

**ПРАВИТЕЛЬСТВО РОССИЙСКОЙ ФЕДЕРАЦИИ
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»**

Факультет компьютерных наук
Образовательная программа «Прикладная математика и информатика»

Отчет о программном проекте

на тему _____

(промежуточный, этап 1)

Выполнил:

студент группы БПМИ18_____

Подпись

И.О. Фамилия

Дата

Принял:

руководитель проекта _____

Имя, Отчество, Фамилия

Должность

Место работы

Дата _____ 2020

Оценка (по 10-тибалльной шкале)

Подпись

Москва 2020

Содержание

1	Задачи проекта	3
2	Алгоритм Affinity Propagation	3
2.1	Использование алгоритма	3
2.2	Шаги алгоритма	3
2.3	Similarity Matrix	4
2.4	Responsibility Matrix	4
2.5	Availability Matrix	4
2.6	Оптимизации	4
3	Алгоритм PrePost	6
3.1	Постановка задачи	6
3.2	Основные понятия	6
3.3	Шаги алгоритма	6
3.4	FP-tree	7
3.5	PPC-Tree	7
3.6	NL-1	8
3.7	NL-k	8
3.8	Способ нахождения NL-k	8
4	Источники	10

1 Задачи проекта

Изучение и сравнение алгоритмов кластеризации и алгоритмов поиска ассоциативных правил. Далее разобраны алгоритмы Affinity Propagation (алгоритм кластеризации), PrePost (алгоритм поиска ассоциативных правил).

2 Алгоритм Affinity Propagation

2.1 Использование алгоритма

- Одной из главных особенностей алгоритма является то, что он не требует указывать количество кластеров на которые надо разбить данные.
- Также у каждого кластера в результате алгоритма есть представитель, являющийся точкой из изначальных данных.
- Датасет не очень большой ($N \leq 10^6$) или в меру большой, но разреженный ($N \leq 10^7$)
- В различных источниках утверждается, что алгоритм успешно применялся при сегментировании изображений, выделении групп в геноме, группировании фильмов на Netflix, моделирование экономических процессов.

2.2 Шаги алгоритма

1. Заполним Similarity Matrix для входных данных алгоритма.
2. Заведём матрицы R (Responsibility) и A (Availability). $A = R = 0$
3. Обновляем матрицы R и A поочередно по формулам $MAX_ITERATIONS$ раз. Обновление может быть прервано, если изменения стали меньше фиксированного значения.
4. Заводим Criterion Matrix: $C = A + R$.
5. В каждой строке матрицы C индекс максимального значения соответствует представителю данной точки. (Соответственно точки с одинаковым индексом максимального значения отнесем к одному кластеру)

2.3 Similarity Matrix

$S(i, j)$ – такая функция, что :

- при $i \neq j$: $S(i, j) = -1 * (\text{расстояние между точками } i \text{ и } j)^2$
- при $i = j$: Элементы на диагонали заполняются одним значением. От выбранного числа зависит количество кластеров, которые получатся. Если выбрано маленькое число, то получится маленькое число кластеров и наоборот. Чаще всего число выбирается минимальным или медианой из уже посчитанных $S(i, j)$, $i \neq j$.

Примечание: за $S(i, j)$ может быть взята любая функция такая, что $S(i, j) > S(i, k)$, если i ближе к j чем к k . Чаще используется именно минус квадрат расстояния между точками.

2.4 Responsibility Matrix

Responsibility Matrix рассчитывается по формуле :

$$r(i, k) \leftarrow s(i, k) - \max_{k' \text{ such that } k' \neq k} \{a(i, k') + s(i, k')\}$$

2.5 Availability Matrix

Availability Matrix рассчитывается по формулам :

1. Элементы на диагонали -
$$a(k, k) \leftarrow \sum_{i' \text{ such that } i' \neq k} \max\{0, r(i', k)\}$$
2. Элементы вне диагонали -
$$a(i, k) \leftarrow \min\left\{0, r(k, k) + \sum_{i' \text{ such that } i' \notin \{i, k\}} \max\{0, r(i', k)\}\right\}$$

2.6 Оптимизации

Можно ввести ограничение на число вершин в дереве:

- В матрицу S добавляют шум для предотвращения вычислительных осцилляций в случаях, когда есть несколько хороших разбиений на кластеры

- При обновлении R и A используется не простое присваивание, а присваивание с экспоненциальным сглаживанием. Авторы советуют использовать параметр сглаживания со значением от 0.5 до 1 (по умолчанию 0.5)
- Часто AP нуждается в постобработке — дополнительной кластеризации лидеров групп. Подходит любой другой метод, может помочь дополнительная информация о задаче.

3 Алгоритм PrePost

3.1 Постановка задачи

- Вход - набор транзакций $T = T_1, T_2 \dots$ (название из-за применимости в поиске ассоциативных правил в чеке клиентов магазинов) и число ϵ ps.
- Выход - Набор itemset'ов, которые часто встречаются в транзакциях.

3.2 Основные понятия

- I - множество $i_1, i_2 \dots$, состоящие из всех встречаемых элементов транзакций.
- ϵ ps - число, поданное на вход алгоритму, по которому мы будем определять встречается ли правило часто или нет.
- itemset : некоторое подмножество I .
- $supp(X) = \frac{|\{t \in T; X \subseteq t\}|}{|T|}$, иными словами - отношение числа транзакций, содержащих множество X , ко всем транзакциям.

X будем считать часто встречаемым, если $supp(x) \geq \epsilon$ ps.

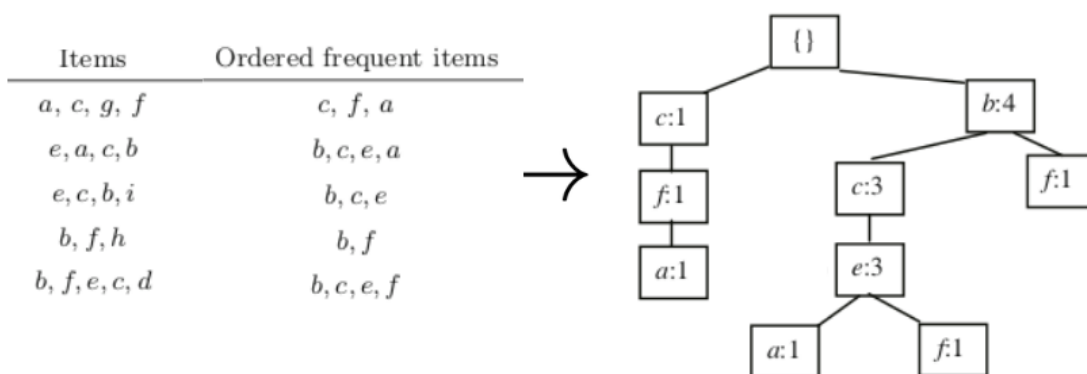
3.3 Шаги алгоритма

1. Построение PPC-tree
 - (a) Построение FP-tree
 - (b) Подсчет pre-order и post-order для каждой вершины дерева.
2. Формирование из дерева NL1
3. Формирование из NL1 NL2
4. Формирование NL-K
5. На каждом шаге, сохраняем часто встречаемые itemsets, для которых $supp(x) \geq \epsilon$ ps.

3.4 FP-tree

Сначала подсчитаем сколько раз каждый из элементов встречается в транзакциях. Отбросим все элементы, которые встречаются меньше $\epsilon * |T|$ раз. Затем внутри каждой транзакции отсортируем оставшиеся элементы по убыванию частоты. После этого строится префиксное дерево по полученным “словам” из транзакций. Также подсчитаем для каждой вершины количество слов, проходящих через неё.

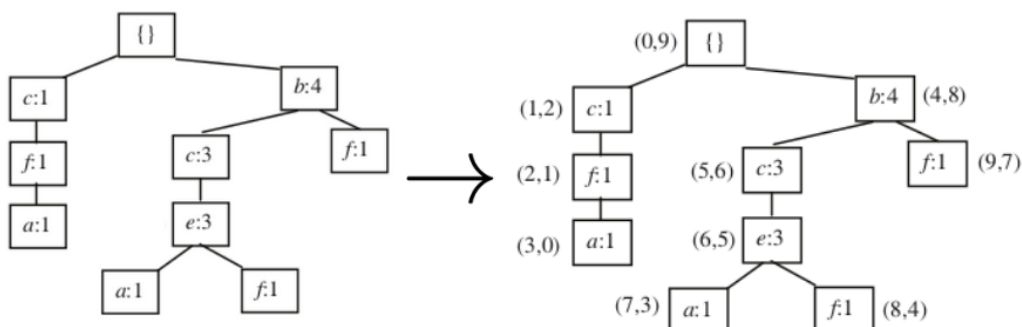
Пример построения FP-Tree по транзакциям:



3.5 PPC-Tree

PPC-Tree - является модификацией FP-Tree. После построения FP-Tree подсчитываются pre-order и post-order для каждой вершины дерева. Pre-order - порядок обхода вершин: вершина, левое поддерево, правое поддерево. Post-order - порядок обхода вершин: левое поддерево, правое поддерево, вершина. Pre-order и post-order легко подсчитываются с помощью любого обхода графа (BFS, DFS).

Пример построения PPC-Tree из FP-Tree:



3.6 NL-1

NL1 для элемента $i \in I$ назовем множество всех вершин в дереве, которые хранят данный символ, представленные в виде троек:

$\langle (\text{pre-order вершины}, \text{post-order вершины}) : \text{кол-во заканчивающихся слов в данной вершине} \rangle$. $\langle x, y \rangle : z$

3.7 NL-k

Пусть мы хотим построить NL-k для подмножества I мощности k - i_1, i_2, \dots, i_k , в котором выполнено $\text{supp}(i_1) \geq \text{supp}(i_2) \geq \dots \geq \text{supp}(i_k)$. Тогда NL-k определяется рекурсивно через NL-(k-1) следующим образом:

- Пусть NL(K-1) для i_2, i_3, \dots, i_k ($P1 = \langle x_{11}, y_{11} \rangle : z_{11}, \dots, \langle x_{1m}, y_{1m} \rangle : z_{1m}$) и NL(K-1) для i_1, i_3, \dots, i_k ($P2 = \langle x_{21}, y_{21} \rangle : z_{21}, \dots, \langle x_{2n}, y_{2n} \rangle : z_{2n}$)
- $\forall i \in [1..m] \forall j \in [1..n]$, если $\langle x_{1i}, y_{1i} \rangle : z_{1i}$ является предком $\langle x_{2j}, y_{2j} \rangle : z_{2j}$, то положим в NL-k $\langle x_{1i}, y_{1i} \rangle : z_{2j}$

Если в NL-k для подмножества $IS = i_1, i_2, \dots, i_k$ просуммировать z_i , то мы получим количество транзакций в которых встречается itemset IS . Таким образом мы будем узнавать часто встречаемые itemset, сравнивая $\text{supp}(IS)$ с $\epsilon * |T|$. А это то, чего мы и хотели добиться.

3.8 Способ нахождения NL-k

Теперь обсудим как же находить NL-k. Мы будем находить их все поуровнево, то есть сначала найдем все NL-1, затем все NL-2 и т. д. Пусть мы нашли все NL для itemsets размера меньших k , и теперь хотим найти все NL-k. Это задача разбивается на два случая.

1 случай. $k = 2$. Заведем квадратную таблицу размером T , равным количеству частых NL-1. Для каждой вершины N из PPC-tree и для каждого предка этой вершины Na , $T[N][Na] +=$ (кол-во слов, проходящих через N ; это число записано в вершине N PPC-tree). Затем для каждой пары элементов i_1, i_2 добавим NL-2 от itemset i_1, i_2 , только если $T[i_1][i_2] \geq \epsilon * |T|$.

2 случай. $k > 2$.

Пусть L_{k-1} - список частых itemsets мощности $k-1$, отсортированных по убыванию supp этих itemsets. Тогда для каждой пары из L_{k-1} , для которых совпадают все элементы, кроме первых построим NL-k. Пусть элементы в паре были $x_1 x_2 \dots x_{k-1}$ и $y_1 y_2 \dots y_{k-1}$ и $\text{supp}(x) \geq \text{supp}(y)$, а их NL-(k-1) были $P1$ и $P2$ соответственно. Тогда будем итерироваться 2 указателями по $P1$ и $P2$.

Если элемент из $P1$ является предком элемента из $P2$, то добавляем в $NL-k$ $\langle x1i, y1i : z2j \rangle$, как это и обсуждалось ранее, иначе если $x1i > x2j$ то сдвигаем второй итератор, а если $y1i < y2j$ то сдвигаем первый итератор. После добавления в $NL-k$ всех вершин, подсчитываем общий $\text{supp}(x_1 x_2 \dots x_{k-1})$ и если он $\geq \epsilon * |T|$, то добавляем $x_1 x_2 \dots x_{k-1}$ в L_k .

4 Источники

Affinity Propagation:

- [Clustering by Passing Messages Between Data Points](#) (Science, 16 FEBRUARY 2007, Brendan J. Frey and Delbert Dueck)
- [Статья на Хабре](#)
- [Статья на Towards Data Science](#)

PrePost:

- [A New Algorithm for Fast Mining Frequent Itemsets Using N-Lists](#)(Sciece China. Information Sciences · September 2012, DENG ZhiHong, WANG ZhongHui JIANG JiaJian)