

Лабораторная работа 0

Получение и обработка данных

Постановка задачи

Требуется сформировать/получить два набора данных соответствующие следующим критериям:

- 1) Один из датасетов должен представлять собой корпус документов. Язык, источник и тематика произвольна
- 2) Второй датасет должен содержать категориальные, количественные признаки. Для данного датасета определить предсказываемые признаки (для задачи регрессии и классификации). Если такого признака нет, спроектировать

Данные датасеты будут в дальнейшем использованы в оставшихся лабораторных работах.

По каждому датасету построить распределения признаков (в случае корпуса документов – построить распределение слов) и объяснить имеющуюся картину. Вычислить статистические характеристики признаков. Обнаружить и решить возможные проблемы с данными. Если решить данную проблему невозможно, объяснить почему.

Требования

- 1) Датасеты должны быть уникальны
- 2) Исходный код должен быть написан в одном код стайле
- 3) Должен быть указан источник данных

Дата выдачи

16.03.2019

Срок выполнения

2 недели