

Выполнил(а) Долматов Д.А., № группы К3221, дата 20.11.2021, оценка \_\_\_\_\_  
ФИО студента не заполнять

<b>Название статьи/главы книги:</b> Классификация медиатекстов с использованием машинного обучения		
<b>ФИО автора статьи:</b> Л.В. Мотовских	<b>Дата публикации:</b> 2020	<b>Размер статьи</b> 7 стр.
<b>Прямая полная ссылка на источник и сокращенная ссылка:</b> <a href="https://cyberleninka.ru/article/n/klassifikatsiya-mediatekstov-s-ispolzovaniem-mashinnogo-obucheniya/viewer">https://cyberleninka.ru/article/n/klassifikatsiya-mediatekstov-s-ispolzovaniem-mashinnogo-obucheniya/viewer</a> <a href="https://clck.ru/Yw8Bc">https://clck.ru/Yw8Bc</a>		
<b>Тэги, ключевые слова или словосочетания:</b> классификация текстов; TF-IDF; «случайный лес»; метод опорных векторов; электронные СМИ.		
<b>Перечень фактов, упомянутых в статье:</b> Методы машинного обучения намного эффективнее оптимизируют работу с большими данными по сравнению с классическими способами алгоритмирования способов получения и обработки информации. Одной из сфер применения алгоритмов нейронных сетей применяется в классифицировании информации в СМИ. В статье описывается два алгоритма, а источником медиатекстов являются внутренние СМИ РФ, которые имеют подразделение на: общество, экономику, политику, спорт и культуру. Сначала собирается выборка, при этом отбрасываются 5% как самых длинных статей, так и 5% самых коротких по количеству символов. Итоговая выборка подразделяется на две в соотношении 85:15 на обучающую и проверочную. Текстовый формат информации представляется в терм-документной таблице - матрицы TF-IDF, после чего применяется алгоритм «Random forest», который случайно распределяет элементы подвыборки каждой матрицы (дереву) и присваивает им случайный набор элементов с количеством деревьев, узлов, вершин, глубин. Точность таких результатов равна 86%. Второй алгоритм «Метод опорных векторов» разделяет подвыборку на несколько гиперплоскостей, которые находятся на максимально возможном расстоянии друг от друга. Предсказание осуществляется через положение элемента относительно оптимальной гиперплоскости. Для присваивания первичного набора значений используется метод случайного поиска, а затем поиск по сетке параметров. Результаты классификации выборки медиатекста оказались равны 88%, что оправдало ожидания исследователей.		
<b>Позитивные следствия и/или достоинства описанной в статье технологии</b> <ol style="list-style-type: none"> <li>1) Поиск по сетке параметров явно превосходит быстрый поиск по точности.</li> <li>2) Точное предсказывание модели классификатора для основных рубрик СМИ.</li> <li>3) Возможность классификации большего количества неразмеченных текстов других СМИ.</li> </ol>		
<b>Негативные следствия и/или недостатки описанной в статье технологии</b> <ol style="list-style-type: none"> <li>1) Наибольшая погрешность классификатора между категориями «Общество» и «Спорт».</li> <li>2) Необходимость больших вычислительных мощностей для построения групп TF-IDF.</li> <li>3) Отсутствие стратегического значения данной технологии (кроме сортировки писем СПАМа).</li> </ol>		
<b>Ваши замечания, пожелания преподавателю или анекдот о программистах</b> Such a great day today!		