

УДК 004.75

Д.Г. БУХАНОВ, В.М. ПОЛЯКОВ, В.Г. РУБАНОВ, В.Г. СИНЮК

Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация

МЕТОД СТРУКТУРИЗАЦИИ В ЗАДАЧАХ САМОДИАГНОСТИРОВАНИЯ РАСПРЕДЕЛЕННЫХ ВЫЧИСЛИТЕЛЬНЫХ СИСТЕМ

Предложено использовать метод кластеризации для структурирования и выбора расположения точек наблюдения при децентрализованном диагностировании распределенных вычислительных систем. Рассмотрены особенности подхода, предложены и исследованы критерии оценки качества кластеризации. Проведены модельные эксперименты, исследовано влияние выбора начальных разбиений на скорость сходимости алгоритма. Исследовано влияние выбора начальных разбиений на скорость сходимости алгоритма. Предложен критерий оценки результатов алгоритма кластеризации.

Ключевые слова: диагностирование, распределенные вычислительные системы, кластеризация.

Введение

Исследования в области диагностирования распределенных вычислительных систем (РВС) показывают, что методы, основанные на централизованном подходе, при котором весь диагностический трафик передается в один узел, ведут к значительному снижению производительности, а также накладывают ряд технологических ограничений, которых можно избежать.

Одним из путей решения задач контроля и диагностирования РВС является применение методов системного децентрализованного самодиагностирования, обоснование которых приведено в [1 – 3].

Основой данного подхода является предположение, что каждая ЭВМ может определить корректный диагностический образ всей системы, базируясь на результатах взаимных проверок других ЭВМ.

Технологии децентрализованного самодиагностирования позволяют снизить избыточность передаваемого трафика в РВС, проводить диагностирование объектов в реальном масштабе времени на основании информации, получаемой в ходе «штатного» режима функционирования.

Системное самодиагностирование может быть представлено в виде повторения циклов, каждый из которых состоит из фаз:

- 1) выполнение ЭВМ всех предписанных им проверок и формирования результатов;
- 2) передача этих результатов в проверяющие ЭВМ;
- 3) анализ и оценка полученных результатов в проверяющих ЭВМ и формирование в них собственного списка элементов, подозреваемых в неисправности;

4) взаимнообмен собственными списками между всеми ЭВМ распределенной системы при помощи некоторого алгоритма;

5) формирование во всех исправных или специально выделенных диагностических станциях подписков глобального списка неисправностей для принятия решений по возможному восстановлению работоспособности.

1. Постановка задачи

Топология РВС представлена в виде графовой модели, описываемой картежем

$$G_{РВС} = (Y \cup D \cup C, U) \quad (\text{рис. 1}).$$

В такой модели $Y = \{y_j\}$ - множество входных узлов; $D = \{d_i\}$ - множество внутренних элементов (узлов) системы; $C = \{c_k\}$ - множество узлов – центров кластеров (или контрольных точек); $U = \{u_q\}$ - множество ребер. В качестве условного расстояния между узлами принято время прямой и обратной передачи информационных пакетов.

Требуется расположить агенты наблюдения, сбора и анализа информации в сети для снижения передаваемого контрольно-диагностирующего трафика, а также для наиболее быстрого определения места возможных неисправностей.

2. Использование алгоритма кластеризации для задания точек контроля в РВС

Рассмотрим использование классического алгоритма k-means[4] для разбиения структуры РВС

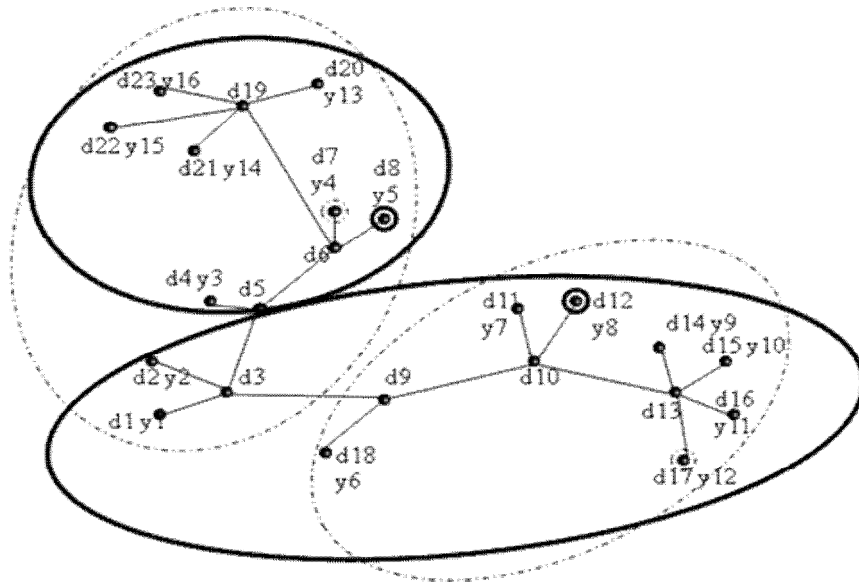


Рис. 1. Пример сетевой структуры

на кластеры с целью организации размещения точек контроля и диагностирования в реальном масштабе времени.

Данный алгоритм реализует следующую последовательность действий:

- 1) выбирается k – произвольных центров системы;
- 2) все элементы разбиваются на k – групп, наиболее близких к одному из центров;
- 3) затем вычисляются новые центры кластеров.

Алгоритм повторяется до тех пор, пока не перестанут меняться центры кластеров и соответственно границы.

Пусть заданы три множества:

$$M = \{m_i\}_1^k \text{ – обучающее множество,}$$

$$C = \{c_j\}_1^e \text{ – центры кластеров,}$$

$$V = \{v_{it}\}_1^e \text{ – матрица разбиения.}$$

Множество M представляет собой набор векторов данных m_i , характеризующих “условное расстояние” между i -м и остальными узлами РВС, полученный измерением времени передачи пакетов; k – количество векторов данных; c_j – центр кластера; e – количество кластеров;

$$2 \leq e \leq \sqrt{|M|}; \tag{1}$$

v_{it} – вектор разбиения, $v_{it} \in \{0;1\}$;

$$\sum_1^e v_{it} = 1;$$

$$0 < \sum_1^k v_{it} < k.$$

При инициализации алгоритма вектор разбиения заполняется значениями в соответствии с предполагаемыми центрами кластеров.

Алгоритм представляет собой следующую итерационную процедуру:

шаг 1: номер итерации $l=0$; инициализировать начальное разбиение, выбрать точность δ ;

шаг 2: определить центры кластеров по формуле:

$$c_j^{(i)} = \frac{\sum_{i=1}^k v_{ij}^{(i-1)} m_j}{\sum_{j=1}^k v_{ij}^{(i-1)}}, \tag{2}$$

где $1 \leq i \leq e$;

шаг 3: обновить матрицу разбиения, чтобы минимизировать квадраты ошибок, используя формулу:

$$v_{ij}^{(l)} = \begin{cases} 1, & \min d(m_j; c_j^l), \\ 0, & \text{в остальных случаях;} \end{cases} \tag{3}$$

шаг 4: проверить условие $|V^l - V^{l-1}| < \delta$; если условие выполняется – завершить процесс, иначе перейти в шаг 2 с номером итерации $l = l + 1$.

Для структуры, изображенной на рис. 1, обучающее множество M будет представлено в виде матрицы 16×16 .

При $e=2$ результат работы алгоритма будет иметь следующий вид:

шаг 1: $l=0$; начальное разбиение выполнено таким образом, что центрами кластера являются узлы y_{12} и y_4 ;

V^1 :

0	0	0	0	0	1	1	1	1	1	1	1	0	0	0	0
1	1	1	1	1	0	0	0	0	0	0	0	1	1	1	1

шаг 2: определим центры кластеров из формулы (2);

 c_1 :

3,14	3,14	3,71	4,85	4,85	0,85	0,85	0,85
0,85	0,85	0,85	0,85	6,42	6,42	6,42	6,42

 c_2 :

3	3	2,28	1,71	2	4,28	5,71	5,71
7,42	7,42	7,42	7,42	2,57	2,57	2,57	2,57

шаг 3: обновить матрицу разбиения согласно формуле (3);

 V^{l+1} :

1	1	0	0	0	1	1	1	1	1	1	1	0	0	0	0
0	0	1	1	1	0	0	0	0	0	0	0	1	1	1	1

шаг 4: матрица разбиения изменилась, следовательно, условие останова не выполнилось ($|V^l - V^{l-1}| \neq 0$). Алгоритм не завершился, переходим к шагу 2 и $l=l+1$. На следующей итерации при центрах кластеров в u_8, u_5 алгоритм завершается, матрица разбиения не меняется. Алгоритм завершился в две итерации.

На выходе алгоритма получится матрица разбиения с количеством кластеров равным e . Для определения рационального количества кластеров используется метод адаптивной кластеризации [4].

Ключевым моментом в адаптивной кластеризации является выбор критерия, по которому будет оцениваться качество кластеризации. Зададим следующие величины: минимальное среднее внутреннее расстояние в кластере

$$R^{BH} = \sum_{k=1}^{|c|} \sum_{j=1}^j d(m_j; c_k) / |c| \rightarrow \min;$$

и максимальное среднее расстояние между кластерами

$$R^{BHE} = \sum_{k=1}^{|c|} d(c_k; c) / |c| \rightarrow \max.$$

Тогда критерием выбора количества точек наблюдений выберем отношение:

$$Kr = R^{BH} / R^{BHE} \rightarrow \min;$$

Из формулы (1) следует, что выбор количества кластеров в примере может быть только 2, 3 и 4. При $e=2$ отношение R^{BH} к R^{BHE} дает значение 0,36. При $e=3$, $Kr=0,21$ при $e=4$, $Kr=0,26$. Из этого следует, что число кластеров равное трем является наиболее предпочтительным.

Таким образом, минимизация отношения R^{BH} к R^{BHE} дает оптимальное разбиение на кластеры с точки зрения расстояния между элементами внутри кластера и расстояния между кластерами.

Основными недостатками такого подхода являются:

- 1) сильная зависимость от выбора начальных значений центров кластера;
- 2) невозможность определения оптимального количества кластеров.

3. Выбор начальных значений центров кластера

Классическая реализация подразумевает случайный выбор центров кластера. Данный подход является неэффективным в связи с тем, что при неудачном выборе происходит увеличение числа итераций в алгоритме.

Рассмотрим ситуацию, при которой центры кластеров будут находиться в самых отдаленных элементах системы. Исключительной является ситуация, когда элемент находится на равном расстоянии от нескольких кластеров.

Для снижения числа итераций можно выбрать максимальные экстремумы функций, например, как показано на рис. 2.

График представляет собой зависимость времени передачи в условных единицах (t_y) до всех узлов (N).

Точка, лежащая на оси абсцисс, означает номер узла, относительно которого измерялось время. В данном случае график показан для 82 узла.

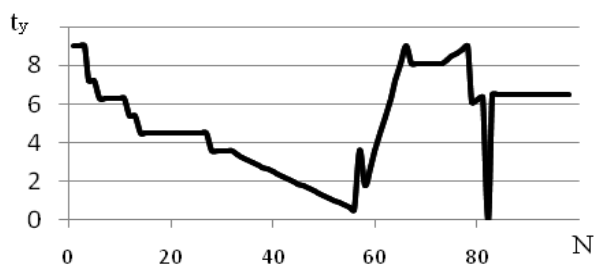


Рис. 2. График удаленности узлов

В ходе модельного експеримента при 98 узлах были получены следующие результаты (табл. 1).

Таблица 1
Количество итераций алгоритма в зависимости от выбора начального разбиения

Количество кластеров	Случайный выбор	Выбор крайних точек
3	12	4
4	7	3
5	11	4

В таблице показано количество итераций в зависимости от выбора начальных центров кластеров. Из результатов видно, что выбор начального разбиения сильно влияет на количество итерации алгоритма.

4. Оценка критерия выбора количества точек наблюдения

Из результатов эксперимента следует (табл. 2), что при использовании адаптивной кластеризации рассогласование значений минимального среднего

внутреннего расстояния $R^{ВН}$ в кластере к максимальному среднему внешнему расстоянию $R^{ВНЕ}$ не происходит.

Таблица 2
Значение критериев при разном количестве кластеров

$ c $	$R^{ВН}$	$R^{ВНЕ}$	Kr
2	3,46	5	0,69
3	1,5	6,3	0,23
4	1,13	6,7	1,7
5	3,67	4,5	0,81
6	3,06	3,86	0,82
7	3,06	3,38	0,9
8	2,97	3	0,99

Для этих наборов данных лучшим по критерию, представленному в статье, является выбор четырех кластеров. При таком количестве кластеров максимальное внешнее расстояние равно 6,7, а среднее внутреннее расстояние 1,13.

Графическая иллюстрация показана на рис. 3.

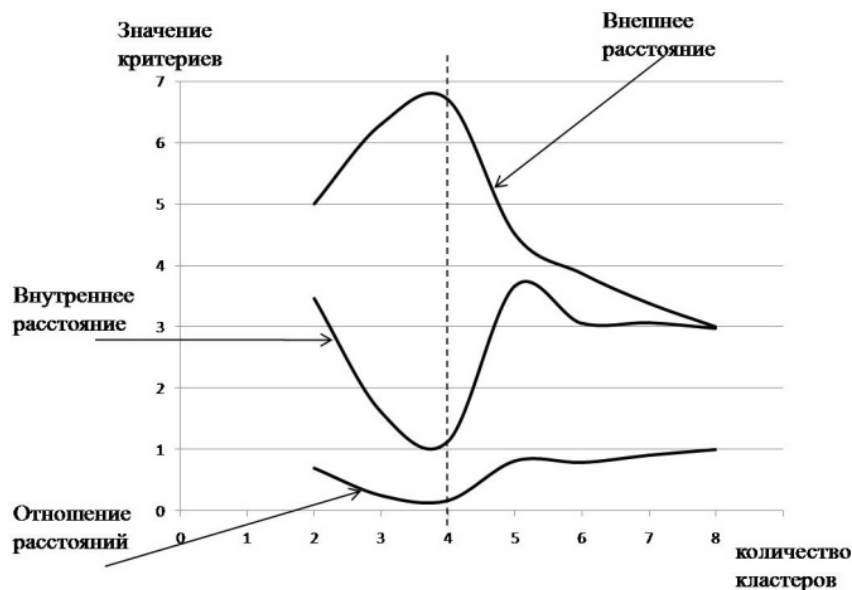


Рис. 3. График функций изменения значений критериев при разном числе кластеров

Из графика видно, что наилучшее разбиение получается в ситуации когда

$$R^{ВНЕ} - R^{ВН} \rightarrow \max.$$

Таким образом, максимизация разности $R^{ВНЕ}$ к $R^{ВН}$ дает оптимальное разбиение на кластеры с точки зрения расстояния между элементами внутри кластера и расстояния между кластерами.

Заключение

Рассмотрен метод выбора расположения элементов контроля в распределенных системах диагностирования на основе кластерного анализа. Использование в качестве алгоритма кластеризации алгоритм k-means влечет за собой две основные проблемы: выбор начального разбиения и оценка качества кластеризации.

Исследовано влияние выбора начальных разбиений на скорость сходимости алгоритма. Предложен критерий оценки результатов алгоритма кластеризации.

Работа выполнена при финансовой поддержке РФФИ, проекты: 12-07-000493-а, №11-01-00359-а.

Литература

1. Che-Liang Yang, *A Fault Identification Algorithm for t_i – Diagnosable System [Text]* / Che-Liang Yang, Gerald M. Masson // *IEEE Trans. Comput.* – 1986. – № 6. – P. 503 – 510.

2. Fabrizio, L. *Diagnosis and fault identification algorithms for large scale computing system [Text]* / Lombardi Fabrizio, Chin-Long Wey // *Supercomputing System Proceedings of the First International Conference, St Petersburg, Florida, Dec. 16-29, 1985.* – P. 404 – 413.

3. Молдованова, О.В. *Децентрализованный алгоритм самодиагностики для крупномасштабных распределенных вычислительных систем различных топологий [Текст]* / О.В. Молдованова // *Пробл. информатики.* – 2012. – № 2. – С. 70 – 75.

4. *Технологии анализа данных [Текст]* / А.А. Барсегян, М.С. Куприянов, В.В. Степаненко, И.И. Холлод. – СПб.: БХВ-Петербург 2007. – С 143 – 172.

Поступила в редакцию 28.02.2013, рассмотрена на редколлегии 13.03.2013

Рецензент: д-р физ.-мат. наук, проф., зав. кафедрой высшей математики А.Г. Брусенцев, Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация.

МЕТОД СТРУКТУРИЗАЦІЇ У ЗАДАЧАХ САМОДІАГНОСТУВАННЯ РОЗПОДІЛЕНИХ ОБЧИСЛЮВАЛЬНИХ СИСТЕМ

Д.Г. Буханов, В.М. Поляков, В.Г. Рубанов, В.Г. Синюк

Запропоновано використовувати метод кластеризації для структурування та вибору розташування місць спостереження при децентралізованому діагностуванні розподілених обчислювальних систем. Розглянуто особливості підходу, запропоновані та досліджені критерії оцінки якості кластеризації. Проведені модельні експерименти, досліджено вплив вибору початкових розподілень на швидкість збіжності алгоритму. Досліджений вплив вибору початкового розбиття на швидкість збіжності алгоритму. Запропонований критерій оцінки результатів алгоритму кластеризації.

Ключові слова: діагностування, розподілені обчислювальні системи, кластеризація.

METHOD OF STRUCTURING IN TASKS OF SELF-DIAGNOSIS OF DISTRIBUTED COMPUTER SYSTEMS

D.G. Bykhanov, V.M. Polakov, V.G. Rubanov, V.G. Sinuk

It's proposed to use a clustering method to structure and select the location of the observation points at the decentralized diagnosis of distributed computing systems. The features of the approach are analyzed. It's proposed and researched criteria of quality evaluation of clustering. Model experiments the effect of the initial partition for the rate of convergence of the algorithm is researched. Influence of choice of the initial breaking up on speed of algorithm. convergence is analyzed. The criterion of estimation of results of algorithm of clusterization is offered.

Key words: diagnostics, distributed computing systems, clustering.

Буханов Дмитрий Геннадьевич – аспирант кафедры программного обеспечения вычислительной техники и автоматизированных систем, Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация, e-mail:db.old.stray@gmail.com.

Поляков Владимир Михайлович – канд. техн. наук, зав. кафедрой программного обеспечения вычислительной техники и автоматизированных систем, Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация.

Рубанов Василий Григорьевич – д-р техн. наук, профессор, зав. кафедрой технической кибернетики, Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация.

Синюк Василий Григорьевич – канд. техн. наук, проф. кафедры программного обеспечения вычислительной техники и автоматизированных систем, Белгородский государственный технологический университет им. В.Г. Шухова, Белгород, Российская Федерация.