

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ
РОССИЙСКОЙ ФЕДЕРАЦИИ
Федеральное государственное бюджетное образовательное
учреждение высшего профессионального образования

САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ

Математико-механический факультет
Кафедра прикладной кибернетики

Путников Семен Андреевич

КУРСОВАЯ РАБОТА

Алгоритмы кластеризации: выявление
профилей пользователей по анализу
активности

Руководитель курсовой работы
_____ Н. В. Кузнецов
«___» _____ 2019 г.

Санкт-Петербург, 2019 г.

Содержание

1	Введение	2
2	Постановка задачи	3
3	Теоретический аспект	4
3.1	Метрики	4
3.2	Feature engineering	5
3.3	Класстеризация	5
4	Практический аспект	8
5	Итог	9

1 Введение

Мы живем в постиндустриальном обществе, где как известно: "кто владеет информацией тот владеет миром". Люди стараются собирать сведения обо всем: погоде, экономике, политике, демографии и во многих других отраслях. Каждая компания старается собрать как можно больше данных, надеясь занять лидирующую позицию. Однако сама по себе информация не представляет особой ценности, это своего рода сырье, которое еще нужно обработать. Только после анализа информация представляет для человека или компании определенную ценность, стоит отметить, что чем глубже и точнее произведен анализ информации, тем ценнее его результат.

В рамках анализа и обработки информации появился термин Data Mining, который означает: добыча данных, интеллектуальный анализ данных, глубинный анализ данных. Достаточно часто Data Mining связывают с Business intelligence, которые в свою очередь широко применяются в бизнесе.

Задачи, решаемые Data Mining:

1. Классификация — отнесение входного вектора (объекта, события, наблюдения) к одному из заранее известных классов.
2. Кластеризация — разделение множества входных векторов на группы (кластеры) по степени «похожести» друг на друга.
3. Сокращение описания — для визуализации данных, упрощения счета и интерпретации, сжатия объемов собираемой и хранимой информации.
4. Ассоциация — поиск повторяющихся образцов. Например, поиск «устойчивых связей в корзине покупателя».
5. Прогнозирование — нахождение будущих состояний объекта на основании предыдущих состояний (исторических данных)
6. Анализ отклонений — например, выявление нетипичной сетевой активности позволяет обнаружить вредоносные программы.
7. Визуализация данных.

Рассмотрим задачу кластеризации. Она является одной из важнейших задач Data Mining. Благодаря ей существует возможность прогнозировать, а так же получать детальную информацию об уже имеющихся кластерах. Например, в рамках государства это означает, что администрация может оперативно и точно получать информацию о социальных слоях общества, сегментах бизнеса, категориях годности граждан и др., это позволяет аппарату управления быстро реагировать на изменения, а так же задавать вектор развития страны.

2 Постановка задачи

В данной работе предлагается рассмотреть следующую задачу - анализ имеющихся алгоритмов кластеризации и создание пайплайна для построения кластеризации по массиву данных. Этот массив представляет собой логи REST-запросов пользователей на портале учета древесины и сделок с ней. Проведем декомпозицию задачи и выделим следующие этапы.

- Анализ существующих алгоритмов кластеризации
- Исследование инструментов Elasticsearch
- Обработка массива данных
- Реализация алгоритма кластеризации для поставленной задачи
- Анализ полученных результатов

3 Теоретический аспект

3.1 Метрики

Для работы алгоритмов кластеризации необходимо ввести метрику на данных.

Определение 1. $\rho(x, y) : X \times X \rightarrow \mathbb{R}$ называют метрикой, если выполняются следующие условия:

1. $\rho(x, y) = 0 \Leftrightarrow x = y$ (аксиома тождества)
2. $\rho(x, y) = \rho(y, x)$ (аксиома симметрии)
3. $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ (неравенство треугольника)

Приведем наиболее распространенные в Data Mining метрики расстояний.

1. Евклидова метрика

$$\rho(x, y) = \sqrt{\sum_i^n (x_i - y_i)^2} \quad (1)$$

2. Квадрат евклидова расстояния

$$\rho(x, y) = \sum_i^n (x_i - y_i)^2 \quad (2)$$

3. Манхэттенская метрика

$$\rho(x, y) = \sum_i^n |x_i - y_i| \quad (3)$$

4. Расстояние Чебышева

$$\rho(x, y) = \max_i |x_i - y_i| \quad (4)$$

5. Степенное расстояние

$$\rho(x, y) = \sqrt[r]{\sum_i^n (x_i - y_i)^p} \quad (5)$$

6. Дискретная метрика

$$\rho(x, y) = \begin{cases} 0, & x = y, \\ 1, & x \neq y. \end{cases} \quad (6)$$

3.2 Feature engineering

Первый этап процесса кластеризации - отбор признаков. Они бывают различных типов: бинарные, вещественные, категориальные, текстовые и другие. В исследуемых нами данных признаки текстовые и категориальные. К категориальным относят признаки, значение которых нельзя сравнивать между собой, можно лишь проверять их равенство. Текстовый признак представляет из себя последовательность слов. Рассмотрим наиболее распространенные алгоритмы их отбора.

- **One-hot encoding или "Мешок слов".** В рамках этого метода строится словарь всех уникальных слов в датасете, а потом каждому слову приобретает уникальный индекс. Тогда каждое предложение можно будет отобразить списком, длина которого равна числу уникальных слов в словаре, а в каждом индексе в этом списке будет храниться, сколько раз данное слово встречается в предложении.

	MARY	IS	HUNGRY	HAPPY	FOR	APPLES	NOT	JOHN	HE	
"Mary is hungry for apples."	1	1	1	0	1	1	0	0	0	→ [1, 1, 1, 0, 1, 1, 0, 0, 0]
"John is happy he is not hungry for apples."	0	2	1	1	1	1	1	1	1	→ [0, 2, 1, 1, 1, 1, 1, 1, 1]

Рис. 1: пример работы one-hot encoding.

- **TF-IDF.** В этом методе оценивается не количество вхождений слов, а их оценки важности для текста. Руководствуемся двумя принципами: чем чаще слово встречается в документе, тем оно важнее, и, чем реже слово встречается в остальных документах, тем оно важнее. Рассмотрим формулу на их основе.

n_{iw} (term frequency) - число вхождений слова w в текст x_i^j

N_w (document frequency) - число текстов, содержащих w

Тогда важность слова w для документа x_i^j исчисляется по формуле:

$$TF - IDF(i, w) = n_{dw} \log(l/N_w)$$

3.3 Кластеризация

Кластеризация - это отдельный класс задач машинного обучения, отличающийся от классификации тем, что у объектов обучающей выборки нет заранее заданных ответов учителя. Задача состоит в выделении отдельных кластеров, состоящих из близких объектов так, чтобы объекты разных кластеров существенно различались. Рассмотрим основные алгоритмы кластеризации.

- **Метод k-средних.** На входе этого метода некоторым образом расставляются k первоначальных центров кластеров. Потом для каждого объекта

вычислится его близость к каждому из кластеров. Затем вычисляется новый центр кластера. И так повторяется пока центры кластеров не стабилизируются.

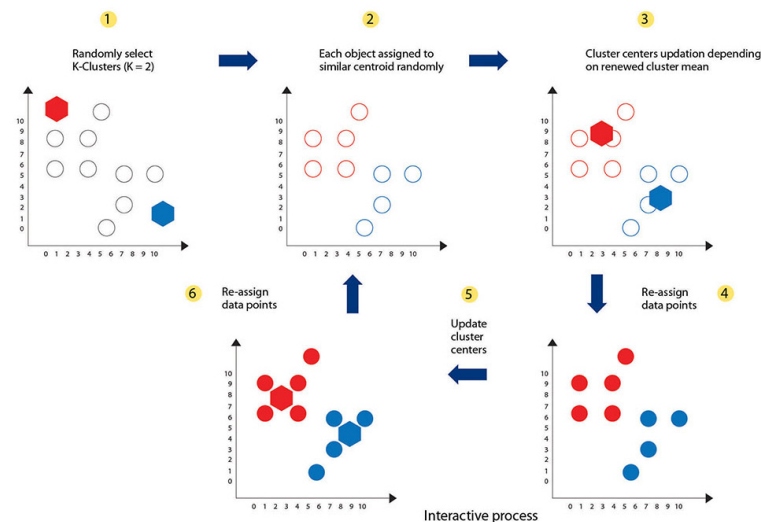


Рис. 2: пример работы метода k-средних.

- **Алгоритм Ланса-Уильямса.** Этот алгоритм относится к классу агломеративных методов иерархической кластеризации. Агломеративность означает то, что метод объединяет мелкие кластеры в более крупные и отображают историю этих объединений в виде дендрограммы.

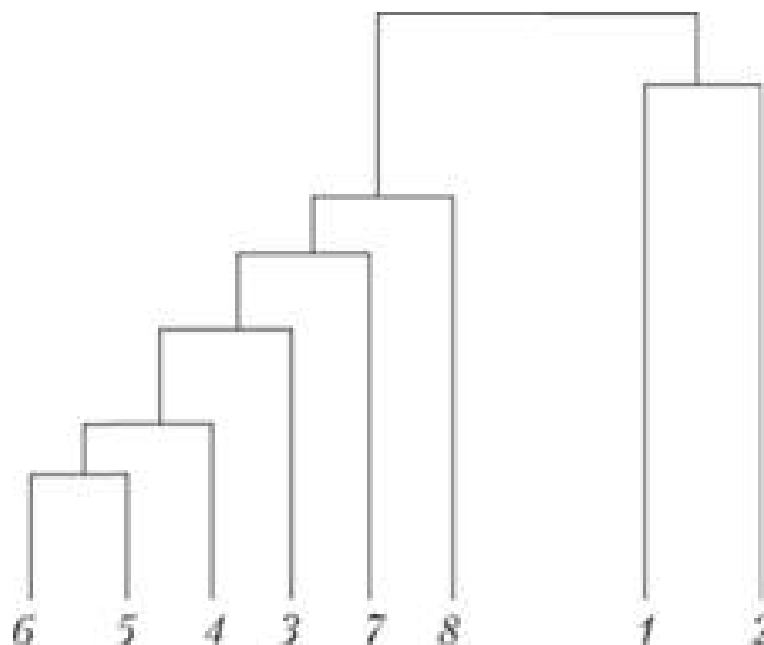


Рис. 3: пример дендрограммы.

Алгоритм основан на том, что сначала мы берём все кластеры одноэлементные, то есть каждый кластер соответствует какой-то одной точке, и кластеров ровно столько, сколько объектов обучающей выборки. А дальше мы находим в текущем множестве кластеров два самых ближайших и

сливаем их в один кластер. Повторяем это ровно до того момента, пока все кластеры не объединятся в один.

Отдельно стоит отметить формулу Ланса-Ульямса, которая используется для нахождения расстояния между кластером W , полученным из объединения U и V , и кластера S , зная расстояния между исходными кластерами.

$$R(U \cup V, S) = \alpha_U R(U, S) + \alpha_V R(V, S) + \beta R(U, V) + \gamma |R(U, S) - R(V, S)|$$

Частные случаи формулы:

1. Расстояние ближайшего соседа:

$$R(W, S) = \min_{w \in W, s \in S} \rho(w, s);$$

где $\alpha_U = \alpha_V = \frac{1}{2}$, $\beta = 0$, $\gamma = -\frac{1}{2}$

2. Расстояние дальнего соседа:

$$R(W, S) = \max_{w \in W, s \in S} \rho(w, s);$$

где $\alpha_U = \alpha_V = \frac{1}{2}$, $\beta = 0$, $\gamma = \frac{1}{2}$

3. Расстояние между центрами:

$$R(W, S) = \rho^2\left(\sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|}\right);$$

где $\alpha_U = \frac{|U|}{|W|}$, $\alpha_V = \frac{|V|}{|W|}$, $\beta = -\alpha_U \alpha_V$, $\gamma = 0$

4. Групповое среднее расстояние:

$$R(W, S) = \frac{1}{|W||S|} \sum_{w \in W} \sum_{s \in S} \rho(w, s);$$

где $\alpha_U = \frac{|U|}{|W|}$, $\alpha_V = \frac{|V|}{|W|}$, $\beta = \gamma = 0$

4 Практический аспект

Во время поиска практического решения задачи. Появился ряд проблем связанных с размером массива данных. Входные данные представляют собой JSON-файл весом более 100 гигабайт. На персональных машинах такие объемы не обработать, поэтому был создан пайплайн сборки датасета по кусочкам.

Первым делом командой в терминале разбиваем исходный файл на кусочки по 10000 строк и перенесем их в отдельную папку.

```
mkdir parts && cd parts && split -d -l 10000 --additional-suffix=.  
json ../data.json ex && cd ..
```

Дальше рассмотрим скрипт, который обработает файлы по кусочкам и соберет в один датасет. Обходя дерево json, из первой строки получаем название всех признаков.

```
def getIndexList(dataJson):  
    indexSet = set()  
    for i in dataJson:  
        for k in i.keys():  
            if(k != '_source'): indexSet.add(k)  
        for j in i['_source'].keys():  
            indexSet.add(j)  
    indexList = list(indexSet)  
    return indexList
```

По полученному списку создаем пустой Data Frame. Следующим шагом мы обходим по строчке файл и записываем данные в frame.

```
def writeData(indexList,dataJson,i):  
    data = getDataFrame()  
    for n in atpbar(range(len(dataJson)), name='LOCAL {}'.format(i)):  
        pattern = defaultdict(list)  
        for g in indexList:  
            pattern[g].append(None)  
        for k in dataJson[n].keys():  
            if(k != '_source'): pattern[k] = dataJson[n][k]  
        for j in dataJson[n]['_source'].keys(): pattern[j] = dataJson[n][  
            |['_source'][j]  
        data = data.append(pd.DataFrame.from_dict(pattern),  
            ignore_index=True, sort = False)  
    return data
```

Как только этот json файл закончится, переходим к следующему файлу и повторяем предыдущие шаги.

5 Итог

Подводя итог, можно сказать, что цели проекта были частично достигнуты. Была исследованы существующие методы отбора признаков и кластеризации данных, а также их дальнейшей валидации. В результате практической работы был создан скрипт обработки и подготовки данных для применения кластеризации.

В процессе работы над проектом были преодолены ряд сложностей связанных с размером входных данных. Изучены методы для обработки таких массивов.

Проект может иметь продолжение с целями оптимизации алгоритмов кластеризации, а также внедрения данного решения в существующие системы тестирования веб-портала по учету древесины и сделок с ней.

Список литературы

[1] Воронцов К.В. Методы кластеризации –

[http : //shad.yandex.ru/lectures/machinelearning.xml](http://shad.yandex.ru/lectures/machinelearning.xml)

[2] Барсегян А. А. и др. Технологии анализа данных: Data Mining, Visual Mining, Text Mining, OLAP. 2-е изд., перераб. и доп. — СПб.: БХВ-Петербург, 2007.

[3] Вагин В. Н., Головина Е. Ю., Загорянская А. А. Достоверный и правдоподобный вывод в интеллектуальных системах. — М.: Физматлит. 2004.

[4] Кацко И. А., Human Н. Б. Практикум по анализу данных на компьютере. — М.: КолосС, 2009.