

УДК 004.422

КЛАССИФИКАЦИЯ ТЕСТОВЫХ СООБЩЕНИЙ, РАЗМЕЩЕННЫХ НА ВЕБ-САЙТАХ СОЦИАЛЬНЫХ СЕТЕЙ, ПО ТИПУ ПРОБЛЕМЫ ДЛЯ ПРИМЕНЕНИЯ В СОС БАНКОВСКОЙ ОРГАНИЗАЦИИ

Сергеев Д.Н. (Университет ИТМО)

Научный руководитель – к.т.н., доцент Воробьева А.А.
(Университет ИТМО)

Аннотация

Работа посвящена вопросам машинного обучения в области классификации текстовых сообщений, размещенных на веб-сайтах социальных сетей, по типу проблемы для применения в СОС банковской организации. Представлены результаты анализа современной научно-технической, нормативной, методической литературы по теме выявления коротких текстовых сообщений о заданной банковской организации, содержащих информацию о технических сбоях и банковском мошенничестве, а также экспериментальным исследованиям.

Введение.

При современном развитии информационных технологий и распространенности коммуникационных сетей, почти у каждой крупной банковской компании существует страница, размещенная на веб-сайтах социальных сетей. Зачастую клиенты банка сталкиваются с техническими сбоями при обслуживании, банковским мошенничеством и т.п. Пользователи сети Интернет имеют возможность сообщить о данных проблемах на веб-ресурсах социальных сетей. Одна из основных задач банка – способность решить проблему клиентов в момент обращения. При этом самым динамично развивающимся каналом для общения клиентов являются социальные сети. Сотрудниками Центра по управлению киберинцидентами (security operation center, SOC) банковской организации ведутся работы по мониторингу Интернет для выявления в социальных сетях и на тематических порталах сообщений о технических сбоях и банковском мошенничестве. Но не всегда администраторам групп удается своевременно обработать данные сообщения. Подобные задержки в обработке сообщений пользователей о сбоях и мошенничестве могут повлечь не только финансовые потери банковской организации, но и снижение ее репутации. Сотрудникам приходится уделить немало времени, чтобы отреагировать на подобные посты пользователей, ведь помимо сообщений об ошибках и банковском мошенничестве, немало комментариев, которые не содержат необходимую информацию. Потенциальными потребителями результатов данного исследования обозначена потребность автоматизации процесса выявления коротких текстовых сообщений, размещенных на веб-сайтах социальных сетей и содержащих информацию о технических сбоях и банковском мошенничестве. Однако, решение данной задачи, путем применения методов автоматического анализа текста, поможет сократить не только временные затраты, но и финансовые потери, обусловленные снижением репутационного ущерба. Все успешно размеченные сообщения пользователей будут перенаправлены сотрудникам, для дальнейшей работы с этими данными

Основная часть.

Целью работы является формирование этапов разработки экспериментального образца, основанного на методах машинного обучения, с целью выявления текстовых сообщений, размещенных на веб-сайтах и социальных сетях Интернета для определения тематики сообщений о банковской организации для применения в СОС банковской организации.

В соответствии с заявленными целями и задачами работы объектом исследования является выявление перспективных методов для определения вида проблемы в текстовых сообщениях, относящихся к банковской организации, а предметом исследования – алгоритм определения типа проблемы в сообщениях о банковской организации.

Методы определения тематики сообщения применяются для первичной группировки сообщений к одной из трех категорий: технический сбой, банковское мошенничество, иная тематика. Применение методов классификации к текстам часто требует предварительной их очистки, предобработки (стемминг, лемматизация и пр.) и векторизации (преобразование текста в векторное представление). Выбор метода предобработки текстовых сообщений и алгоритмов машинного обучения для их классификации является непростой задачей.

В большинстве проанализированных работ использовались следующие алгоритмы предобработки текста и машинного обучения для дальнейшей классификации: метод опорных векторов, наивный байесовский классификатор, метод k-ближайших соседей, деревья решений, линейная регрессия, сверточная нейронная сеть.

По результатам проведенного исследования были рассмотрены различные подходы к классификации текста. Был разработан алгоритм классификации текстовых сообщений для определения типа проблемы на основе нейронной сети. Перед проведением экспериментальных исследований тексты были очищены от цифр и других символов, не являющимися буквами русского алфавита. Далее данные были разделены на тренировочную, тестовую и валидационную части.

Выводы.

В ходе экспериментальных исследований проведено сравнение различные методы машинного обучения, которые были реализованы в рассматриваемых литературных источниках. Задачей экспериментальных исследований являлось установление зависимости показателей результирующей точности классификации сообщений, относящихся к банковской организации и содержащих информацию о банковском мошенничестве, сбое банковский сервисов или же отсутствие определенных проблем по всем классам. Использовались следующие методы методов машинного обучения: стохастический градиентный спуск, случайный лес, логистическая регрессия, дерево принятия решений, метод k-ближайших соседей, LSTM-нейронная сеть. Результаты экспериментов показали, что наиболее высокие результаты определения тематики сообщения достигаются с использованием нейросети.

Сергеев Д.Н.

Воробьева А.А.