

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ  
”ВЫСШАЯ ШКОЛА ЭКОНОМИКИ”»

**Московский институт электроники и математики**

Юткин Дмитрий Игоревич, группа БИВ-144  
Вдовкин Василий Алексеевич, группа БИВ-144

**Классификация и агрегация новостных статей, используя методы и модели  
обработки естественного языка**

Междисциплинарная курсовая работа  
по направлению 09.03.01 Информатика и вычислительная техника  
студентов образовательной программы бакалавриата  
«Информатика и вычислительная техника»

Студент \_\_\_\_\_ Д.И. Юткин  
Студент \_\_\_\_\_ В.А. Вдовкин

Руководитель  
старший преподаватель  
Л.Л. Волкова  
\_\_\_\_\_

Москва 2017 г.

## **Аннотация**

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

## **Abstract**

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

## Оглавление

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Введение</b>                                       | <b>4</b> |
| <b>2</b> | <b>Цель и задачи курсовой работы</b>                  | <b>4</b> |
| <b>3</b> | <b>Сбор и подготовка данных</b>                       | <b>5</b> |
| 3.1      | Получение данных с новостных источников               | 5        |
| 3.2      | Преоброессинг новостных статей                        | 5        |
| <b>4</b> | <b>Классификация новостных статей</b>                 | <b>5</b> |
| 4.1      | Линейная модель на TF-IDF признаках                   | 5        |
| 4.2      | Градиентный бустинг деревьев на word2vec              | 5        |
| <b>5</b> | <b>Агрегация новостных статей</b>                     | <b>5</b> |
| 5.1      | Кластеризация с помощью алгоритмов машинного обучения | 5        |
| 5.2      | Объединение в связанные компоненты графа              | 5        |
| <b>6</b> | <b>Разработка веб-приложения</b>                      | <b>5</b> |
| 6.1      | Back-end часть  | 5        |
| 6.2      | Front-end часть                                       | 5        |
| <b>7</b> | <b>Заключение</b>                                     | <b>5</b> |
|          | <b>Список литературы</b>                              | <b>6</b> |
|          | <b>Приложение</b>                                     | <b>7</b> |

## 1 Введение

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

## 2 Цель и задачи курсовой работы

Целью курсовой работы является разработка веб-сервиса, который

- а) в реальном времени получает статьи из новостных источников;
- б) классифицирует полученные статьи по общим темам;
- в) агрегирует по схожести содержания статьи из различных источников.

Последние два пункта должны быть автоматизированы с помощью алгоритмов анализа данных и машинного обучения.

Для достижения поставленной цели, должны быть выполнены следующие задачи:

- а) Изучить подходы автоматической обработки текста и естественного языка;
- б) Собрать набор данных для обучения моделей машинного обучения;
- в) Исследовать различные подходы к классификации новостных статей;
- г) Исследовать различные подходы к агрегации новостных статей;
- д) Разработать веб-сервис.

### **3 Сбор и подготовка данных**

#### **3.1 Получение данных с новостных источников**

#### **3.2 Препроцессинг новостных статей**

### **4 Классификация новостных статей**

#### **4.1 Линейная модель на TF-IDF признаках**

#### **4.2 Градиентный бустинг деревьев на word2vec**

### **5 Агрегация новостных статей**

#### **5.1 Кластеризация с помощью алгоритмов машинного обучения**

#### **5.2 Объединение в связные компоненты графа**

### **6 Разработка веб-приложения**

#### **6.1 Back-end часть**

#### **6.2 Front-end часть**

### **7 Заключение**

## **Список литературы**

## Приложение