

МЕТОДИКА АДАПТИВНОЙ КЛАСТЕРИЗАЦИИ

ФАКТОГРАФИЧЕСКИХ ДАННЫХ

Цель работы:

Повышение эффективности процесса кластеризации фактографических данных с использованием методов интеллектуального анализа данных.

Назначение работы:

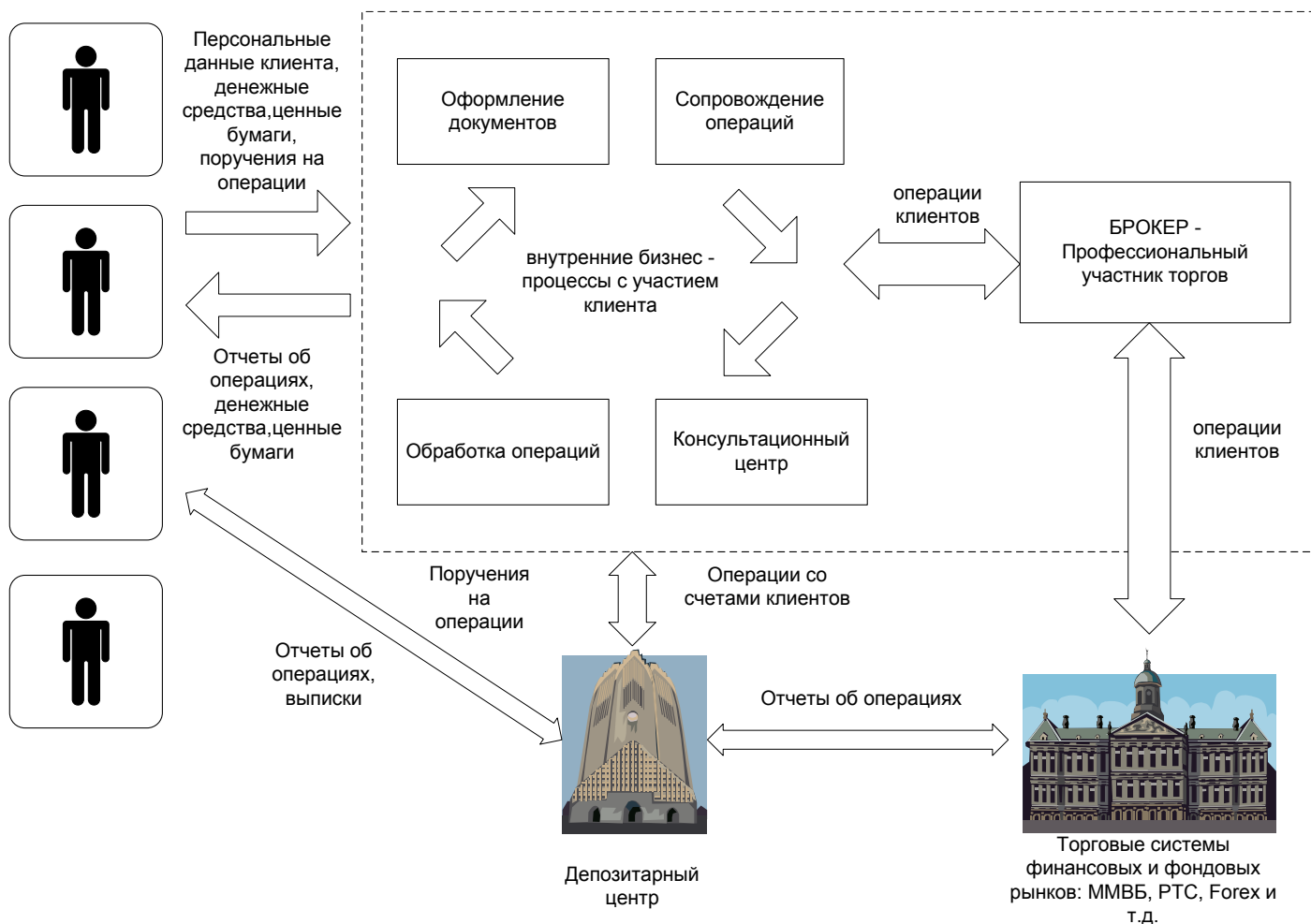
Оптимизация списка услуг брокерского обслуживания для выделенных групп клиентов.

Задачи, подлежащие решению:

1. Исследование методов интеллектуального анализа данных, используемых для кластеризации фактографических данных.
2. Исследование деятельности кредитной организации с целью выявления разновидностей показателей, характеризующих бизнес – процессы.
3. Разработка критериев оценки/выбора исследованных методов кластерного анализа. Обобщение исходных данных алгоритмов с целью оценки возможности адаптации алгоритмов к различным предметным областям.
4. Разработка адаптивного алгоритма кластеризации фактографических данных.
5. Разработка программного комплекса интеллектуального анализа данных.
6. Апробация методов и анализ экспериментальных данных.

ОПИСАНИЕ ПРЕДМЕТНОЙ ОБЛАСТИ

Брокерское обслуживание клиентов

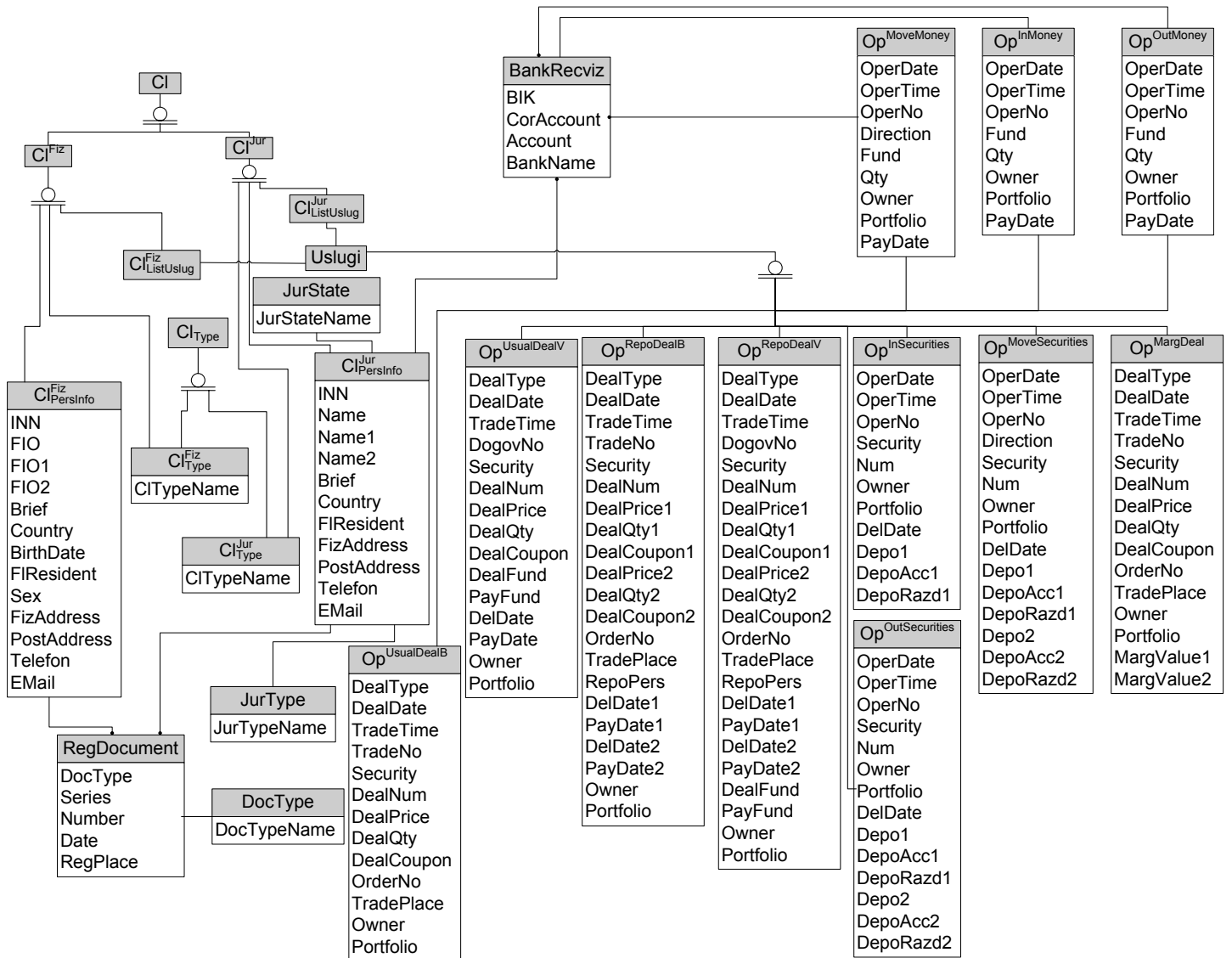


СУЩЕСТВУЮЩИЕ РЕШЕНИЯ И ИХ НЕДОСТАТКИ:

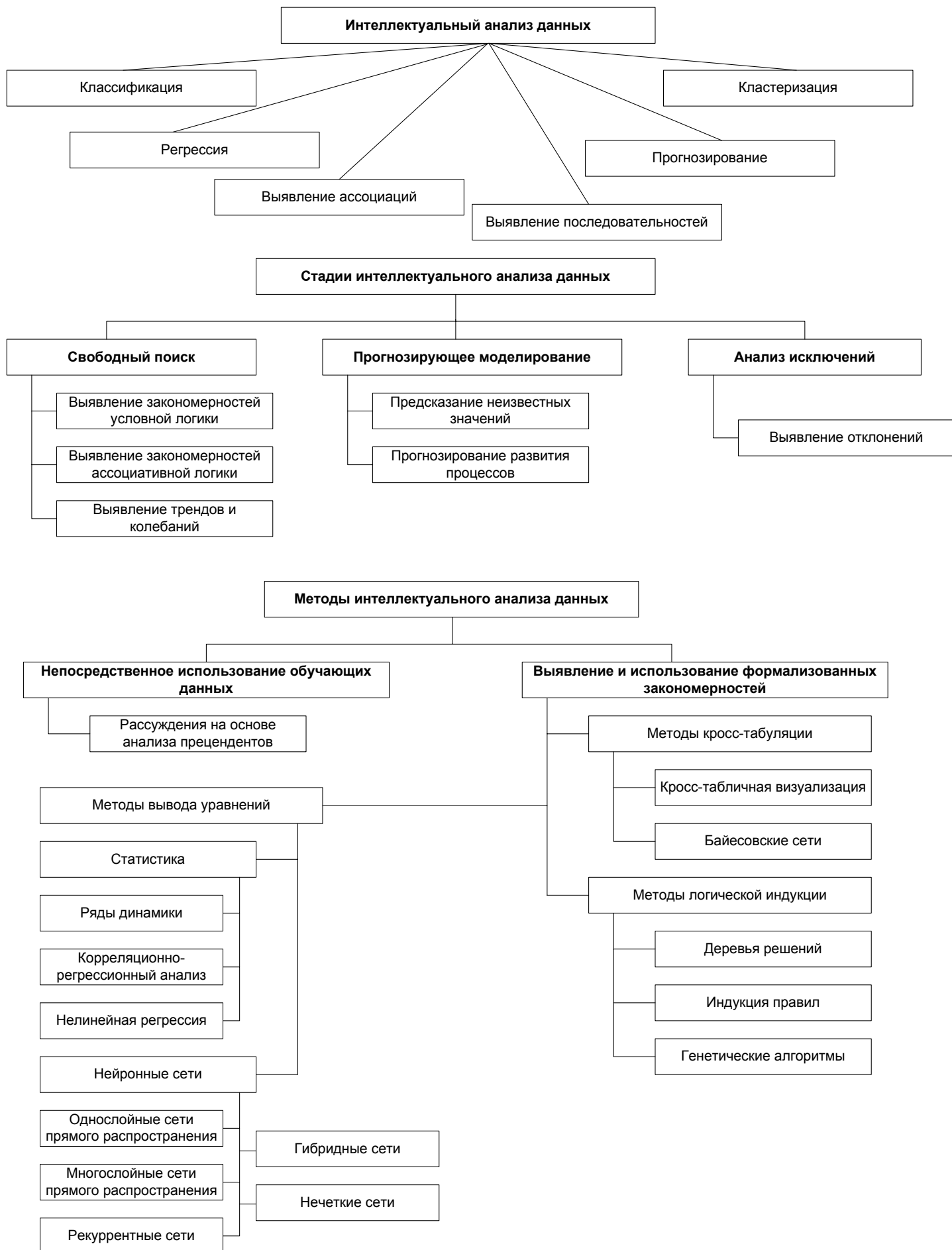
Продукт/ Характеристики	Информационно – аналитическая система «Арион»	Информационно – аналитическая система «Прогноз»	Аналитическая платформа «Deductor»	Business Object Enterprise	R&K Cognos Analyst
Тип системы	Информационно - аналитическая система	Информационно - аналитическая система	Инструментарий для создания на базе аналитической платформы эффективной аналитической системы	Инструментарий для формирования аналитической отчетности	Инструментарий для построения системы принятия решения
Функциональное назначение	Обработка и представление информации на естественном языке, выявление логических взаимосвязей и развернутый поиск	Аналитическая система поддержки принятия решения с использованием имитационных, оптимизационных и статистических методов	Корпоративная отчетность, прогнозирование, сегментация, поиск закономерностей	Корпоративная отчетность	Анализ финансового состояния предприятия
Недостатки	отсутствие комплексного подхода при обработке данных				
	узкая специализация на предметной области небольшое количество реализованных методов интеллектуального анализа данных			узкая специализация на предметной области небольшое количество реализованных методов интеллектуального анализа данных	

ФОРМАЛИЗОВАННАЯ МОДЕЛЬ ПРЕДМЕТНОЙ ОБЛАСТИ

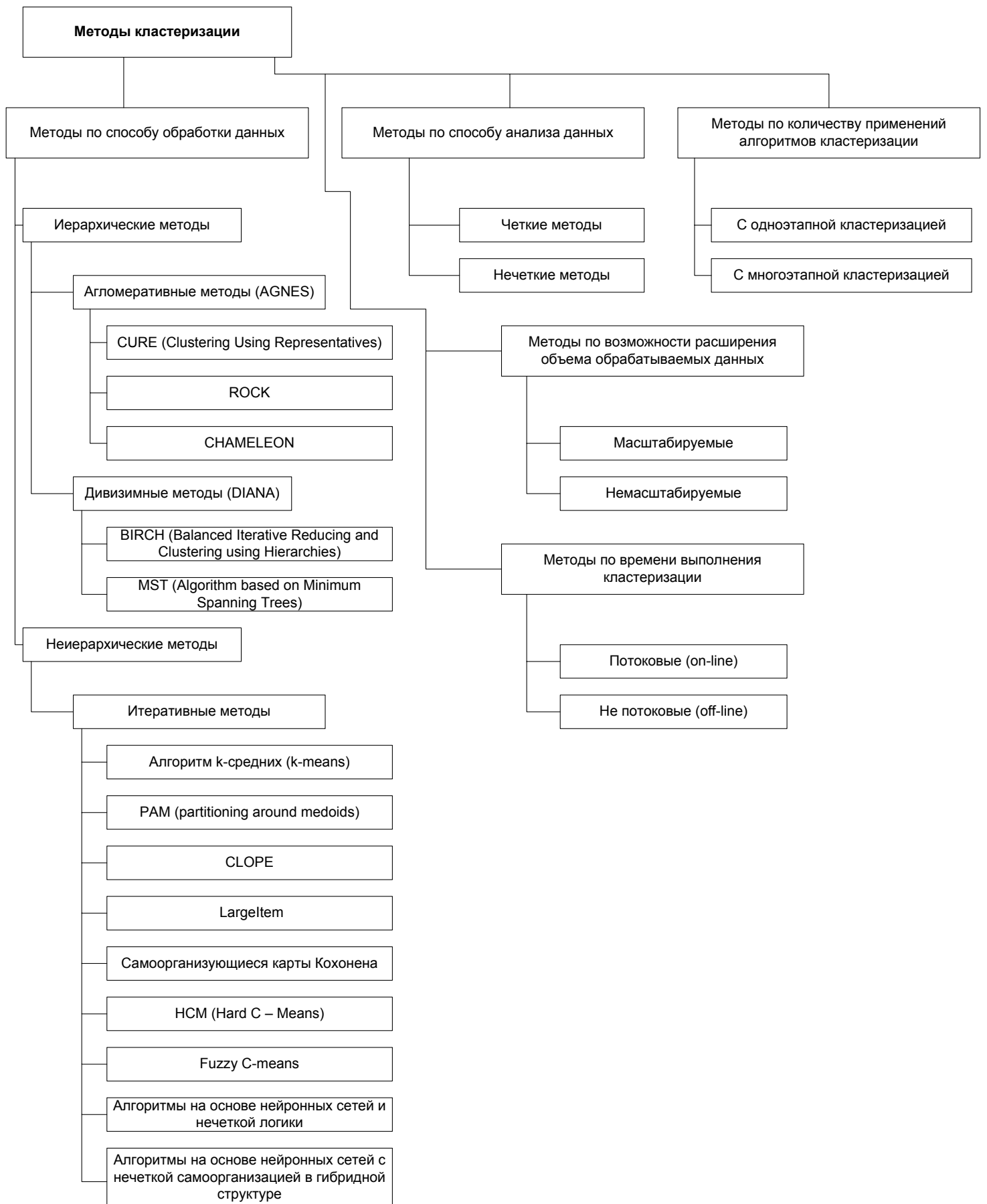
$CI = \{ CI^{Fiz}, CI^{Jur} \}$, $|CI^{Fiz}| \geq 1$, $|CI^{Jur}| \geq 1$, $CI^{Fiz} = \{ CI_{PersInfo}^{Fiz}, CI_{Type}^{Fiz}, CI_{ListUslug}^{Fiz} \}$, $CI^{Jur} = \{ CI_{PersInfo}^{Jur}, CI_{Type}^{Jur}, CI_{ListUslug}^{Jur} \}$
 $CI_{PersInfo}^{Fiz} = \{ INN, FIO, FIO1, FIO2, Brief, Country, BirthDate, RegDocument, FIResident, Sex, FizAddress, PostAddress, Telefon, EMail \}$
 $CI_{PersInfo}^{Jur} = \{ INN, Name, Name1, Name2, Brief, Country, JurType, BankRecviz, JurState, RegDocument, FIResident, FizAddress, PostAddress, \}$
 $RegDocument = \{ DocType, Series, Number, Date, RegPlace \}$, $DocType = \{ \text{Паспорт РФ, Свидетельство о регистрации} \}$, $FIResident = \{ \text{Да, Нет} \}$
 $JurType = \{ \text{Юр.лицо, Банк, РКЦ, Отделение банка} \}$, $Sex = \{ \text{Женский, Мужской} \}$, $BankRecviz = \{ \text{БИК, CorAccount, Account, BankName} \}$
 $JurState = \{ \text{Головная организация, Филиал, Дочерняя организация, Зависимая компания, Подразделение} \}$
 $Uslugi = \{ Op^{UsualDealB}, Op^{UsualDealV}, Op^{RepoDealB}, Op^{RepoDealV}, Op^{InSecurities}, Op^{OutSecurities}, Op^{MoveSecurities}, Op^{MargDeal}, Op^{InMoney}, Op^{OutMoney}, Op^{MoveMoney} \}$
 $Op^{UsualDealB} = \{ DealType, DealDate, TradeTime, TradeNo, Security, DealNum, DealPrice, DealQty, DealCoupon, OrderNo, TradePlace, Owner, Portfolio \}$
 $Op^{UsualDealV} = \{ DealType, DealDate, TradeTime, DogovNo, Security, DealNum, DealPrice, DealQty, DealCoupon, \}$
 $\quad \quad \quad \{ DealFund, PayFund, DelDate, PayDate, Owner, Portfolio \}$
 $Op^{RepoDealB} = \{ DealType, DealDate, TradeTime, TradeNo, Security, DealNum, DealPrice1, DealQty1, DealCoupon1, \}$
 $\quad \quad \quad \{ DealPrice2, DealQty2, DealCoupon2, OrderNo, TradePlace, RepoPers, DelDate1, PayDate1, DelDate2, PayDate2, Owner, Portfolio \}$
 $Op^{RepoDealV} = \{ DealType, DealDate, TradeTime, DogovNo, Security, DealNum, DealPrice1, DealQty1, DealCoupon1, \}$
 $\quad \quad \quad \{ DealPrice2, DealQty2, DealCoupon2, RepoPers, DelDate1, PayDate1, DelDate2, PayDate2, DealFund, PayFund, Owner, Portfolio \}$
 $Op^{InSecurities} = \{ OperDate, OperTime, OperNo, Security, Num, Owner, Portfolio, DelDate, Depo1, DepoAcc1, DepoRazd1 \}$
 $Op^{OutSecurities} = \{ OperDate, OperTime, OperNo, Security, Num, Owner, Portfolio, DelDate, Depo1, DepoAcc1, DepoRazd1, \}$
 $\quad \quad \quad \{ Depo2, DepoAcc2, DepoRazd2 \}$
 $Op^{MoveSecurities} = \{ OperDate, OperTime, OperNo, Direction, Security, Num, Owner, Portfolio, DelDate, \}$
 $\quad \quad \quad \{ Depo1, DepoAcc1, DepoRazd1, Depo2, DepoAcc2, DepoRazd2 \}$
 $Op^{MargDeal} = \{ DealType, DealDate, TradeTime, TradeNo, Security, DealNum, DealPrice, DealQty, DealCoupon, \}$
 $\quad \quad \quad \{ OrderNo, TradePlace, Owner, Portfolio, MargValue1, MargValue2 \}$
 $Op^{InMoney} = \{ OperDate, OperTime, OperNo, Fund, Qty, Owner, Portfolio, PayDate, BankRecviz1 \}$, где $BankRecviz1 \in \{ BankRecviz \}$
 $Op^{OutMoney} = \{ OperDate, OperTime, OperNo, Fund, Qty, Owner, Portfolio, PayDate, BankRecviz1, BankRecviz2 \}$, где
 $\quad \quad \quad \{ BankRecviz1, BankRecviz2 \} \subset \{ BankRecviz \}$
 $Op^{MoveMoney} = \{ OperDate, OperTime, OperNo, Direction, Fund, Qty, Owner, Portfolio, PayDate, BankRecviz1, BankRecviz2 \}$, где
 $\quad \quad \quad \{ BankRecviz1, BankRecviz2 \} \subset \{ BankRecviz \}$
 $CI_{ListUslug}^{Fiz} = \{ U_i \}_{n_1}$, $U_i \in Uslugi$, $n_1 \geq 1$, $CI_{ListUslug}^{Jur} = \{ U_i \}_{n_2}$, $U_i \in Uslugi$, $n_2 \geq 1$, $CI_{ListUslug} = \{ CI_{ListUslug}^{Fiz}, CI_{ListUslug}^{Jur} \}$, $|CI_{ListUslug}| \geq 1$
 $CI_{Type}^{Fiz} = \{ CI_{Type_1}^{Fiz}, CI_{Type_2}^{Fiz}, \dots, CI_{Type_m}^{Fiz} \}$, $m \geq 0$, $CI_{Type}^{Jur} = \{ CI_{Type_1}^{Jur}, CI_{Type_2}^{Jur}, \dots, CI_{Type_k}^{Jur} \}$, $k \geq 0$, $CI_{Type} = CI_{Type}^{Fiz} \cup CI_{Type}^{Jur}$, $|CI_{Type}| \geq 1$
 $CI_{Type} = \gamma \left(\left\langle CI_{PersInfo}, CI_{ListUslug} \right\rangle \right)$, $\gamma = \{ \gamma_1, \gamma_2, \dots, \gamma_l \}$, $l > 0$
 $n_1 - ? n_2 - ? m - ? k - ? CI_{Type}^{Fiz} - ? CI_{Type}^{Jur} - ? CI_{Type} - ? \gamma_1 \cdot \gamma_2 \cdot \dots \cdot \gamma_l - ?$, $\phi \left(CI_{ListUslug} \mid CI_{Type}, \text{Доходность}, \text{Экон. риски} \right) = \text{optimum}$, $\gamma \rightarrow \phi$



КЛАССИФИКАЦИЯ МЕТОДОВ ИНТЕЛЛЕКТУАЛЬНОГО АНАЛИЗА ДАННЫХ



КЛАССИФИКАЦИЯ МЕТОДОВ КЛАСТЕРИЗАЦИИ



КРИТЕРИИ ВЫБОРА МЕТОДОВ КЛАСТЕРИЗАЦИИ

Критерий 1 (Cr_1): Объем информации (количество строчек или по занимаемому месту) по отношению к времени обработки.

Критерий 2 (Cr_2): Размерность информации (количество атрибутов в строке) в порядковом выражении.

Критерий 3 (Cr_3): Типы атрибутов: числовой дискретный, числовой непрерывный, строковый.

Критерий 4 (Cr_4): Чувствительность к равномерности информации (наличие аномалий-выбросов во входном наборе данных).

Критерий 5 (Cr_5): Априорное (экспертное) представление о форме получаемых кластеров.

Критерий 6 (Cr_6): Априорное (экспертное) представление о количестве кластеров.

Критерий 7 (Cr_7): Априорное (экспертное) представление о перекрываемости кластеров.

Критерий 8 (Cr_8): Качество кластеризации.

Критерий Метод	Критерий 1 (Cr_1)	Критерий 2 (Cr_2)	Критерий 3 (Cr_3)	Критерий 4 (Cr_4)	Критерий 5 (Cr_5)	Критерий 6 (Cr_6)	Критерий 7 (Cr_7)	Критерий 8 (Cr_8)
CURE	линейная зависимость, большой объем данных	невысокая размерность	числовые	нечувствителен	сложная форма	входной параметр	без пересечений	высокое
BIRCH	линейная зависимость, большой объем данных	невысокая размерность	числовые	чувствителен	сферическая форма	входной параметр	без пересечений	среднее
MST	линейно-логарифмическая зависимость, большой объем данных	средняя размерность	числовые	чувствителен	сложная форма, в т.ч. выпуклой и вогнутой форм	входной параметр	без пересечений	высокое
k-средних	нелинейная зависимость, небольшой объем данных	низкая размерность	числовые	очень чувствителен	сферическая форма	входной параметр	без пересечений	низкое
PAM	нелинейная зависимость, небольшой объем данных	низкая размерность	числовые	чувствителен	сложная форма	входной параметр	без пересечений	среднее
CLOPE	линейная зависимость, огромный объем данных	средняя размерность	числовые	чувствителен	сложная форма	вычисляемая величина	без пересечений	высокое
Самоорганизующиеся карты Кохонена	нелинейная зависимость, большой объем данных	высокая размерность	числовые	низкая чувствительность	сложная форма	входной параметр	без пересечений	среднее
HCM	нелинейная зависимость, большой объем данных	средняя размерность	числовые	чувствителен	сложная форма	входной параметр	без пересечений	среднее
Fuzzy C-means	нелинейная зависимость, большой объем данных	высокая размерность	числовые	чувствителен	сложная форма	входной параметр	с пересечениями	среднее

АДАПТАЦИЯ МЕТОДОВ КЛАСТЕРИЗАЦИИ

Par_{Dist}^{Birch} - способ определения дистанции между кластерами

Par_{Qual}^{Birch} - метод оценки качества кластеризации

$Par_{ThrQual}^{Birch}$ - порог для метода оценки качества кластеризации

Par_{Thr1}^{Birch} - начальный порог (Фаза 1)

$Par_{OutlinePerc1}^{Birch}$ - процент аномалий (выбросов) в общем объеме (Фаза 1)

$Par_{ClustAlg3}^{Birch}$ - алгоритм выполнения дополнительной кластеризации (Фаза 3)

$Par_{RefPass4}^{Birch}$ - количество очисток кластеров (Фаза 4)

$Par_{OutlinePerc4}^{Birch}$ - процент аномалий (выбросов) в общем объеме (Фаза 4)

Par_k^{Birch} - количество кластеров

Par_c^{HCM} - количество кластеров

$Par_{StopThr}^{HCM}$ - пороговая величина для остановки алгоритма

Par_k^{PAM} - количество кластеров

$Par_{MethodCh}^{PAM}$ - способ выбора начальных центров

$Par_{MaxIterNum}^{PAM}$ - максимальное количество итераций

$Par_c^{Fuzzy\ C-means}$ - количество кластеров

$Par_m^{Fuzzy\ C-means}$ - экспонентциальный вес

$Par_{StopThr}^{Fuzzy\ C-means}$ - пороговая величина для остановки алгоритма

$Par_k^{Cure} \equiv Par_k^{Birch} \equiv Par_k^{MST} \equiv Par_k^{k-средних} \equiv Par_k^{PAM} \equiv Par_k^{Kohonen} \equiv Par_c^{HCM} \equiv Par_c^{Fuzzy\ C-means} \equiv Par_c$

$Par_{VolPerc1}^{Kohonen} + Par_{VolPerc2}^{Kohonen} + Par_{VolPerc3}^{Kohonen} = 100\%$; $Par_{MethodCh}^{k-средних} \equiv Par_{MethodCh}^{PAM} \equiv Par_{MethodCh}$; $Par_{MaxIterNum}^{k-средних} \equiv Par_{MaxIterNum}^{PAM} \equiv Par_{MaxIterNum}$

$Par_k^{k-средних}$ - количество кластеров

$Par_{MethodCh}^{k-средних}$ - способ выбора начальных центров

$Par_{MaxIterNum}^{k-средних}$ - максимальное количество итераций

Par_k^{MST} - количество кластеров

$Par_k^{Kohonen}$ - количество кластеров

$Par_{VolPerc1}^{Kohonen}$ - объем обучающего множества

$Par_{VolPerc2}^{Kohonen}$ - объем валидационного множества

$Par_{VolPerc3}^{Kohonen}$ - объем тестового множества

$Par_{DivMethod}^{Kohonen}$ - способ разделения множества

$Par_{StopThr}^{Kohonen}$ - порог остановки алгоритма обучения и тестирования нейронов

$Par_{TrainRate}^{Kohonen}$ - скорость обучения нейронов

Par_k^{Clope} - количество кластеров

Par_k^{Cure} - количество кластеров

Par_c^{Cure} - количество точек, которое захватывается при обработке за один раз

Par_p^{Cure} - количество предварительных разделов

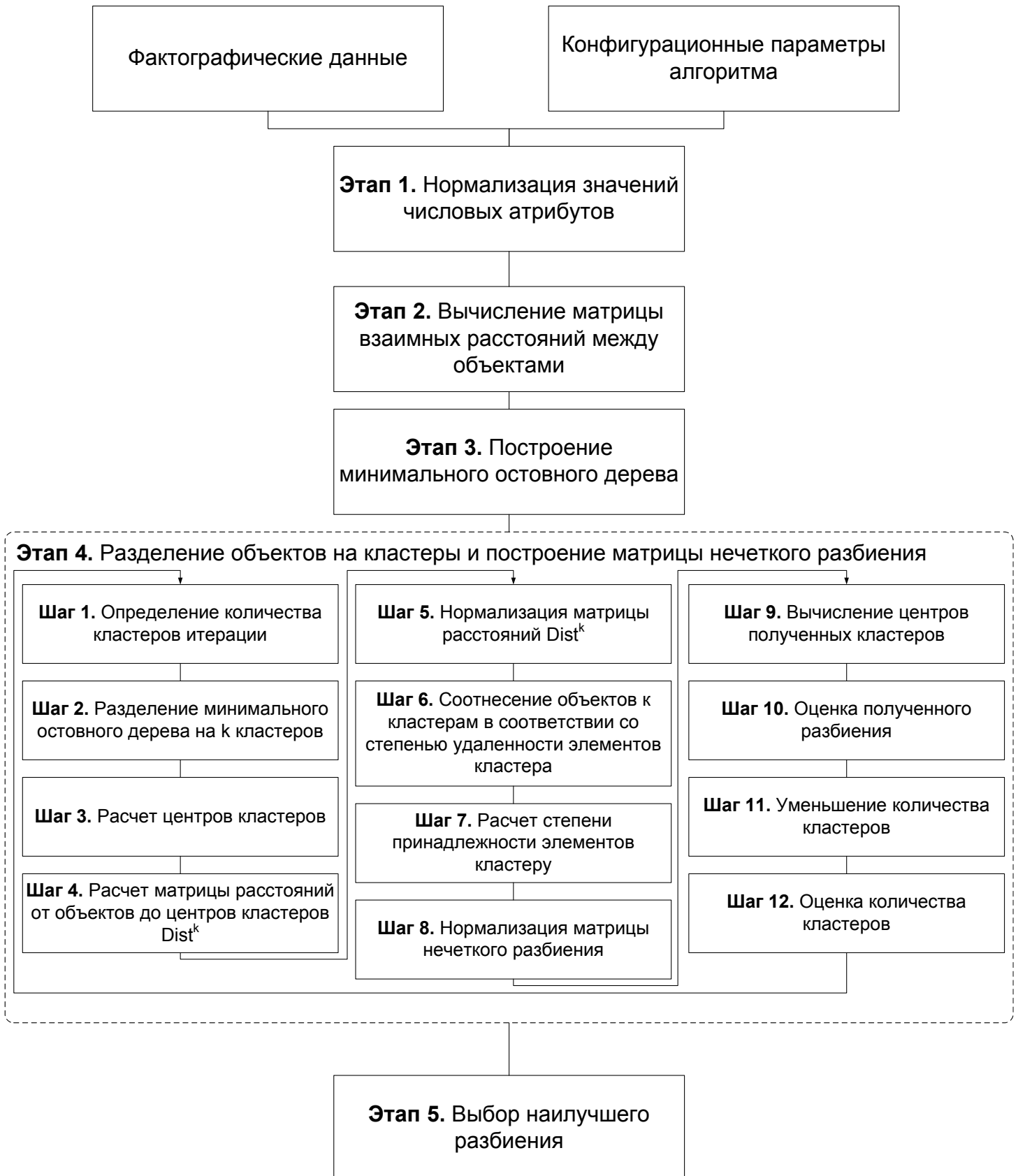
Par_q^{Cure} - сокращающий коэффициент для предварительных разделов

Par_a^{Cure} - коэффициент удаленности

Метод \ Параметр	Параметр 1	Параметр 2	Параметр 3	Параметр 4	Параметр 5	Параметр 6	Параметр 7	Параметр 8	Параметр 9
CURE	Par_k^{Cure}	Par_c^{Cure}	Par_q^{Cure}	Par_a^{Cure}	Par_p^{Cure}				
BIRCH	Par_k^{Birch}	Par_{Dist}^{Birch}	Par_{Qual}^{Birch}	$Par_{ThrQual}^{Birch}$	Par_{Thr1}^{Birch}	$Par_{OutlinePerc1}^{Birch}$	$Par_{ClustAlg3}^{Birch}$	$Par_{RefPass4}^{Birch}$	$Par_{OutlinePerc4}^{Birch}$
MST	Par_k^{MST}								
k-средних	$Par_k^{k-средних}$	$Par_{MethodCh}^{k-средних}$	$Par_{MaxIterNum}^{k-средних}$						
PAM	Par_k^{PAM}	$Par_{MethodCh}^{PAM}$	$Par_{MaxIterNum}^{PAM}$						
CLOPE	Par_r^{Clope}								
Само-организующиеся карты Кохонена	$Par_k^{Kohonen}$	$Par_{VolPerc1}^{Kohonen}$	$Par_{VolPerc2}^{Kohonen}$	$Par_{VolPerc3}^{Kohonen}$	$Par_{DivMethod}^{Kohonen}$	$Par_{StopThr}^{Kohonen}$	$Par_{TrainRate}^{Kohonen}$		
HCM	Par_c^{HCM}	$Par_{StopThr}^{HCM}$							
Fuzzy C-means	$Par_c^{Fuzzy\ C-means}$	$Par_m^{Fuzzy\ C-means}$	$Par_{StopThr}^{Fuzzy\ C-means}$						

АДАПТИВНЫЙ АЛГОРИТМ КЛАСТЕРИЗАЦИИ ФАКТОГРАФИЧЕСКИХ ДАННЫХ

Укрупненная схема



АДАПТИВНЫЙ АЛГОРИТМ КЛАСТЕРИЗАЦИИ ФАКТОГРАФИЧЕСКИХ ДАННЫХ (Этапы 1 - 4)

Входные данные алгоритма:

$D = \{u_1, u_2, \dots, u_m\}$, где u_i – объекты кластеризации, m – количество объектов кластеризации, $i = \overline{1, m}$;
 $u_i = \{(Value_{i1}, t_1), (Value_{i2}, t_2), \dots, (Value_{in}, t_n)\}$, где $Value_{ij}$ – значение j^{oo} атрибута i^{oo} объекта кластеризации, t_j – тип атрибута объекта кластеризации, n – количество атрибутов объекта кластеризации, $j = \overline{1, n}$;
 $t_j = \{ValueType_j, FieldType_j\}$, где $ValueType_j$ – тип значения атрибута, $ValueType_j \in ValueTypes$, $FieldType_j$ – вид значения атрибута, $FieldType_j \in FieldTypes$; $FieldTypes = \{Входное, Идентифицирующее, Информационное\}$;
 $ValueTypes = \{\text{Целочисленный тип}, \text{Денежный тип}, \text{Лингвистический тип}\}$, где $\text{Целочисленный тип} \subset \mathbb{Z}$,
 $\text{Денежный тип} \subset \mathbb{R}$, $\text{Лингвистический тип} \subset \text{Словарная система}$;
 $\text{Словарная система} = \{\text{Лингв.тип1}, \text{Лингв.тип2}, \dots, \text{Лингв.типS}\}$, где Лингв.тип_i – объект словарной системы, характеризующий оценочные/качественные показатели объектов кластеризации;
 q – максимальное количество кластеров, $q \leq m$; p – размазанность кластеров, $p \in (0; 10]$;
 $K = \{K_1, K_2, \dots, K_n\}$, где K_i – весовой коэффициент влияния атрибута объекта, $K_i \in [0; 1]$;
 w – степень удаленности элементов, $w \in (0; 1]$;
 $Metric$ – способ определения расстояния между объектами, $Metric \in Metrics$;
 $Metrics = \left\{ \begin{array}{l} \text{Евклидово расстояние}, \text{Квадрат Евклидова} \\ \text{расстояния}, \text{расстояние Чебышева} \end{array} \right\}$;
 $OstTreeMethod$ – способ построения минимального остовного дерева, $OstTreeMethod \in OstTreeMethods$;
 $OstTreeMethods = \left\{ \begin{array}{l} \text{Алгоритм Борувки}, \text{Алгоритм} \\ \text{Крускала}, \text{Алгоритм Прима} \end{array} \right\}$;
 $NormMethod$ – способ проведения нормализации значений числовых атрибутов, $NormMethod \in NormMethods$;
 $NormMethods = \{\text{Линейная нормализация}, \text{Статистическая нормализация}\}$

Выходные данные алгоритма: $C = \{C_1, C_2, \dots, C_c \mid O^c \rightarrow \max, c \leq q, C_1 \cup C_2 \cup \dots \cup C_c = D\}$, $u_i \in C_j, i = \overline{1, m}, j = \overline{1, c}$, $\mu_{ij} \in [0; 1]$

Описание алгоритма:

Этап 1. Нормализация значений числовых атрибутов.

В случае линейной нормализации выполняется следующее:

$$Value_{ij} := \left\{ \begin{array}{l} \frac{Value_{ij}}{\text{Max}(Value_{ij})} \mid \text{Max}(Value_{ij}) \neq 0, \\ \text{Max}(Value_{ij}) \mid t_j \in \{\text{Целочисленный тип}, \text{Денежный тип}\} \end{array} \right\}$$

В случае статистической нормализации выполняется

$$\text{следующее: } Value_{ij} := \left\{ \begin{array}{l} \frac{Value_{ij} - \frac{\sum_{i=1}^m Value_{ij}}{m}}{\sqrt{\frac{\sum_{i=1}^m (Value_{ij})^2}{m} - \left(\frac{\sum_{i=1}^m Value_{ij}}{m}\right)^2}} \mid \text{Max}(Value_{ij}) \neq 0, \\ \frac{\sum_{i=1}^m (Value_{ij})^2}{m} - \left(\frac{\sum_{i=1}^m Value_{ij}}{m}\right)^2 \mid t_j \in \{\text{Целочисленный тип}, \text{Денежный тип}\} \end{array} \right\}$$

Этап 2. Вычисление матрицы взаимных расстояний между объектами.

$Dist_{ij} = \|u_i - u_j\| = Metric(u_i, u_j)$, где $Metric$ – способ определения расстояния между объектами.

Если $Metric = \text{Евклидово расстояние}$, то $Dist_{ij} = \sqrt{\sum_w ([Value_{iw} - Value_{jw}] * K_w)^2}$,

Если $Metric = \text{Квадрат Евклидова расстояния}$, то $Dist_{ij} = \sum_w ([Value_{iw} - Value_{jw}] * K_w)^2$,

Если $Metric = \text{Расстояние Чебышева}$, то $Dist_{ij} = \text{Max}_w [Value_{iw} - Value_{jw}] * K_w$,

где $i, j \in [1, m]$, $w = \overline{1, n}$ при условии $FieldType[w] = \text{"Входное"}$

Этап 3. Построение минимального остовного дерева по выбранному способу построения дерева с использованием матрицы взаимных расстояний между объектами $Dist$.

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

Матрица нечеткого разбиения: $F = [\mu_{ij}]$, $\mu_{ij} \in [0, 1]$, $i \leq q$, $j = \overline{1, m}$, где μ_{ij} – степень принадлежности i^{oo} объекта к

j^{ny} кластеру. Матрица разбиения обладает следующими свойствами: $\sum_{i=1}^k \mu_{ij} = 1, j = \overline{1, m}, 0 < \sum_{j=1}^m \mu_{ij} \leq m, i = \overline{1, k}$.

АДАПТИВНЫЙ АЛГОРИТМ КЛАСТЕРИЗАЦИИ ФАКТОГРАФИЧЕСКИХ ДАННЫХ (Этапы 4 - 5)

Шаг 1. Определение количества кластеров итерации: $k := q$.

Шаг 2. $Dist_{ij}^k := \left\{ 0 \mid Dist_{ij}^k = \text{Max} \right\}$.

Шаг 3. Расчет центров выделенных кластеров V_i^k .

$V_i^k = \text{Avg} \left(\left\{ u_j \mid u_j \in C_i^k \right\} \right)$, где Avg – оператор вычисления среднего значения показателей объектов, входящих в кластер k , $i = \overline{1, k}$, $j = \overline{1, m}$.

Для числовых типов оператор Avg определяется выражением:

$$\text{Avg}[r] = \frac{\sum_{u_j \in V_i^k} \left\{ Value_{jr} \mid FieldType[w] = \text{"Входное"} \right\}}{|V_i^k|}, \quad j = \overline{1, m}, \quad r = \overline{1, n}.$$

Для лингвистических типов оператор Avg определяется выражением:

$$\text{Avg}[r] = \left\{ \begin{array}{l} Value_{jr} \\ \sum_{\substack{u_j \in V_i^k, u_l \in V_i^k \\ \|Value_{jr} - Value_{lr}\| = \min}} \varphi = \max_{u_j \in V_i^k, u_l \in V_i^k} \end{array} \mid FieldType[w] = \text{"Входное"} \right\}, \quad r = \overline{1, n}, \quad j = \overline{1, m}, \quad l = \overline{1, m}, \quad \text{где } \varphi - \text{частота значения}$$

атрибута в пределах кластера V_i^k .

Шаг 4. Расчет матрицы расстояний от объектов до центров кластеров V_i^k : $Dist_{ij}^k = \|V_i^k - u_j\| = \text{Metric}(V_i^k, u_j)$, $i = \overline{1, k}$, $j = \overline{1, m}$, где Metric – способ определения расстояния между объектами.

Шаг 5. Нормализация матрицы расстояний от объектов до центров кластеров V_i^k :

$$Dist_{ij}^{k'} = \begin{cases} \frac{Dist_{ij}^k}{\text{Max}(Dist_{ij}^k)}, & \text{Max}(Dist_{ij}^k) \neq 0 \\ 1, & \text{Max}(Dist_{ij}^k) = 0 \end{cases}, \quad i = \overline{1, k}, \quad j = \overline{1, m}.$$

Шаг 6. Соотнесение объектов к кластерам в соответствии со степенью удаленности элементов кластера (w):

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

Шаг 7. Расчет степени принадлежности кластеру: $\mu_{ij} = (1 - Dist_{ij}^{k'})^2$, $i = \overline{1, k}$, $j = \overline{1, m}$.

Шаг 8. Нормализация матрицы нечеткого разбиения: $\mu_{ij} = \frac{\mu_{ij}}{\sum_{i=1}^m \mu_{ij}}$, $j = \overline{1, m}$.

Шаг 9. Вычисление центров полученных кластеров с использованием матрицы нечеткого разбиения:

$$V_i^{k'} = \frac{\sum_{j=1}^m \mu_{ij}^p * u_j}{\sum_{j=1}^m \mu_{ij}^p}, \quad i = \overline{1, k}. \quad \text{Для лингвистических атрибутов центра кластера вычисление производится с}$$

использованием выражения: $V_i^{k'}[r] = Value_{jr} \mid_{\mu_{ij} = \text{Max}(\mu_{ij})}$.

$$\sum_{i=1, k} \frac{|V_i^{k'}| * \sum_{j=1}^m \mu_{ij}^p * \|V_i^{k'} - u_j\|}{\text{Min}_{i \neq j} (\|V_i^{k'} - u_j\|) * \text{Max}_{u_j \in V_i^{k'}} (\|V_i^{k'} - u_j\|) * \sum_{j=1}^m \|V_i^{k'} - u_j\| * k}$$

Шаг 10. Оценка качества разбиения: $O^k = \frac{\sum_{i=1, k} \frac{|V_i^{k'}| * \sum_{j=1}^m \mu_{ij}^p * \|V_i^{k'} - u_j\|}{\text{Min}_{i \neq j} (\|V_i^{k'} - u_j\|) * \text{Max}_{u_j \in V_i^{k'}} (\|V_i^{k'} - u_j\|) * \sum_{j=1}^m \|V_i^{k'} - u_j\| * k}}{m * k^2}$, где

$|V_i^{k'}|$ – количество элементов в кластере i ;

$\|V_i^{k'} - u_j\| = \text{Metric}(V_i^{k'}, u_j)$ – расстояние от центра кластера i до элемента u_j ;

$u_j \in V_i^{k'}$ – отражение условия о принадлежности элемента кластеру.

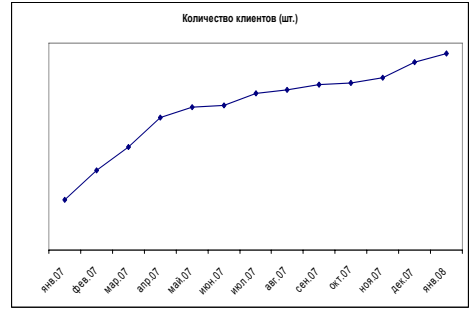
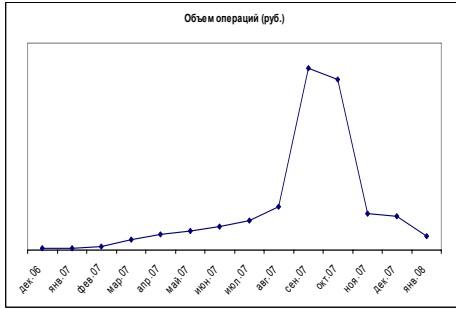
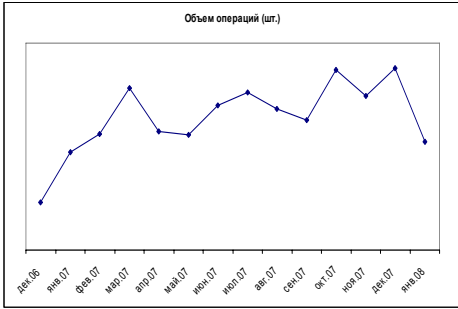
Шаг 11. $k := k - 1$.

Шаг 12. Если $k > 0$, то переход на шаг 2.

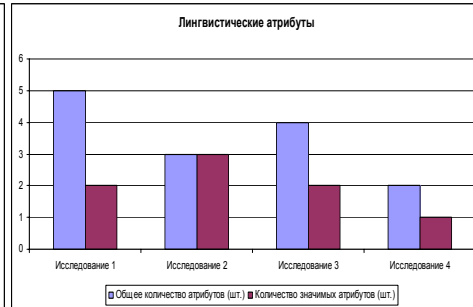
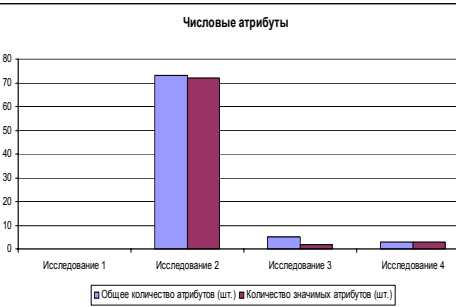
Этап 5. Выбор наилучшего разбиения: $O_{Onm} = \text{MAX}_{i=1, q} (O^i)$.

АНАЛИЗ ЭКСПЕРИМЕНТАЛЬНЫХ ДАННЫХ

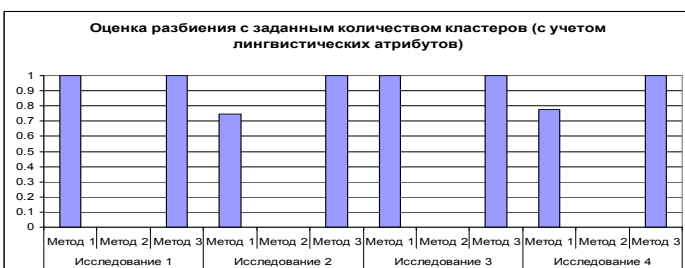
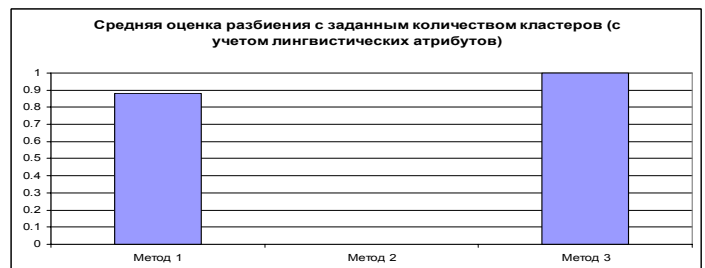
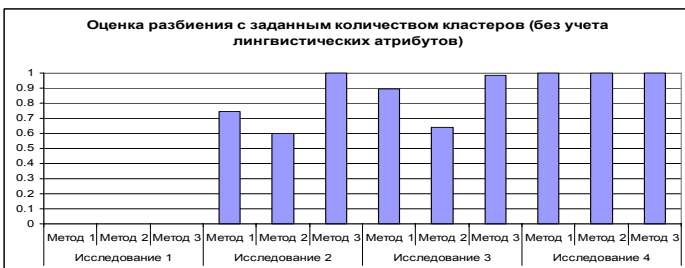
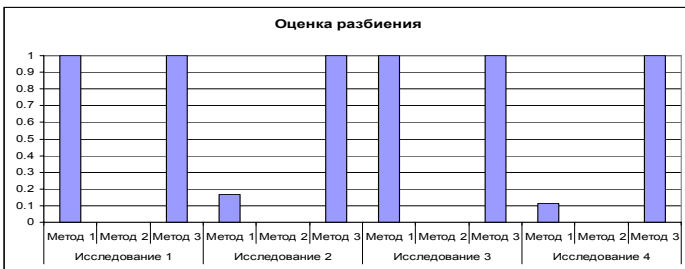
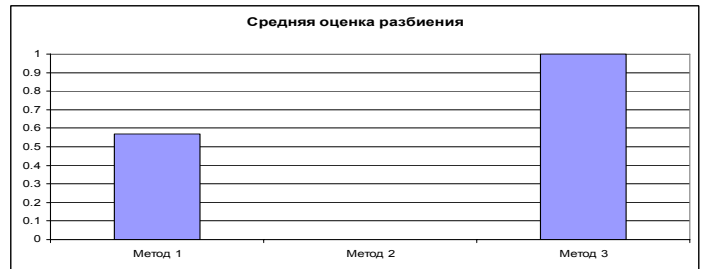
Показатели операций брокерского обслуживания



Характеристики исследуемых данных



ОЦЕНКА ЭФФЕКТИВНОСТИ МЕТОДИКИ



СРАВНЕНИЕ ЭФФЕКТИВНОСТИ МЕТОДИК

- Исследование 1: выделение секторов инвестирования на основе анализа показателей финансовых инструментов;
- Исследование 2: выделение групп клиентов на основе данных об оборотах за период, частоты проведения операций, количестве финансовых инструментов, группы используемых финансовых инструментов;
- Исследование 3: выявление категорий финансовых инструментов для оценки эффективности операций;
- Исследование 4: выделение классов автомобилей на основе данных о максимальной скорости, цвете кузова, воздушном сопротивлении, массе.

Метод 1 – карта Кохонена, Метод 2 – алгоритм k – средних, Метод 3 – ADAKL

Общая оценка разбиений

Показатель Метод		Количество полученных кластеров	Оценка разбиения	Оценка разбиения с заданным количеством кластеров (без учета лингвистических атрибутов)	Оценка разбиения с заданным количеством кластеров (с учетом лингвистических атрибутов)
1	Метод 1	30	1.0000	-	1.0000
	Метод 2	-	-	-	-
	Метод 3	30	1.0000	-	1.0000
2	Метод 1	12	0.1667	0.7467	0.7467
	Метод 2	-	-	0.6000	-
	Метод 3	3	1.0000	1.0000	1.0000
3	Метод 1	15	1.0000	0.8953	1.0000
	Метод 2	-	-	0.6388	-
	Метод 3	10	1.0000	0.9857	1.0000
4	Метод 1	9	0.1111	1.0000	0.7778
	Метод 2	-	-	1.0000	-
	Метод 3	3	1.0000	1.0000	1.0000

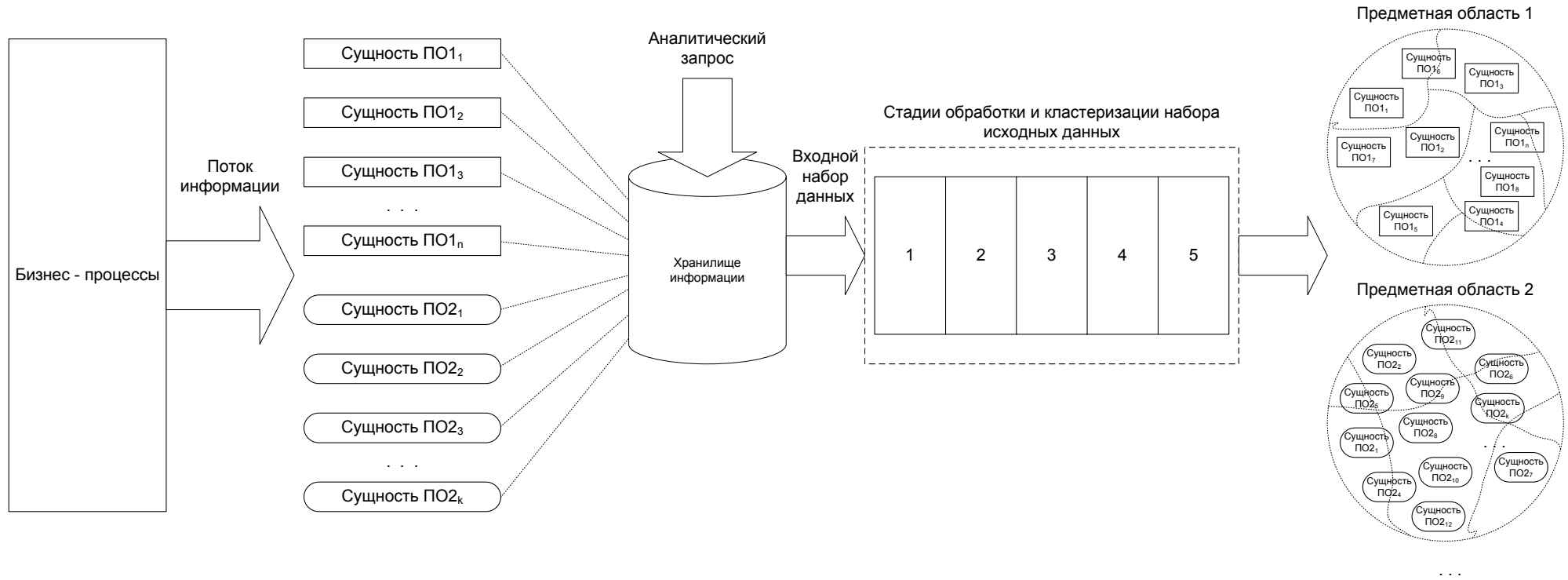
Средняя оценка разбиений

Оценка Метод	Оценка разбиения	Оценка разбиения с заданным количеством кластеров (без учета лингвистических атрибутов)	Оценка разбиения с заданным количеством кластеров (с учетом лингвистических атрибутов)	Итоговая оценка
Метод 1	0.5695	0.8807	0.8811	0.7771
Метод 2	-	0.7463	-	0.2488
Метод 3	1.0000	0.9952	1.0000	0.9984

Оценка разбиения выполняется на основе показателей выполненной кластеризации:
$$O = \frac{r}{n} \times \begin{cases} q/k, & q \leq k \\ k/q, & q > k \end{cases}$$

- где
- q – количество кластеров по итогам кластеризации;
 - r – количество элементов, правильно распределенных по соответствующим кластерам;
 - k – исходное количество кластеров;
 - n – количество объектов кластеризации.

Укрупненная схема обработки потока информации



Укрупненная схема архитектуры программного комплекса

