

УДК 004.9: 004.021: 004.032.26

## МЕТОД АДАПТИВНОЙ КЛАСТЕРИЗАЦИИ ТЕКСТОВЫХ ДОКУМЕНТОВ

**А.А.Мальков**, к.т.н., доцент, Тверской государственной университет, kja227@list.ru

**В.К.Иванов**, к.т.н., доцент, Тверской государственной университет, mtivk@mail.ru

**Аннотация:** В статье предложен метод адаптивной кластеризации текстовых документов – результатов работы поисковой системы. Реализация метода предполагает, что для настройки параметров кластеризации должна использоваться не только информация, полученную от пользователя, но и полученную в результате поиска документов. Идея заключается в использовании нечеткого алгоритма кластеризации Гюстафсона-Кесселя. Для решения задачи определения количества кластеров при инициализации алгоритма предлагается использовать самоорганизующиеся карты Кохонена с динамически изменяемыми размерами. Приведено описание используемых алгоритмов и положительные результаты апробации метода на модельной задаче об ирисах Фишера. Показано, что на основе предложенного решения может быть построен список рубрик, объединяющих семантически связанные источники информации.

**Ключевые слова:** документ карта Кохонена, кластеризация, нейронная сеть, нечеткий алгоритм, поиск.

### Введение

В настоящее время наблюдается экспоненциальный рост потребности в информации, которая, как правило, слабо структурированный характер (достаточно вспомнить текстовые информационные ресурсы Интернет). В этой связи можно выделить несколько важных проблем [1].

Во-первых, это большое количество источников поиска информационных объектов (документов). При этом происходит постоянное пополнение информационных источников, хранилища которых зачастую пересекаются.

Во-вторых, проблемой является сужение области поиска не только источников, но и самих объектов. Результаты запроса в общем случае – это перечисление источников, отсортированных, как правило, по морфологической и синтаксической близости документа к запросу. Возможность отбора документов, семантически близких указанному из всех найденных, существует, однако пользователю приходится просматривать в этом случае сотни объектов. То есть, речь должна идти о действительно семантическом поиске документов.

В-третьих, поиск нужной информации становится все более сложным, трудоемким и неэффективным технологическим процессом. Пользователь, переходя от списка к списку документов, обязан уточнять критерии поиска и доводить свой запрос до некоторого оптимального набора слов, по которому он часто получает во многом знакомый ему перечень документов. Процесс поиска закичивается, а время поиска значительно возрастает.

В статье предложен метод адаптивной кластеризации текстовых документов – результатов работы поисковой системы. Реализация метода предполагает, что для настройки параметров кластеризации должна использоваться не только информация, полученную от пользователя, но и полученную в результате поиска документов. Идея заключается в использовании модифицированного нечеткого алгоритма кластеризации Гюстафсона-Кесселя. Для решения задачи определения количества кластеров при инициализации алгоритма предлагается использовать самоорганизующиеся карты Кохонена с динамически изменяемыми размерами. Приведено описание используемых алгоритмов и положительные результаты апробации метода на модельной задаче об ирисах Фишера. Показано, что на основе предложенного решения может быть построен список рубрик, объединяющих семантически связанные источники информации.

Работы выполнены при финансовой поддержке РФФИ, проект № 13-07-00342.

### **Теоретическое обоснование**

С точки зрения пользователя система должна обладать хорошим соотношением, таких, в некотором смысле противоположных, показателей, как полнота и точность поиска, а также выполнять новый поиск на основе полученных ранее результатов. Найденные и рубрицированные документы в большинстве случаев могут пересекаться по рубрикам, но у пользователя должен быть выбор между четкой и нечеткой рубрикацией. Кроме того, от пользователя, как от инициатора поиска, могла бы быть получена некоторая информация, необходимая для настройки параметров системы, но ее объем не должен превосходить разумного предела.

С точки зрения реализации система должна использовать не только информацию, полученную от пользователя, но и полученную в результате поиска документов. В частности, по запросу пользователя определенным образом строится список найденных документов и, используя обработанный определенным образом запрос, вычисляется матрица весов документов –  $tf \cdot idf$  [2]. На основе этой информации необходимо получить список рубрик и распределение документов по ним.

Для решения этой задачи предлагается выполнять нечеткую кластеризацию векторов документов матрицы весов. Проблема здесь – это динамически изменяющееся количество компонент этих векторов в зависимости от видоизменения запроса. Также нужно учитывать и возможность классификации из-за возможного появления новых источников информации.

Методов, выполняющих кластеризацию документов достаточно много. Однако не существует универсального метода. Каждый из них хорош (или даже оптимален) при выполнении ряда условий, но имеются и недостатки.

Метод *SuffixTreeClustering* [3] предполагает повторную обработку текстов документов. *LSA/LSI* [4] выполняет огромное количество вычислений, и в результате своей работы выдает непересекающиеся кластеры. Для метода *ConceptIndexing* [5] необходимо наличие обучаемого варианта, задание количества кластеров. В *SingleLink*, *CompleteLink*, *GroupAverage* [6] также полученные кластеры не пересекаются. Главным параметром *K-means* и *Fuzzy-C-Means* [6]

является задание количества кластеров. Для карт Кохонена SOM[6] основной недостаток, как правило, - это длительный процесс обучения.

### Постановка задачи

Пусть документ  $d$  представляется набором терминов  $\{t_l\}_{l=1}^M$ , которые влияют на отнесение документа к какой-либо рубрике, т.е. документ – это вектор  $d = \{t_1, t_2, \dots, t_M\}$ . Тогда мощность  $M$  пространства терминов  $T$  будет представлять размерность пространства документов. Координатами вектора документа, который представляется на обработку, будут величины значимости конкретного термина для этого документа. Если в некотором документе отсутствует термин  $t_l$ , то  $l$ -я координата вектора принимается равной 0. В данной работе за основу взята мера  $tf*idf$ , при этом координаты документов в пространстве терминов преобразуются в матрицу весов документов.

Кластером  $C$  в пространстве терминов  $T$  назовем множество схожих, в смысле некоторой меры  $\rho$ , векторов весов терминов документов. Тогда семантическую схожесть документов будет определять выбранная мера схожести  $\rho$  представляющих их векторов. Таким образом, необходимо провести кластеризацию множества документов  $D = \{d_i\}_{i=1}^N$ . В итоге необходимо получить разбиение  $C = C_1 \cup C_2 \cup \dots \cup C_K$  множества  $D$  на возможно пересекающиеся группы, где каждому документу будет приписана степень принадлежности каждому кластеру, которая будет определять вес значимости документа для каждого кластера-рубрики.

Важным вопросом является определение «качества» кластеризации – разбиение множества документов на «правильное» количество кластеров. На сегодняшний день не существует единого подхода для определения наилучшего количества кластеров. Известны лишь верхняя и нижняя оценки на количество кластеров  $K$ . В качестве ограничений на количество кластеров обычно берутся границы интервала  $K \in [2, \sqrt{N/2}]$ , где  $N$  – количество документов.

### Описание предлагаемого метода кластеризации

Для решения указанной задачи в различных источниках [7], [8] предлагается вычислять так называемые функционалы качества. Ниже представлены некоторые из них и условия для нахождения оптимального количества кластеров.

$$V_{PC} = \frac{1}{N} \sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^2 \rightarrow \max$$

$$V_{CI} = (K \cdot V_{PC} - 1) / (K - 1) \rightarrow \max$$

$$V_{CE} = \frac{-1}{N} \sum_{i=1}^K \sum_{j=1}^N \mu_{ij} \log \mu_{ij} \rightarrow \min$$

$$V_{XB} = \frac{\sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^m \|d_j - C_i\|^2}{N \cdot \min_{p \neq i} \{\|C_p - C_i\|^2\}} \rightarrow \min$$

$$V_{Know} = \frac{\sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^m \|d_j - C_i\|^2 + \frac{1}{K} \sum_{i=1}^K \|C_i - \bar{C}\|^2}{N \cdot \min_{p \neq i} \{\|C_p - C_i\|^2\}} \rightarrow \min$$

$$V_T = \frac{\sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^m \|d_j - C_i\|^2 + \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{p=1, p \neq i}^K \|C_i - C_p\|^2}{N \cdot \min_{p \neq i} \{\|C_p - C_i\|^2\} + \frac{1}{K}} \rightarrow \min$$

$$V_{FS} = \sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^m \|d_j - C_i\|^2 - \sum_{i=1}^K \sum_{j=1}^N \mu_{ij}^m \|C_i - \bar{C}\|^2 \rightarrow \min$$

где  $\bar{C} = \sum_{i=1}^K \frac{C_i}{K}$ .

Исходя из такой постановки задачи, для поиска количества кластеров и соответствующего разбиения можно следовать следующему алгоритму.

Задать  $K_{min}, K_{max}$

$p = K_{min}$

*repeat*

Провести кластеризацию  $C = C_1 \cup C_2 \cup \dots \cup C_p$

Вычислить индекс(-ы) для полученного варианта кластеризации

*if* индекс(-ы) достиг(-ли) оптимального значения *then*

*return*  $p$

*else*

$p = p + 1$

*endif*

*until*  $p \leq K_{max}$

Здесь на каждом шаге алгоритма необходимо запускать некоторый метод кластеризации и вычислять значения функционалов. При этом нужно следить за «скачками» этих функционалов или разработать некоторый алгоритм для вычисления оптимума того или иного функционала, что в указанной постановке весьма проблематично, так как сами функционалы содержат информацию о разбиении. Также нужно обращать внимание на варианты кластеризации, при которых произошли эти «скачки». На следующем этапе необходимо просматривать выделенные варианты кластеризации и выбрать наилучший. Очевидно, что вычислительная сложность такого подхода резко возрастает с увеличением объема входных данных. Также нужно отметить, что при

изменении объема входных данных нужно заново запускать этот алгоритм, что также увеличивает временную сложность работы алгоритма.

Апробация вышеописанного алгоритма была проведена на тестовых данных, известных как ирисы Фишера. В качестве алгоритма кластеризации применялся классический алгоритм Гюстафсона-Кесселя, учитывающий формы кластеров. Количество кластеров изменялось от 2 до 7. Результат работы алгоритма приведен в таблице 1.

Таблица 1. Значения индексов

Число кластеров	$V_{PC}$	$V_{CI}$	$V_{CE}$	$V_{XB}$	$V_{Know}$	$V_T$	$V_{ES}$
2	0,9345	0,869	0,1184	0,0538	8,3266	9,0766	-8,7751
3	0,8758	0,8137	0,2256	0,1345	21,4482	23,9867	-10,1209
4	0,8392	0,7856	0,3208	0,153	24,7771	27,8223	-10,5627
5	0,7981	0,7476	0,3884	0,35	60,0911	71,4726	-10,3254
6	0,7148	0,6577	0,576	0,428	71,2044	81,0086	-9,6332
7	0,7135	0,6657	0,6227	0,2734	44,5073	49,1686	-9,863

Видно, что с точки зрения оптимальности функционалов наилучшим можно считать вариант разбиения исходного множества на 2 кластера. В качестве подтверждения правильности оценки количества кластеров, указанной выше, можно отметить тот факт, что при выборе большего количества кластеров значения функционалов принципиально не отличаются от значений, полученных для последнего варианта кластеризации.

Одним из вариантов решения задачи определения количества кластеров может быть модифицированный алгоритм SOM [9]:

1. Инициализация. Для исходных векторов семантических весов  $w_j(0)$  выбираются случайные значения. Единственным требованием здесь является различие векторов для разных значений  $j = 1..K$ , где  $K$  – общее количество нейронов в решетке. При этом рекомендуется сохранять малую амплитуду значений. Например, можно рандомизировать аргумент функции синуса. Веса нейронов рекомендуется нормализовать.

2. Подвыборка. Из входного пространства выбирается вектор  $d$  с определенной вероятностью. Размерность вектора равна  $M$ .

3. Поиск максимального подобия. Осуществляется поиск нейрона-победителя  $i_d(x, y)$  на шаге  $t$ , используя критерий минимума Евклидова расстояния:

$$i_d(x, y) = \underset{j}{\operatorname{argmin}} \|d - w_j\|, j = 1..K$$

где  $(x, y)$  – координаты нейрона в решетке. На этом шаге необходимо учитывать проблему «мертвых» нейронов. Для ее решения использовались следующие положения. Учитывая тот факт, что после победы нейрон «отдыхает» необходимо «ограничить» его активность на следующем этапе обучения [20]. Для этого можно вести учет активности нейронов:

$$p_j(t + 1) = \begin{cases} p_j(t) + \frac{1}{K}, j \neq i_d(x, y) \\ p_j(t) - p_{min}, j = i_d(x, y) \end{cases}$$

где  $p_{min}$  – минимальный «потенциал», разрешающий участие нейрону в конкурентной борьбе. На практике хороший результат дает  $p_{min} = 0.75$ . Кроме того, количество побед нейронов учитывается [9] при поиске нейрона-победителя, что позволяет задействовать часть нейронов из области пространства, где отсутствуют данные или их количество ничтожно мало.

$$i_d(x, y) = \underset{j}{\operatorname{argmin}}(NW_j \cdot \|d - w_j\|), j = 1..K$$

где  $NW_j$  – количество «побед» нейрона  $j$ .

4. Коррекция. Корректируются векторы семантических весов всех активных нейронов, используя формулу:

$$w_j(t + 1) = w_j(t) + \eta(t)h_{j,i_d(x,y)}(t)(d - w_j(t))$$

где  $\eta(t)$  – параметр скорости сходимости;  $h_{j,i_d(x,y)}(t)$ – функция окрестности нейрона-победителя  $i_d(x, y)$ . Оба этих параметра динамически изменяются:

$$\eta(t) = \eta_0 e^{\left(\frac{-t}{\tau_2}\right)}$$

$$h_{j,i_d(x,y)}(t) = e^{\left(\frac{-d_{j,i}^2}{2\sigma^2(t)}\right)}$$

$$\sigma(t) = \sigma_0 e^{\left(\frac{-t}{\tau_1}\right)}$$

Здесь:  $\eta_0$  – начальное значение скорости сходимости, рекомендуемое значение 0.1, при этом  $\eta(t)$  не должно быть менее 0.01;  $\sigma(t)$  - ширина топологической окрестности нейрона, на начальном этапе  $\sigma_0$  полагают равной радиусу решетки, что означает активность всех нейронов сети на начальном этапе обучения;  $\tau_2 = 1000, \tau_1 = \frac{\tau_2}{\log \sigma_0}$  – временные параметры.

5. Продолжение. Возврат к шагу 2. Вычисления продолжают до тех пор, пока в карте признаков не перестанут происходить заметные изменения.

Преимущество этого алгоритма заключается в том, что он самостоятельно может определить первоначальное распределение центров кластеров [9]. Процесс адаптации нейрона-победителя позволяет так настроить его веса, что он становится «ковариационным центром» найденного кластера.

На данном этапе открытым остается вопрос о размере решетки нейронов. В данном случае задача поиска количества кластеров «преобразуется» в задачу определения размера нейронной сети. Для ее решения была позаимствована идея иерархических методов кластеризации, где разбиение начинается с единственного кластера, который в процессе кластеризации дробится на некоторое конечное число кластеров в соответствии с некоторым условием.

Примененный алгоритм определения размеров карты Кохонена предусматривает, что для каждого кластера вычисляется величина:

$$v_{ij} = \frac{1}{N_c} \sum_{d_i \in C_{ij}} \|d_i - w_{ij}\|$$

где  $C_{ij}$  – множество векторов, отнесенных к  $j$ -му кластеру,  $N_c$  – количество документов в  $C_{ij}$ .

Весь процесс обучения, в данном случае, необходим для определения количества кластеров, то есть возможно изменение карты. Для определения ее новых размеров ищется нейрон  $w^*$ , для которого  $v_{ij}$  максимален.

$$v_{max} = \max \left( \frac{1}{N_c} \sum_{d_i \in C_{ij}} \|d_i - w_{ij}\| \right)$$

Чтобы не нарушать «обученность» нейронов вся карта не перестраивается, а добавляется только строка и столбец нейронов между нейроном  $w^*$  и нейроном, для которого расстояние, в смысле выбранной метрики, наибольшее. Далее инициализируются веса добавленных нейронов, для чего использовался самый простой способ – вычисление средневзвешенного весов нейронов из топологической окрестности.

В качестве критерия остановки изменения размеров карты предложен следующий критерий:

$$N = \frac{1}{K^2} \sum_{i,j=1..K} v_{ij}$$

где  $N < 0.55 v_{p_{max}}$ ,  $v_{p_{max}}$  – величина, определенная на предыдущем шаге конкретизации карты.

Для дальнейших исследований необходимо в условие остановки увеличения размера сети включить семантическую «обособленность» кластеров.

На последнем этапе формируются кластеры при помощи алгоритма Гюстафсона-Кесселя.

### **Заключение**

В статье описан метод адаптивной кластеризации текстовых документов и показана его применимость к решению задачи построения списка рубрик, объединяющих семантически связанные источники информации.

Мы предполагаем, что предложенный метод будет использован в разработке модуля классификации и идентификации связей текстовых документов и мультимедийных объектов интеллектуальной системы информационной поддержки инноваций в науке и образовании [2]. В результате для кластеров семантически связанных данных должны быть выполнены оценка качества классификации, отбор объектов с высокой степенью релевантности и интерпретация результатов. Эти модули завершают генерацию решений, имеющих инновационный потенциал, которая осуществляется в заданном тематическом сегменте или по заданному объекту (промышленному изделию, технологии, продукту).

В настоящее время готовится программная реализация описанного в статье метода, предназначенная для включения в программный продукт GeneticAlgorithmFramework (GAF), представленном в [10].

### **Литература**

1. Маннинг, Кристофер Д., Рагхаван, Прабхакар, Щютце, Хайнрих. Введение в информационный поиск. : Пер. с англ. – М. : ООО «И.Д. Вильямс», 2011. – 528 с.

2. Палюх Б.В., Иванов В.К., Сотников А.Н.. Архитектура интеллектуальной системы информационной поддержки инноваций в науке и образовании // Программные продукты и системы. – 2013. – № 4. – С. 197-202.
3. Oren Eli Zamir. A Phrase-Based Method for Grouping Search Engine Results. University of Washington, Department of Science & Engineering. 1999.
4. Susan T. Dumains, George W. Furnas, Thomas K.Landauer. Indexing by Latent Semantic Analysis. Bell Communications Research 435 South St. Morristown, NJ 07960. Richard Rashman: University Of Western Ontario. 2008.
5. Eui-Hong (Sam) Han and George Kapiris. Centroid-Based Document Classification: Analysis & Experimental Results. University of Minnesota, Department of Computer Science / Army HPC Research Center,4-192 EECS Bldg., 200 Union St. SE, Minneapolis, MN 55455 USA. 2000.
6. Барсегян, А. А. Анализ данных и процессов: учеб. пособие / А. А. Барсегян, М. С. Куприянов, И. И. Холод, М. Д. Тесс, С. И. Елизаров. - 3-е изд., - СПб.: БХВ-Петербург, 2009. -512 с.: ил.
7. C. Duo et al., An Adaptive Cluster Validity Index for the Fuzzy CMeans. International Journal of Computer Science and Network Security, 7(2)146-156, 2007.
8. Yuangang Tang, Fuchun Sun, Zengqi Sun, Improved Validation Index for Fuzzy Clustering, in American Control Conference, 2005. 1120-1125.
9. Виноградов Г.П., Мальков А.А. Эволюционные методы кластеризации, использующие нечеткие отношения и субъективные оценки. Сборник трудов Международной научно-технической конференции AIS'08, CAD-2008, «Интеллектуальные системы», «Интеллектуальные САПР», т.1, М.: Физматлит., с.7-15.
10. Рутковская, Д., Пилиньский, М., Рутковский, Л. Нейронные сети, генетические алгоритмы и нечеткие системы: Пер. с польск. И. Д. Рудинского. – 2-е изд., стереотип. – М.: Горячая линия–Телеком, 2013. – 384 с.
11. Иванов, В.К., Мескин П.И. Реализация генетического алгоритма для эффективного документального тематического поиска // Программные продукты и системы. - Тверь, 2014. - № 4. - С. 118-126. - DOI:10.15827/0236-235X.108.118-126.