11/686166.html>

Город	Ср. мера основанная на графе соц. связей	Ср. мера о лайках репостах
Санкт-Петербург	0,017	0,022
Краснодар	0,008	0,019

Таблица 1: Средние меры радикальности регионов

Результаты применения модели показывают, что критерий основанный на социальных связях показывает большое отличие между регионами, однако критерий основанный на активности пользователей, их лайках и репостах, не демонстрирует сильного отличия. Возможная причина слабого отличия этих мер – неточное выделение радикальных постов.

Данный механизм выделения группы пользователей можно обобщить и применять для других видов групп.

## Заключение

В рамках представленной работы были показаны подходы к выделению группы радикальных футбольных фанатов из социальной сети *Vk.com*. Они могут быть применены к другим группам пользователей.

В дальнейшем планируется опробовать для решения данной задачи модель случайных марковских полей. Для улучшения результатов определения радикальности постов в социальной сети планируется воспользоваться латентным семантическим анализом. Так же предлагается использовать данные твиттера и сопоставить их с данными пользователей из вконтакте. Это предполагает решение задачи о подтверждении их подлинности [10], вычисляя аналогичные признаки вероятно можно будет улучшить результаты.

## Литература

1. Doyle A. et al. Forecasting significant societal events using the Embers streaming predictive analytics system //Big data. – 2014. – Т. 2. – №. 4. – С. 185-195.
2. <http://arstechnica.co.uk/security/2016/02/the-nsas-skynet-program-may-be-killing-thousands-of-innocent-people/>
3. Barberá P. et al. Tweeting From Left to Right Is Online Political Communication More Than an Echo Chamber? //Psychological science. – 2015. – С. 0956797615594620.
4. Zheleva E., Getoor L. To join or not to join: the illusion of privacy in social networks with mixed public and private user profiles //Proceedings of the 18th international conference on World wide web. – ACM, 2009. – С. 531-540.
5. Trusov M., Bodapati A. V., Bucklin R. E. Determining influential users in internet social networks //Journal of Marketing Research. – 2010. – Т. 47. – №. 4. – С. 643-658.
6. Bond R. M. et al. A 61-million-person experiment in social influence and political mobilization //Nature. – 2012. – Т. 489. – №. 7415. – С. 295-298.
7. Wang D. et al. Human mobility, social ties, and link prediction //Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining. – ACM, 2011. – С. 1100-1108.
8. Wijzen F. ‘There are radical Muslims and normal Muslims’: an analysis of the discourse on Islamic extremism //Religion. – 2013. – Т. 43. – №. 1. – С. 70-88.
9. Bergsma S. et al. Broadly Improving User Classification via Communication-Based Name and Location Clustering on Twitter //HLT-NAACL. – 2013. – С. 1010-1019.
10. Peled O. et al. Matching Entities Across Online Social Networks //arXiv preprint arXiv:1410.6717. – 2014.