

**Методика адаптивной кластеризации  
фактографических данных на основе  
интеграции методов минимального  
остовного дерева и нечетких K-средних**

Нейский И.М.

Научный руководитель:  
к.т.н., доцент, Филиппович А.Ю.

# Задачи интеллектуального анализа данных



Кластеризация – разбиение совокупности объектов на однородные группы (кластеры)

# Проблемы кластерного анализа

1. Выбор метода исследования

[И.Д. Мандель]

3. Выбор значения параметра «Количество кластеров»

[А.А. Баргесян, М.С. Куприянов, В.В. Степаненко, И.И. Холод, Н.Б. Паклин]

2. Оценка качества полученного разбиения

[Н.Б. Паклин, А. К. Jain, R.C. Dubes, Б. Дюран, П. Оделл]

4. Постоянно увеличивающиеся объемы данных

**Адаптивная кластеризация – кластеризация, для которой количество кластеров является результатом исследования**

|                    | <b>MST</b>  | <b>Fuzzy C-means</b>   |
|--------------------|---|--|
| <b>Назначение</b>  | Кластеризация больших наборов числовых данных                             | Кластеризация больших наборов числовых данных  |
| <b>Достоинства</b> | Выделяет кластеры произвольной формы, в т.ч. кластеры выпуклой и вогнутой | Плавное разбиение объектов   |
| <b>Недостатки</b>  | Работа только с числовыми данными, четкое распределение объектов          | Вычислительная сложность, задание количества кластеров, возникает неопределенность с объектами, которые удалены от центров всех кластеров, выделение кластеров сферической формы |

# Примеры практических задач

| Показатель                                       | Количественная характеристика | Качественная характеристика |
|--|-------------------------------|-----------------------------|
| Количество исходных объектов (шт.)               | 500 - 50 000                  | -                           |
| Количество значимых характеристик объектов (шт.) | 70 - 150                      | -                           |
| Типы характеристик                               | два вида                      | числовые, лингвистические   |
| Форма получаемых кластеров                       | -                             | сложная, с пересечениями    |
| Количество получаемых кластеров (шт.)            | 5 - 30                        | результат исследования      |

- Разработка тарифной сетки;
- Группировка перечня предлагаемых услуг;
- Формирование потребительской корзины;
- Сегментирование сферы деятельности;
- Распознавание изображений;

- Тематический анализ библиотеки документов;
- Оптимизация использования складских помещений;
- Выделение потенциальных покупателей;
- Построение показательной выборки и др.

# Цель и задачи

## Цель

Разработка методики адаптивной кластеризации фактографических данных на основе интеграции методов минимального остовного дерева и нечетких K-средних

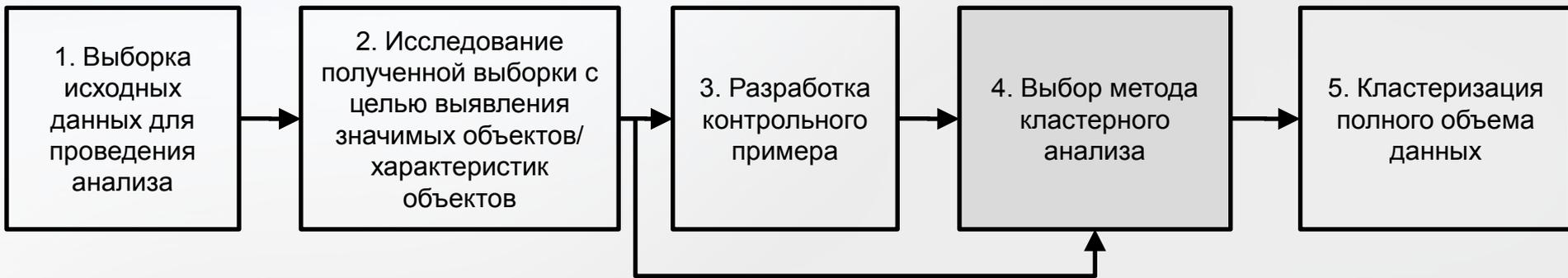
## Назначение

Предназначена для аналитиков, специалистов по анализу данных, разработчиков систем класса Data Mining

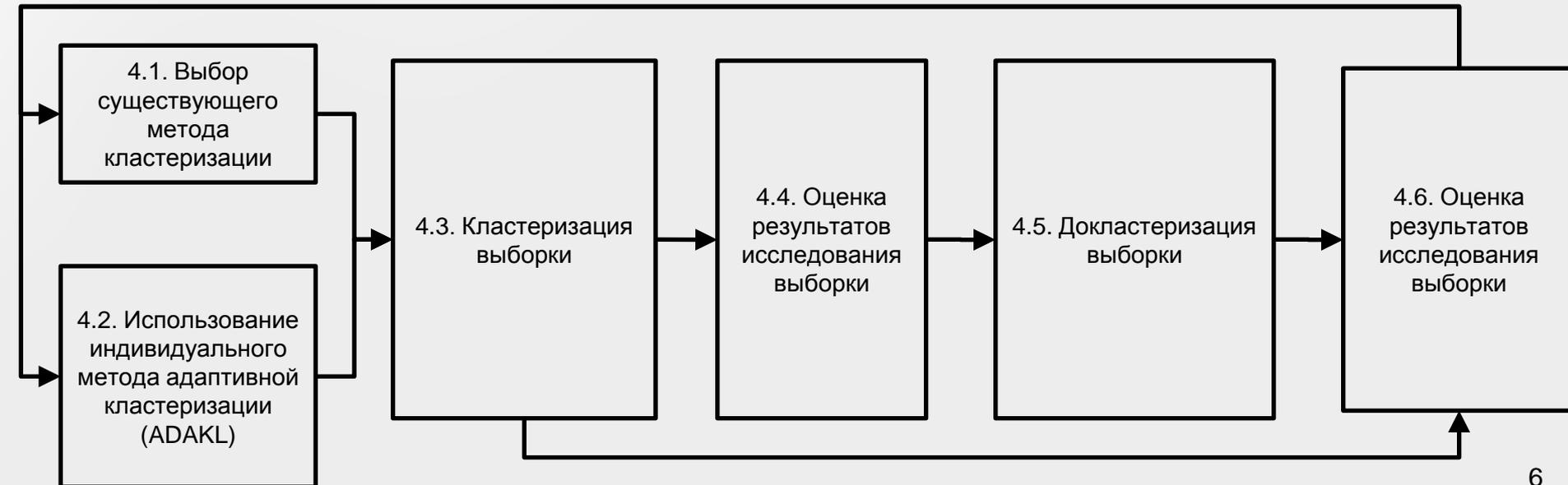
## Задачи

1. Исследование методов и подходов интеллектуального анализа данных, используемых для кластеризации фактографических данных.
2. Разработка методики адаптивной кластеризации фактографических данных.
3. Разработка рекомендаций по выбору существующих методов кластеризации.
4. Разработка метода адаптивной кластеризации фактографических данных.
5. Разработка метода докластеризации.
6. Разработка программного комплекса для автоматизации этапов метода кластеризации.
7. Оценка эффективности разработанного метода с помощью экспериментальных исследований.

# Методика адаптивной кластеризации



## Выбор метода кластерного анализа



# Выбор существующего метода кластеризации

## Шаг 1. Выбор метода кластерного анализа

1. На основе существующих рекомендаций по исследованию предметных областей и задач.
2. На основе критериев
3. Общий алгоритм

## Шаг 2. Настройка параметров выбранного метода кластерного анализа

- объем обучающего множества;
- объем валидационного множества;
- объем тестового множества;
- количество атрибутов входного набора данных;
- тип атрибутов входного набора данных;
- используемость атрибутов входного набора данных.

Характеристические параметры

- количество кластеров;
- алгоритм выполнения дополнительной кластеризации;
- пороговое значение остановки работы алгоритма;
- способ выбора начальных центров;
- максимальное количество итераций;
- количество одновременно обрабатываемых данных;
- количество предварительных разделов;
- коэффициент удаленности.

Итерационные параметры

- способ определения расстояния между кластерами;
- метод оценки качества кластеризации;
- пороговое значение для метода оценки качества кластеризации;
- начальное пороговое значение алгоритма;
- процент аномалий (выбросов) в полном объеме;
- разделяющая функция;
- скорость обучения сети.

Экспертные параметры

## Шаг 3. Анализ массива фактографических данных и оценка разбиения

Кластеризация массива

Аналитическая оценка

Индекс оценки<sub>1</sub>

Индекс оценки<sub>2</sub>

...

Индекс оценки<sub>k</sub>

Результат

Качество?

Да

Нет

# Выбор метода кластеризации на основе рекомендаций

| Метод                   | CURE  | BIRCH  | MST  | k-средних  | PAM  | CLOPE   | SOM  | HCM  | Fuzzy C-Means   |
|-------------------------|---|--|--|--|--|---|--|--|---|
| <b>Рекомендации</b>     | Выявляет кластеры произвольной формы, метод менее чувствителен к выбросам, чем MST. Время работы алгоритма незначительное | Метод предназначен для очень больших наборов данных. Работает с произвольным количеством оперативной памяти. Получаемое разбиение обладает высоким качеством | Лучше всего подходит для выделения кластеров произвольной формы                    | Показывает хорошие результаты при работе с данными, которые распределены по компактным группам сферической формы | Показывает хорошие результаты при работе с данными, которые распределены по компактным группам сферической формы | Кластеризация больших объемов категорийных данных       | Поиск и анализ закономерностей   | Показывает высокие результаты при работе с данными, которые распределены по компактным группам сферической формы | Относит объект к разным кластерам на основе степени принадлежности элемента к кластерам, выделяет кластеры сферической формы  |
| <b>Противопоказания</b> | Не использовать для исследования объектов с большим количеством атрибутов, требует задания пороговых значений             | Не использовать для получения несферических форм кластеров, требует задания пороговых значений   | Очень чувствителен к выбросам и может медленно работать на больших массивах данных | Очень чувствителен к выбросам и может медленно работать на больших массивах данных                               | Чувствителен к выбросам и может медленно работать на больших массивах данных                                     | Требуется подбор оптимального коэффициента отталкивания | Требуется минимизация размеров карты, проблема с аналитическим обоснованием результатов исследования | Чувствителен к выбросам и может медленно работать на больших массивах данных                                     | Высокие требования к вычислительной мощности используемого аппаратного обеспечения, не работает с объектами, которые удалены от всех кластеров, и с вложенными кластерами |

# Критерии выбора методов кластеризации

|                       |  |
|-----------------------|--|
| Критерий 1 ( $Cr_1$ ) | <i>Объем информации</i>  |
| Критерий 2 ( $Cr_2$ ) | <i>Размерность информации</i>  |
| Критерий 3 ( $Cr_3$ ) | <i>Типы атрибутов</i>  |
| Критерий 4 ( $Cr_4$ ) | <i>Чувствительность к равномерности информации</i>                       |
| Критерий 5 ( $Cr_5$ ) | <i>Априорное (экспертное) представление о форме получаемых кластеров</i> |
| Критерий 6 ( $Cr_6$ ) | <i>Априорное (экспертное) представление о количестве кластеров</i>       |
| Критерий 7 ( $Cr_7$ ) | <i>Априорное (экспертное) представление о перекрываемости кластеров</i>  |
| Критерий 8 ( $Cr_8$ ) | <i>Качество кластеризации</i>  |

# Выбор метода кластеризации на основе критериев

|   |
|---|
| <b>Если</b> Критерий 6 = «Неизвестно» и Критерий 8 <> «Высокое»   |
| <b>То</b> Выдать сообщение «Метод кластеризации отсутствует»  |
| <b>Или Если</b> Критерий 8 = «Низкое» и Критерий 7 = «Без пересечений» и Критерий 5 = «Сферическая форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 = «Низкая размерность» |
| <b>То</b> Метод := «к-средних»  |
| <b>Или если</b> Критерий 8 = «Среднее»  |
| <b>Начало блока 1</b>   |
| <b>Если</b> Критерий 7 = «С пересечениями» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 <= «Высокая размерность»                               |
| <b>То</b> Метод := Fuzzy C-means  |
| <b>Или если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сферическая форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 = «Невысокая размерность»                      |
| <b>То</b> Метод := VIRCH  |
| <b>Или если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные неравномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 <= «Высокая размерность»                         |
| <b>То</b> Метод := Самоорганизующиеся карты Кохонена  |
| <b>Или если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 <= «Средняя размерность» <b>То</b>                 |
| <b>То</b> Метод := HCM  |

|   |
|---|
| <b>Или если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 = «Низкая размерность» и Критерий 1 = «Небольшой объем данных»             |
| <b>То</b> Метод := PAM  |
| <b>Или</b>  |
| Выдать сообщение «Метод кластеризации отсутствует»  |
| <b>Конец блока 1</b>  |
| <b>Или если</b> Критерий 8 = «Высокое»  |
| <b>Начало блока 2</b>   |
| <b>Если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные неравномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 = «Низкая размерность»   |
| <b>То</b> Метод := CURE   |
| <b>Или Если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 <= «Средняя размерность» и Критерий 1 <= «Большой объем данных»            |
| <b>То</b> Метод := MST  |
| <b>Или Если</b> Критерий 7 = «Без пересечений» и Критерий 5 = «Сложная форма» и Критерий 4 = «Данные равномерны» и Критерий 3 = «Числовые атрибуты» и Критерий 2 <= «Средняя размерность» и Критерий 1 <= «Огромный объем данных» <b>То</b> |
| <b>То</b> Метод := CLOPE  |
| <b>Или</b>  |
| Выдать сообщение «Метод кластеризации отсутствует»  |
| <b>Конец блока 2</b>  |
| <b>Или</b>  |
| Выдать сообщение «Метод кластеризации отсутствует»  |

# Метод адаптивной кластеризации ADAKL



# Описание метода ADAKL

## Оценка информационного расстояния между объектами

$$Dist_{ij} = \|u_i - u_j\| = Metric(u_i, u_j)$$

$$Metric = \frac{\sum_l (|Value_{il} - Value_{jl}| * K_l)}{\sum_l K_l}$$

$$Dist_{ij} = \sqrt{\sum_l (|Value_{il} - Value_{jl}| * K_l)^2}$$

$$Metric = \frac{\sum_l (|Value_{il} - Value_{jl}| * K_l)}{\sum_l K_l}$$

$$Dist_{ij} = \sum_l (|Value_{il} - Value_{jl}| * K_l)^2$$

$$Metric = \frac{\sum_l (|Value_{il} - Value_{jl}| * K_l)}{\sum_l K_l}$$

$$Dist_{ij} = Max_l (|Value_{il} - Value_{jl}| * K_l)$$

$$K = \{K_1, K_2, \dots, K_n\} \quad K_i \in [0;1]$$

## Первичное разделение на кластеры

$$Dist_{ij}^l := \begin{cases} 0 & | \\ Dist_{ij}^l = Max & , i, j = \overline{1, m}, l = \overline{1, k} \end{cases}$$

## Расчет центров кластеров

$$V_i^k = Avg \left( \{u_j \mid u_j \in C_i^k\} \right)$$

$$Avg[r] = \frac{\sum_{u_j \in V_i^k} \{Value_{jr} \mid FieldType[w] = "Attribute\} \quad i = \overline{1, k} \quad j = \overline{1, m}}{|C_i^k|} \quad r = \overline{1, n} \quad l = \overline{1, m}$$

$$Avg[r] = \left\{ \begin{array}{l} Value_{jr} \quad \mid \quad FieldType[w] = "Attribute" \\ \sum_{\substack{u_j \in C_i^k, u_l \in C_l^k \\ \|Value_{jr} - Value_{lr}\| = \min \quad \varphi = \max}} \end{array} \right\}$$

## Отнесение объектов к кластерам

$$u_j \in V_i^k \mid Dist_{ij}^{k'} \leq w \quad Dist_{ij}^{k'} = Min_i (Dist_{ij}^{k'}) \quad i = \overline{1, k} \quad j = \overline{1, m}$$

## Расчет центров кластеров (2)

$$V_i^{k'} = \frac{\sum_{j=1}^m \mu_{ij}^p * u_j}{\sum_{j=1}^m \mu_{ij}^p}, i = \overline{1, k}$$

$$V_i^{k'}[r] = Value_{jr} \quad \mu_{ij} = Max(\mu_{ij})$$

## Степень принадлежности

$$\mu_{ij} = (1 - Dist_{ij}^{k'})^2 \quad i = \overline{1, k} \quad j = \overline{1, m}$$

# Выбор лучшего разбиения

$$O^k = \frac{\sum_{i=1, k} \left( \underbrace{|V_i^{k'}|}_{1} * \sum_{j=1}^m \underbrace{\mu_{ij}^p * \|V_i^{k'} - u_j\|}_{3} \right)}{\underbrace{\text{Min}_{i \neq j} \left( \|V_i^{k'} - u_j\| \right) * \text{Ma}_{u_j \in V_i^{k'}} \left( \|V_i^{k'} - u_j\| \right) * \sum_{j=1}^m \|V_i^{k'} - u_j\| * k}_{2 \quad m * k^2}}$$

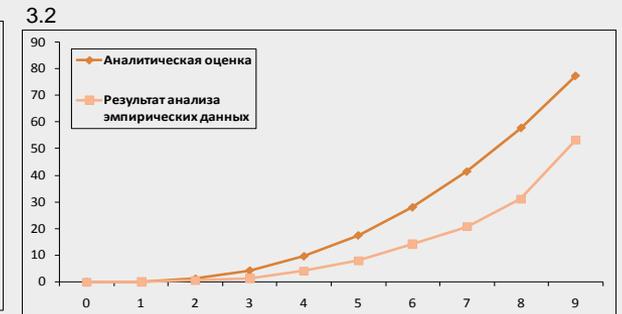
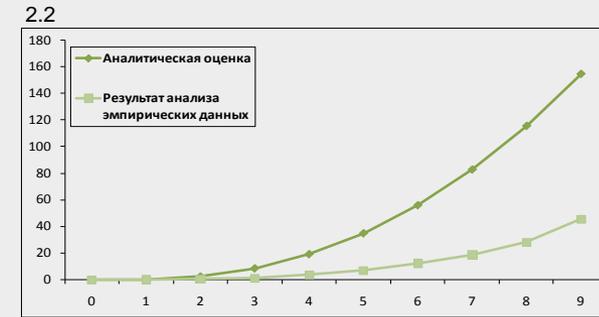
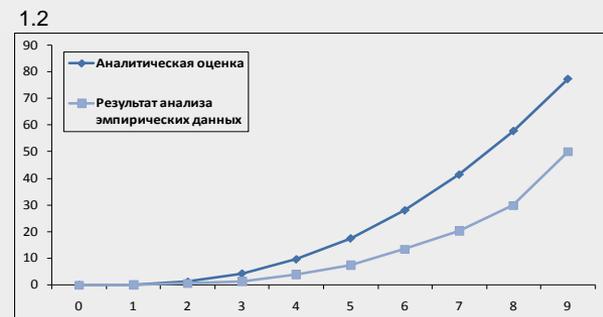
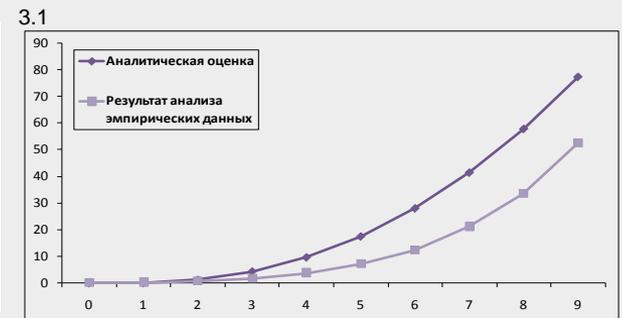
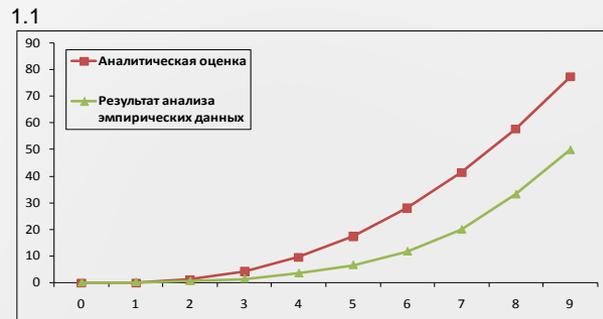
$$O_{\hat{i} \ddot{i} \dot{o}} = \text{MAX}_{i=1, q} \left( O^i \right)$$

# Аналитическая оценка сложности метода

| Алгоритм построения MST     | Алгоритм Борувки               | Алгоритм Крускала                | Алгоритм Прима                 |
|-----------------------------|--------------------------------|----------------------------------|--------------------------------|
| Способ нормализации         |                                |                                  |                                |
| Линейная нормализация       | $O(m^2 * [a + b + Lg(m) + q])$ | $O(m^2 * [a + b + Lg(m^2) + q])$ | $O(m^2 * [a + b + Lg(m) + q])$ |
| Статистическая нормализация | $O(m^2 * [a + b + Lg(m) + q])$ | $O(m^2 * [a + b + Lg(m^2) + q])$ | $O(m^2 * [a + b + Lg(m) + q])$ |

## Итоговая аналитическая оценка метода

$$O(m^2 * [a + b + Lg(m) + q]) ; O(m^2 * [a + b + Lg(m^2) + q])$$



# Исследуемые наборы данных

| Характеристика<br>Исследование | Общее количество<br>исходных данных<br>(шт. записей) | Количество<br>атрибутов числового<br>типа [количество<br>значимых атрибутов]<br>(шт.) | Количество<br>атрибутов текстового<br>атрибута [количество<br>значимых атрибутов]<br>(шт.) | Общее количество<br>атрибутов [общее<br>количество значимых<br>атрибутов] (шт.) |
|--------------------------------|--|---|--|---|
| 1                              | 267  | 0 [0]   | 5 [2]  | 5 [2]   |
| 2                              | 533  | 73 [72]   | 3 [3]  | 76 [75]   |
| 3                              | 267  | 5 [2]   | 4 [2]  | 9 [4]   |
| 4                              | 450  | 3 [3]   | 2 [1]  | 5 [4]   |

| SecBrief              | SecType   | SecBranch    | Nominal    | Volume            |
|-----------------------|-----------|--------------|------------|-------------------|
| 46014                 | Облигации | Госзайм      | 1 000.00   | 10 000 000 000.00 |
| 46018                 | Облигации | Госзайм      | 1 000.00   | 10 000 000.00     |
| 48001                 | Облигации | Госзайм      | 1 000.00   | 24 099 483.00     |
| DCL ASSETS Int. Ltd.  | Облигации | ЗапФинСектор | 100 000.00 | 1 380.00          |
| IMPEX07               | Облигации | ЗапФинСектор | 2 000.00   | 100 000 000.00    |
| RusFed05              | Облигации | ВнешГосзайм  | 1 000.00   | 1 000 000.00      |
| RusFed18              | Облигации | ВнешГосзайм  | 1 000.00   | 1 000 000.00      |
| RusFed30              | Облигации | ВнешГосзайм  | 1.00       | 1 000 000.00      |
| SISFIN08              | Облигации | ЗапФинСектор | 1 000.00   | 350 000.00        |
| ГМК Норникель АО в.5  | Акции     | Металлургия  | 1.00       | 70 000 000.00     |
| Газэнергпромпанк АОИ  | Акции     | Банк         | 1.00       | 100 000 000.00    |
| Казаньоргсинтез АПИ 2 | Акции     | ХимПром      | 1.00       | 25 000 000.00     |
| Конкордия-АВАНТАЖ     | Акции     | ЖилСектор    | 0.00       | 3 000.00          |

| Vehicle | TopSpeed | Colour | AirResistance | Weight   | Type          |
|---------|----------|--------|---------------|----------|---------------|
| V1      | 220      | red    | 0.30          | 1 300.00 | Sport         |
| V2      | 230      | black  | 0.32          | 1 400.00 | Sport         |
| V3      | 260      | red    | 0.29          | 1 500.00 | Sport         |
| V4      | 140      | grey   | 0.35          | 800.00   | Medium market |
| V5      | 155      | blue   | 0.33          | 950.00   | Medium market |
| V6      | 130      | white  | 0.40          | 600.00   | Medium market |
| V7      | 100      | black  | 0.5           | 3000     | Lorry         |
| V8      | 105      | red    | 0.60          | 2 500.00 | Lorry         |
| V9      | 110      | grey   | 0.55          | 3 500.00 | Lorry         |

| ClientID    | ClientType            | BuyCnt | SellCnt | BuyTurn | SellTurn | Freq    | FinInstrCnt | SecGroupType |
|-------------|-----------------------|--------|---------|---------|----------|---------|-------------|--------------|
| Клиент1 - 1 | Долгосрочный инвестор | 2      | 0       | 2000    | 0        | 0.001   | 3           | Акции        |
| Клиент2 - 2 | Спекулянт             | 10     | 10      | 150     | 200      | 1       | 9           | Акции        |
| Клиент3 - 3 | Потребитель           | 0      | 3       | 0       | 30       | 0.01    | 2           | Акции        |
| Клиент4 - 4 | Паевой фонд           | 25     | 25      | 28000   | 32000    | 1       | 15          | Акции        |
| Клиент5 - 5 | Акционер              | 1      | 0       | 8000    | 0        | 0.00001 | 1           | Акции        |
| Клиент6 - 1 | Долгосрочный инвестор | 3      | 0       | 1355    | 0        | 0.003   | 2           | Акции        |
| Клиент7 - 2 | Спекулянт             | 18     | 8       | 570     | 560      | 0.999   | 5           | Акции,       |

# Экспериментальная оценка метода адаптивной кластеризации

Исследование 1: выделение секторов инвестирования на основе анализа показателей финансовых инструментов;

Исследование 2: выделение групп клиентов на основе статистических данных о деятельности клиентов за период;

Исследование 3: выявление категорий финансовых инструментов для оценки эффективности операций;

Исследование 4: выделение классов автомобилей на основе данных об их характеристиках.

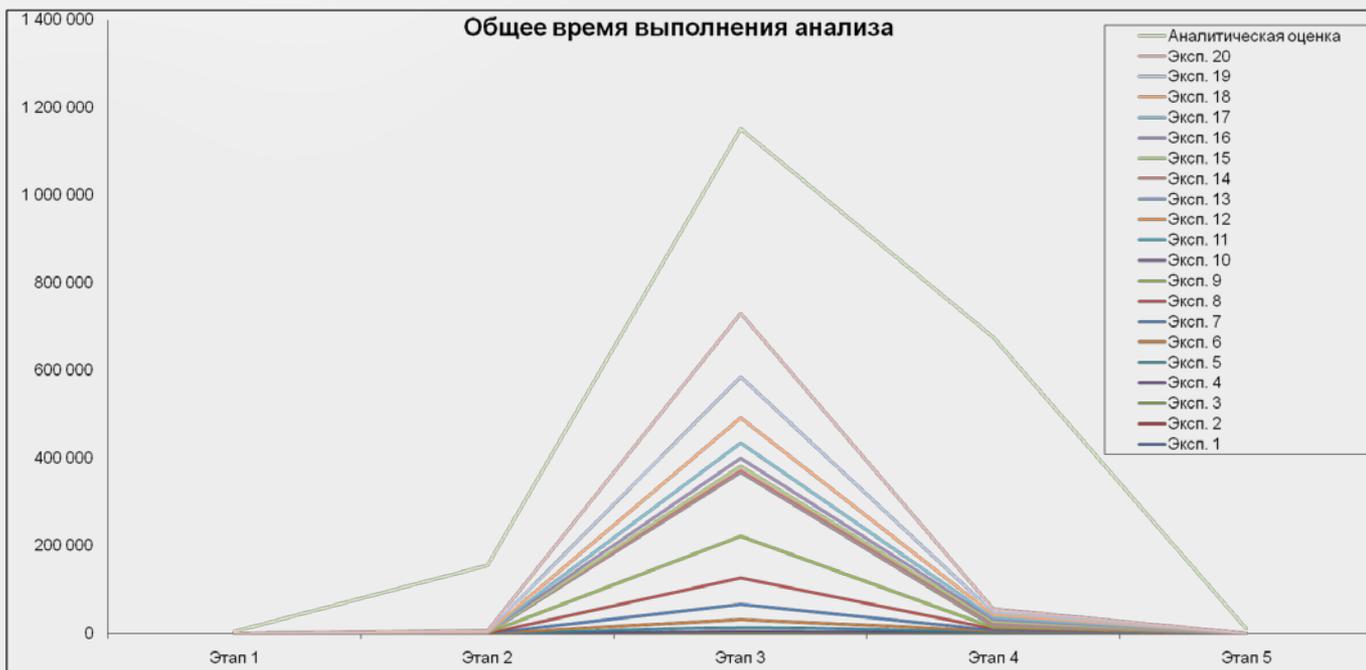
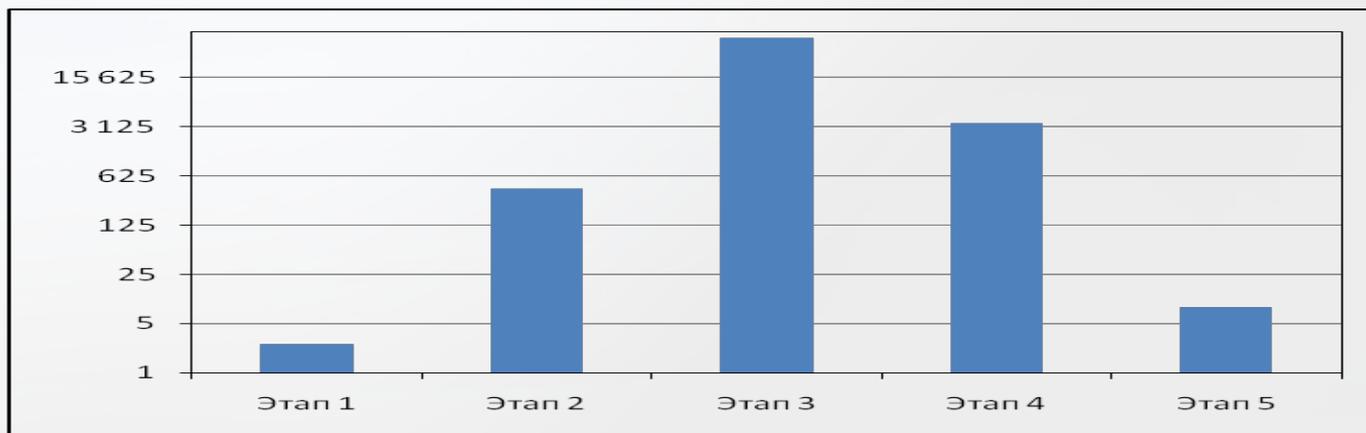
| Показатель<br>Оценка | Средневзвешенная<br>оценка разбиения | Средневзвешенная<br>оценка разбиения с<br>заданным<br>количеством<br>кластеров (без учета<br>лингвистических<br>атрибутов) | Средневзвешенная<br>оценка разбиения с<br>заданным<br>количеством<br>кластеров (с учетом<br>лингвистических<br>атрибутов) | Итоговая<br>оценка |
|----------------------|--------------------------------------|--|---|--------------------|
| Метод 1              | 0.7913                               | 0.9150   | 0.9237  | 0.8767             |
| Метод 2              | -                                    | 0.8232   | -   | 0.8232             |
| Метод 3              | 0.9762                               | 0.9981   | 0.9990  | 0.9911             |

Метод 1 – самоорганизующиеся карты Кохонена

Метод 2 – метод k-средних

Метод 3 – метод ADAKL

# Общее время анализа по этапам



# Метод докластеризации

Набор данных 1

ADAKL

**Этап 1.** Нормализация значений числовых атрибутов

**Этап 2.** Вычисление матрицы взаимных расстояний между объектами

**Этап 3.** Построение минимального остовного дерева

**Этап 4.** Разделение объектов на кластеры и построение матрицы нечеткого разбиения

**Этап 5.** Выбор наилучшего разбиения

Итоговые кластеры 1

Набор данных 2

Докластеризация

**Этап 1.** Нормализация значений числовых атрибутов с использованием метода основного алгоритма

**Этап 2.** Расчет оценочной функции для набора данных 1

**Этап 3.** Расчет оценочной функции для набора данных 2

**Этап 4.** Оценка разницы оценочной функции наборов 1, 2

НЕТ

Входит в доверительный интервал

ДА

**Этап 5.** Расчет ближайших элементов для набора данных 2 из набора данных 1

**Этап 6.** Распределение элементов набора данных 2 по кластерам в соответствии с кластерами элементов набора данных 1

Итоговые кластеры 2

Оценка исходных наборов данных

$$D_1 = \sqrt{\frac{\sum_{i=1}^r \|A_i - Avg[A]\|^2}{r}}$$

$$D_2 = \sqrt{\frac{\sum_{i=1}^o \|B_i - Avg[B]\|^2}{o}}$$

$$|D_1 - D_2| \leq \delta$$

Распределение объектов по кластерам

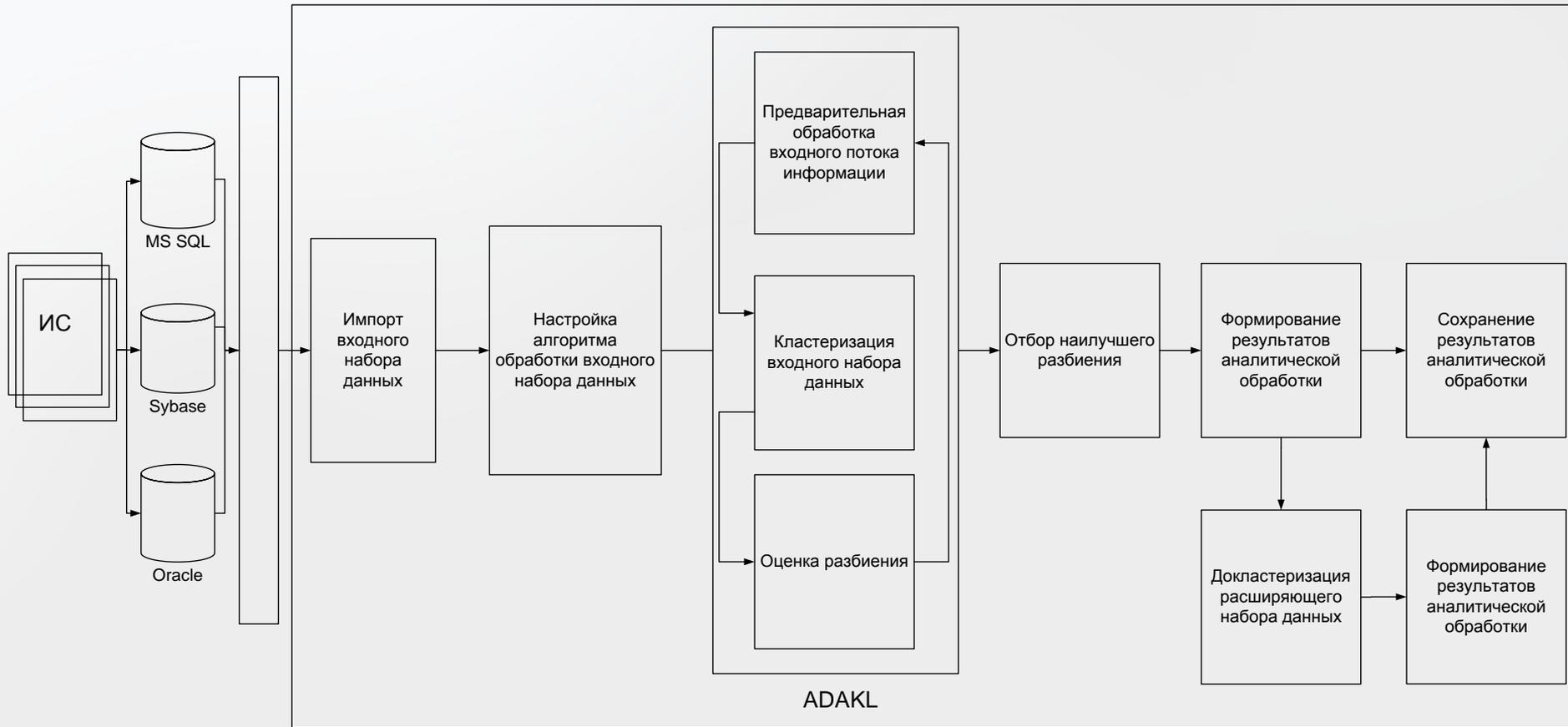
$$\mu_{ij} = (1 - Dist_j^{Norm}) * \mu_{il} =$$

$$\left[ 1 - Dist_{lj} / \text{Max}(Dist_{sj}) \right] * \mu_{il}$$

# Программное решение (1)

Подготовка входного набора данных

Этапы функционирования программного решения



# Программное решение (2)

fmDataQuery

```
select r trin(Brief) SecBrief, case when SecType = 0 then 'Акции' else 'Облигации' end SecType, Branch SecBr
```

| SecBrief             |
|----------------------|
| 46014                |
| 46018                |
| 48001                |
| DCL ASSETS Int. Ltd. |
| IMPEX07              |
| RusFed05             |
| RusFed18             |
| RusFed30             |
| SISFIN08             |
| ГМК Норникель АО в.5 |
| Газэнергпромпанк АОИ |

Отменить    Очистить текст    Выполнить запрос    Импорт данных

fmConfigDataSet

PortfolioID

Вид значения атрибута:

Тип значения атрибута:

Наименование атрибута:

Весовой коэффициент атрибута:

Макс. количество кластеров:

Размазанность кластеров:

Степень удаленности элементов:

Способ определения расстояния:

Способ построения мин. остовного дерева:

Очистить виды/типы значений    Закрыть

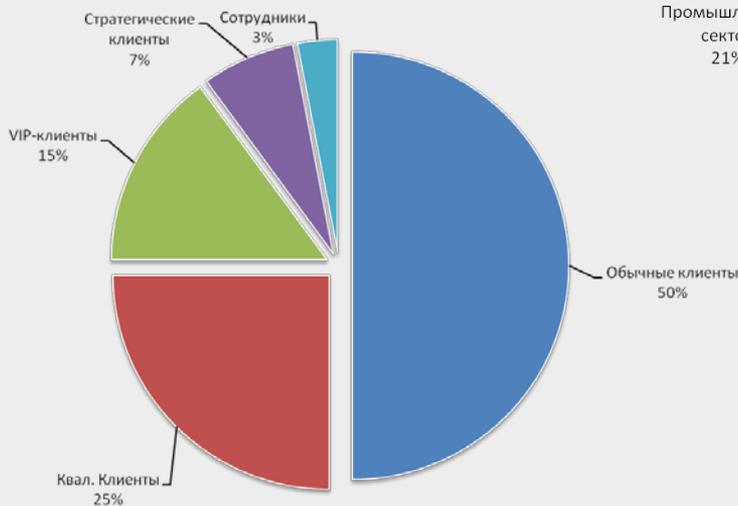
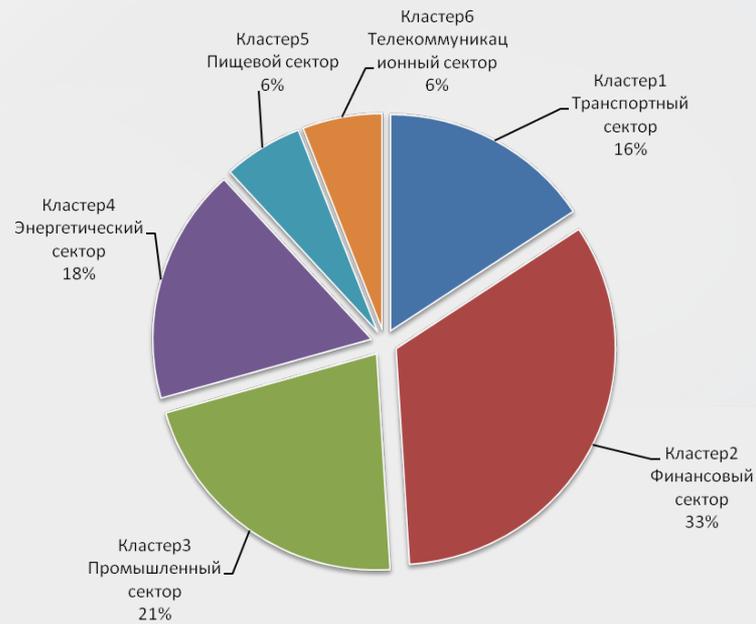
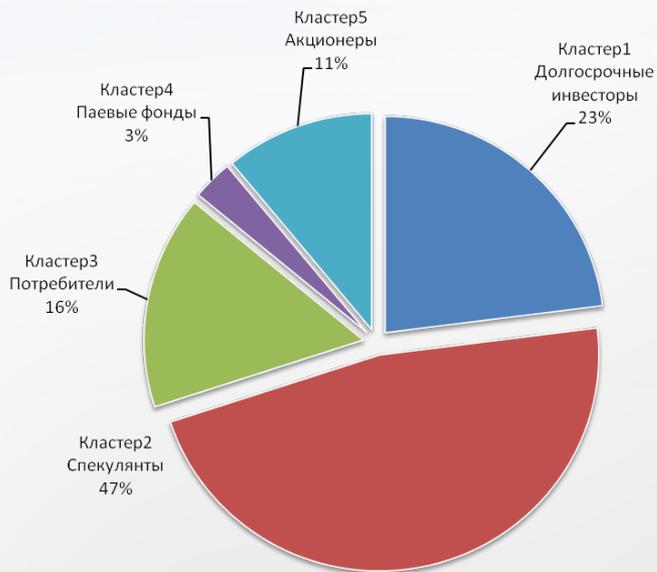
fmResults

| N  | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       |
|----|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 1  | 0.000000 | 2.351880 | 1.452179 | 1.378286 | 1.415267 | 1.706679 | 1.376950 | 1.405524 | 1.290798 | 1.379021 | 1.604524 |
| 2  | 2.351880 | 0.000000 | 1.397319 | 1.710017 | 1.560639 | 1.157141 | 1.721131 | 1.599861 | 1.640260 | 1.995235 | 1.093818 |
| 3  | 1.452179 | 1.397319 | 0.000000 | 0.340590 | 0.225283 | 0.702109 | 0.380000 | 0.261687 | 0.272932 | 0.674564 | 0.326712 |
| 4  | 1.378286 | 1.710017 | 0.340590 | 0.000000 | 0.153178 | 0.838709 | 0.118118 | 0.117207 | 0.074377 | 0.525762 | 0.631212 |
| 5  | 1.415267 | 1.560639 | 0.225283 | 0.153178 | 0.000000 | 0.717327 | 0.197978 | 0.043005 | 0.090110 | 0.629761 | 0.485712 |
| 6  | 1.706679 | 1.157141 | 0.702109 | 0.838709 | 0.717327 | 0.000000 | 0.784366 | 0.735816 | 0.794204 | 1.175497 | 0.606812 |
| 7  | 1.376950 | 1.721131 | 0.380000 | 0.118118 | 0.197978 | 0.784366 | 0.000000 | 0.158302 | 0.153434 | 0.521858 | 0.667812 |
| 8  | 1.405524 | 1.599861 | 0.261687 | 0.117207 | 0.043005 | 0.735816 | 0.158302 | 0.000000 | 0.064215 | 0.607118 | 0.526812 |
| 9  | 1.290798 | 1.640260 | 0.272932 | 0.074377 | 0.090110 | 0.754204 | 0.153434 | 0.064215 | 0.000000 | 0.558540 | 0.559812 |
| 10 | 1.379021 | 1.995235 | 0.674564 | 0.525762 | 0.629761 | 1.175497 | 0.521858 | 0.607118 | 0.558540 | 0.000000 | 0.390212 |
| 11 | 1.604524 | 1.093818 | 0.326712 | 0.631212 | 0.485712 | 0.606812 | 0.667812 | 0.526812 | 0.559812 | 0.990377 | 0.000000 |
| 12 | 1.425326 | 1.449742 | 0.211072 | 0.301091 | 0.179705 | 0.642615 | 0.322509 | 0.206101 | 0.241446 | 0.735516 | 0.398412 |
| 13 | 1.445987 | 2.095887 | 0.868575 | 0.786858 | 0.873640 | 1.358553 | 0.784318 | 0.857124 | 0.810433 | 0.289083 | 1.164212 |
| 14 | 1.893646 | 0.709880 | 0.778562 | 1.097652 | 0.954913 | 0.873180 | 1.140176 | 0.996926 | 1.026500 | 1.411270 | 0.473212 |
| 15 | 1.372492 | 1.745585 | 0.393929 | 0.083295 | 0.200730 | 0.826546 | 0.049943 | 0.159399 | 0.138973 | 0.509581 | 0.680812 |
| 16 | 1.440223 | 1.456973 | 0.215610 | 0.276733 | 0.151324 | 0.578351 | 0.262430 | 0.171995 | 0.222147 | 0.680319 | 0.424212 |
| 17 | 1.671776 | 0.956938 | 0.477264 | 0.762191 | 0.614054 | 0.433434 | 0.765025 | 0.650903 | 0.694781 | 1.093674 | 0.223612 |
| 18 | 1.694128 | 1.097045 | 0.320350 | 0.635061 | 0.492094 | 0.632232 | 0.675910 | 0.534148 | 0.563925 | 0.981391 | 0.046412 |
| 19 | 1.599120 | 1.092371 | 0.371607 | 0.633649 | 0.467133 | 0.400335 | 0.630897 | 0.522257 | 0.560547 | 0.978668 | 0.220212 |
| 20 | 1.479542 | 1.379953 | 0.269952 | 0.394912 | 0.273057 | 0.455976 | 0.292327 | 0.345273 | 0.765577 | 0.409812 | 0.000000 |
| 21 | 1.384452 | 2.043094 | 0.778220 | 0.673768 | 0.767307 | 1.277169 | 0.670193 | 0.740453 | 0.700356 | 0.218676 | 1.005212 |
| 22 | 2.203615 | 0.794917 | 1.330308 | 1.551403 | 1.417659 | 0.744983 | 1.507670 | 1.443568 | 1.497066 | 1.825021 | 1.116412 |

fmResults

| N  | Кластер 1 | Кластер 2 | Кластер 3 | Кластер 4 | Кластер 5 |
|----|-----------|-----------|-----------|-----------|-----------|
| 1  | 0.000587  | 0.002954  | 0.000141  | 0.000181  | 0.000056  |
| 2  | 0.007476  | 0.002153  | 0.000209  | 0.000184  | 0.000055  |
| 3  | 0.000850  | 0.001319  | 0.000000  | 0.000000  | 0.000000  |
| 4  | 0.002345  | 0.010523  | 0.000196  | 0.000535  | 0.000174  |
| 5  | 0.002896  | 0.005899  | 0.000462  | 0.000430  | 0.000135  |
| 6  | 0.001515  | 0.007819  | 0.000029  | 0.000528  | 0.000094  |
| 7  | 0.000193  | 0.000163  | 0.016120  | 0.007972  | 0.007444  |
| 8  | 0.000197  | 0.000420  | 0.000796  | 0.016230  | 0.006373  |
| 9  | 0.000000  | 0.000000  | 0.006196  | 0.005189  | 0.016919  |
| 10 | 0.000587  | 0.002954  | 0.000141  | 0.000181  | 0.000056  |
| 11 | 0.007476  | 0.002153  | 0.000209  | 0.000184  | 0.000055  |
| 12 | 0.000850  | 0.001319  | 0.000000  | 0.000000  | 0.000000  |
| 13 | 0.002345  | 0.010523  | 0.000196  | 0.000535  | 0.000174  |
| 14 | 0.002896  | 0.005899  | 0.000462  | 0.000430  | 0.000135  |
| 15 | 0.001515  | 0.007819  | 0.000029  | 0.000528  | 0.000094  |
| 16 | 0.000193  | 0.000163  | 0.016120  | 0.007972  | 0.007444  |
| 17 | 0.000197  | 0.000420  | 0.000796  | 0.016230  | 0.006373  |
| 18 | 0.000000  | 0.000000  | 0.006196  | 0.005189  | 0.016919  |
| 19 | 0.000587  | 0.002954  | 0.000141  | 0.000181  | 0.000056  |
| 20 | 0.007476  | 0.002153  | 0.000209  | 0.000184  | 0.000055  |
| 21 | 0.000850  | 0.001319  | 0.000000  | 0.000000  | 0.000000  |
| 22 | 0.002345  | 0.010523  | 0.000196  | 0.000535  | 0.000174  |

# Практическое исследование



# Основные результаты работы

1. Проведено исследование существующих методов и подходов интеллектуального анализа данных, используемых для кластеризации фактографических данных.
2. Проведен анализ аналитических программных комплексов с выделением назначения программного комплекса и основных функциональных возможностей.
3. Разработана общая методика адаптивной кластеризации, которая состоит из пяти этапов: выборка исходных данных, исследование полученной выборки с целью выявления значимых для разбиения характеристик, разработка контрольного примера, выбор метода кластеризации, кластеризации полного объема данных.
4. Для выбора метода кластеризации на основе литературных источников выделено восемь критериев.
5. Разработан критерий для оценки качества разбиения, который позволяет проводить оценку и сравнение результатов исследований на основе сравнения итоговых и ожидаемых количественных показателей разбиения.
6. Разработан метод адаптивной кластеризации (ADAKL) на основе интеграции методов минимального остовного дерева и нечетких K-средних, определяющий количество кластеров на основе локального критерия, обладающий двухэтапностью, восемью входными параметрами настройки, нечеткостью при распределении объектов по кластерам, возможностью использования объектов с разными типами атрибутов, приемлемым временем работы и конечностью результата.
7. Разработан локальный критерий оценки разбиения множества на кластеры, который учитывает характеристики практической задачи, лежащей в основе научного исследования: выделение кластеров с наименьшими взаимными расстояниями и наибольшим количеством элементов в кластере, минимизация количества кластеров, минимизация взаимных расстояний между получаемыми центрами кластеров и распределяемыми объектами.
8. Разработан метод докластеризации, позволяющий расширять исследованные массивы фактографических данных и уменьшающий затраты времени на проведение исследования за счет выявления взаимных связей между исследованными объектами и добавляемыми объектами.
9. Разработанный метод ADAKL реализован в виде программного решения, который подтверждает аналитическую оценку.
10. На основе программного решения проведены экспериментальные исследования и оценка состоятельности разработанного метода в сравнении с имеющимися методами (k – средние, карты Кохонена).