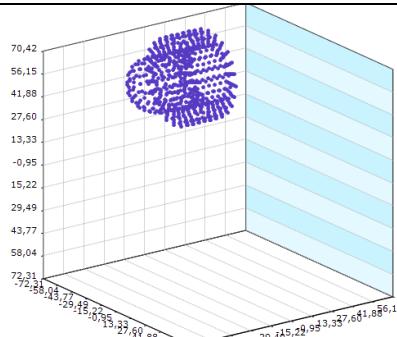
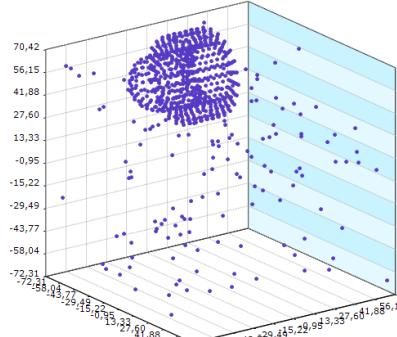


Додаток А

Таблиця А.1 – Експериментальні та реальні дані.

Кількість точок	% шуму	Характеристики вибірки	Просторове зображення
722	0	AveExpectation= 2,10 maxExpectation= 7,08 minExpectation= -0,77 AveStdDeviation= 1600092,39 maxStdDeviation= 2387575,17 minStdDeviation= 321328,11 AveRange= 59,49 maxRange= 72,31 minRange= 52,44 AveCorrelation= 3,03 maxCorrelation= 9,02 minCorrelation= -1,58	
	20	AveExpectation= 0,19 maxExpectation= 0,58 minExpectation= 0,00 AveStdDeviation= 3325480,53 maxStdDeviation= 4479679,53 minStdDeviation= 1069160,12 AveRange= 140,25 maxRange= 141,31 minRange= 138,42 AveCorrelation= 0,67 maxCorrelation= 1,98 minCorrelation= -0,35	

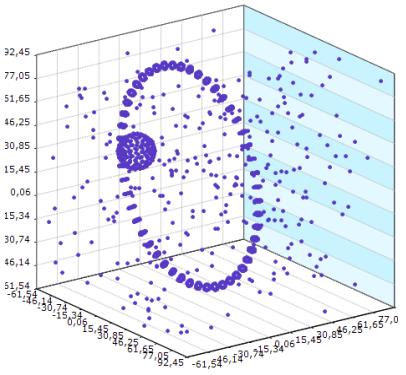
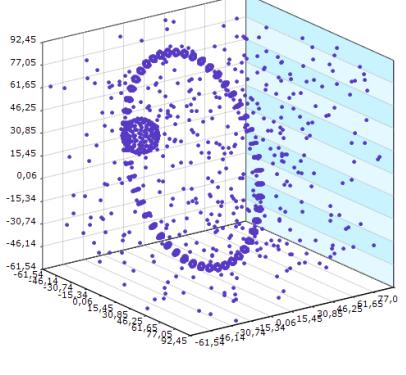
Продовження таблиці А.1.

	40	<p>AveExpectation= -0,05 maxExpectation= 0,53 minExpectation= -0,68 AveStdDeviation= 4625156,99 maxStdDeviation= 6082670,76 minStdDeviation= 1798336,22 AveRange= 141,58 maxRange= 142,42 minRange= 141,00 AveCorrelation= 0,31 maxCorrelation= 0,91 minCorrelation= -0,15</p>	
	60	<p>AveExpectation= -0,05 maxExpectation= 0,49 minExpectation= -0,65 AveStdDeviation= 5941195,32 maxStdDeviation= 7593426,39 minStdDeviation= 2672303,41 AveRange= 141,25 maxRange= 142,42 minRange= 140,31 AveCorrelation= 0,23 maxCorrelation= 0,67 minCorrelation= -0,09</p>	

Продовження таблиці А.1.

782	0	<p>AveExpectation= -0,39 maxExpectation= 0,00 minExpectation= -0,67 AveStdDeviation= 4241167,92 maxStdDeviation= 6547451,32 minStdDeviation= 82153,88 AveRange= 103,39 maxRange= 144,09 minRange= 22,27 AveCorrelation= 0,01 maxCorrelation= 0,02 minCorrelation= 0,00</p>	
	20	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 6901696,96 maxStdDeviation= 10357365,41 minStdDeviation= 771019,95 AveRange= 152,00 maxRange= 152,54 minRange= 151,00 AveCorrelation= -0,02 maxCorrelation= 0,01 minCorrelation= -0,06</p>	

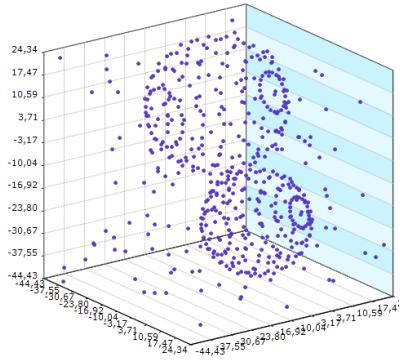
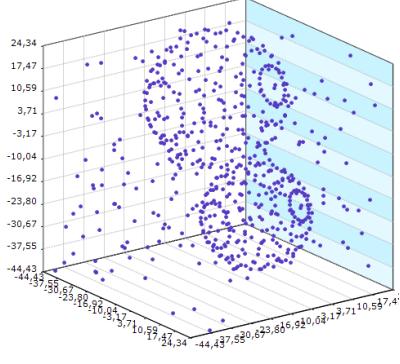
Продовження таблиці А.1.

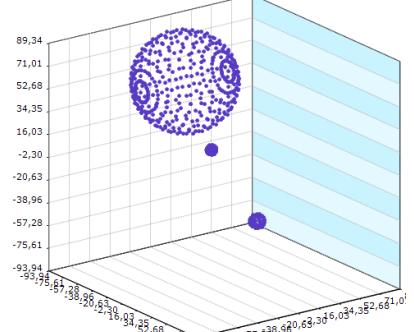
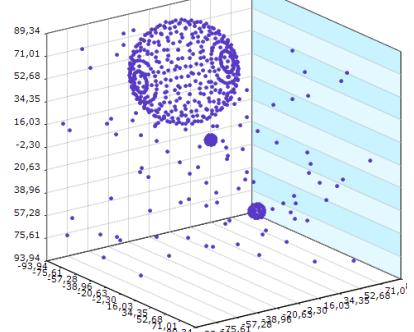
	40	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 9418044,41 maxStdDeviation= 13890736,90 minStdDeviation= 1521326,92 AveRange= 152,33 maxRange= 153,45 minRange= 151,00 AveCorrelation= -0,03 maxCorrelation= 0,00 minCorrelation= -0,06</p>	
	60	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 11845187,01 maxStdDeviation= 17273110,25 minStdDeviation= 2607450,88 AveRange= 152,33 maxRange= 153,45 minRange= 151,54 AveCorrelation= -0,05 maxCorrelation= 0,00 minCorrelation= -0,10</p>	

Продовження таблиці А.1.

382	0	<p>AveExpectation= 0,10 maxExpectation= 0,70 minExpectation= -0,31 AveStdDeviation= 766331,19 maxStdDeviation= 1405571,80 minStdDeviation= 77981,37 AveRange= 52,79 maxRange= 65,17 minRange= 39,69 AveCorrelation= -0,18 maxCorrelation= 0,10 minCorrelation= -0,69</p>	
	20	<p>AveExpectation= -0,61 maxExpectation= 0,00 minExpectation= -1,50 AveStdDeviation= 1144620,84 maxStdDeviation= 1940495,72 minStdDeviation= 225749,65 AveRange= 66,26 maxRange= 68,34 minRange= 63,43 AveCorrelation= -0,26 maxCorrelation= 0,00 minCorrelation= -0,62</p>	

Продовження таблиці А.1.

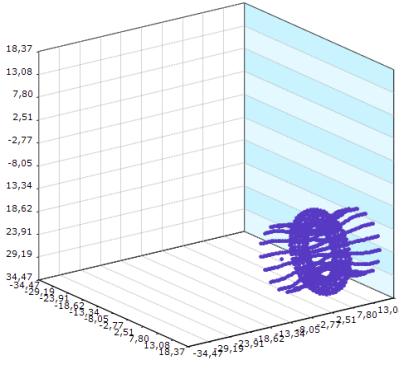
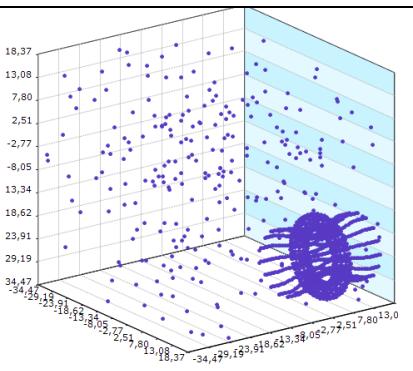
	40	<p>AveExpectation= -0,10 maxExpectation= 0,00 minExpectation= -0,24 AveStdDeviation= 1496430,74 maxStdDeviation= 2447026,02 minStdDeviation= 336807,99 AveRange= 67,59 maxRange= 68,34 minRange= 67,00 AveCorrelation= -0,25 maxCorrelation= 0,00 minCorrelation= -0,55</p>	
	60	<p>AveExpectation= -0,10 maxExpectation= 0,00 minExpectation= -0,24 AveStdDeviation= 1821133,68 maxStdDeviation= 2894346,77 minStdDeviation= 442973,54 AveRange= 67,59 maxRange= 68,34 minRange= 67,00 AveCorrelation= -0,26 maxCorrelation= 0,00 minCorrelation= -0,48</p>	

556	0	<p>AveExpectation= -10,72 maxExpectation= 0,24 minExpectation= -32,11 AveStdDeviation= 2464777,25 maxStdDeviation= 4209462,00 minStdDeviation= 500369,72 AveRange= 114,75 maxRange= 134,30 minRange= 79,37 AveCorrelation= 0,51 maxCorrelation= 1,54 minCorrelation= 0,00</p>	
	20	<p>AveExpectation= -0,14 maxExpectation= 0,17 minExpectation= -0,38 AveStdDeviation= 3973166,20 maxStdDeviation= 6541042,98 minStdDeviation= 948258,37 AveRange= 178,09 maxRange= 180,94 minRange= 175,00 AveCorrelation= 0,25 maxCorrelation= 0,76 minCorrelation= -0,06</p>	

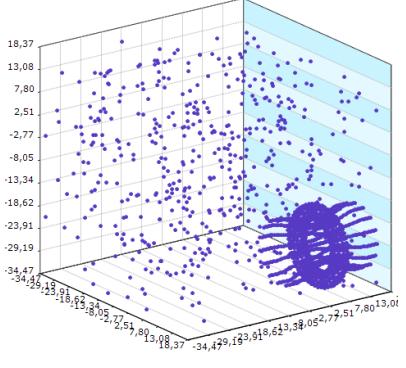
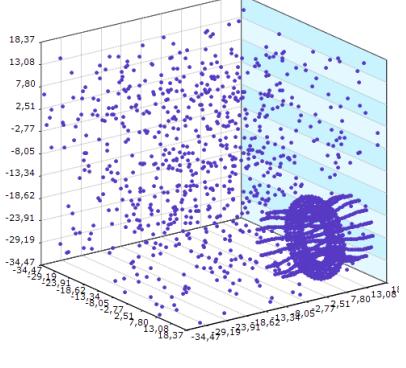
Продовження таблиці А.1.

	40	<p>AveExpectation= -0,02 maxExpectation= 0,09 minExpectation= -0,17 AveStdDeviation= 5522130,71 maxStdDeviation= 8696946,19 minStdDeviation= 1748233,48 AveRange= 180,09 maxRange= 180,34 minRange= 179,94 AveCorrelation= 0,17 maxCorrelation= 0,50 minCorrelation= -0,01</p>	
	60	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 7203610,91 maxStdDeviation= 10850564,28 minStdDeviation= 2779506,90 AveRange= 181,42 maxRange= 182,34 minRange= 180,94 AveCorrelation= 0,09 maxCorrelation= 0,27 minCorrelation= -0,02</p>	

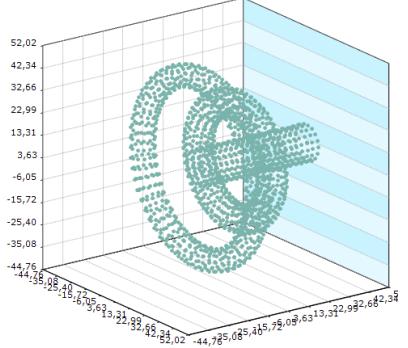
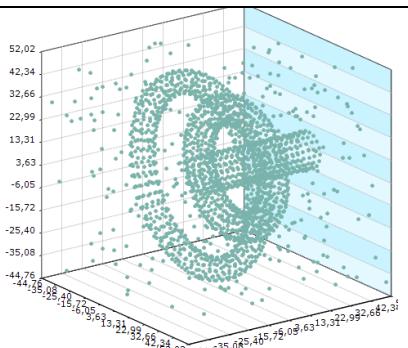
Продовження таблиці А.1.

1292	0	<p>AveExpectation= 0,91 maxExpectation= 2,81 minExpectation= -0,10 AveStdDeviation= 4212128,29 maxStdDeviation= 10029279,67 minStdDeviation= 349007,62 AveRange= 24,43 maxRange= 34,47 minRange= 18,81 AveCorrelation= 4,49 maxCorrelation= 9,42 minCorrelation= -2,38</p>	
	20	<p>AveExpectation= -0,06 maxExpectation= 0,00 minExpectation= -0,17 AveStdDeviation= 5999719,90 maxStdDeviation= 13782016,65 minStdDeviation= 917930,93 AveRange= 51,61 maxRange= 52,37 minRange= 51,00 AveCorrelation= 0,69 maxCorrelation= 1,32 minCorrelation= -0,24</p>	

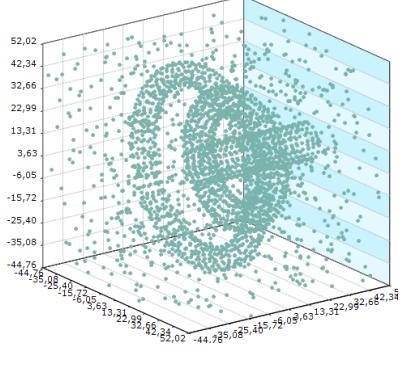
Продовження таблиці А.1.

	40	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00</p> <p>AveStdDeviation= 7457980,14 maxStdDeviation= 16997770,61 minStdDeviation= 1210568,25</p> <p>AveRange= 51,61 maxRange= 52,37 minRange= 51,00</p> <p>AveCorrelation= 0,26 maxCorrelation= 0,51 minCorrelation= -0,06</p>	
	60	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00</p> <p>AveStdDeviation= 8915388,80 maxStdDeviation= 20191761,83 minStdDeviation= 1526109,33</p> <p>AveRange= 51,61 maxRange= 52,37 minRange= 51,00</p> <p>AveCorrelation= 0,10 maxCorrelation= 0,21 minCorrelation= -0,01</p>	

Продовження таблиці А.1.

1751	0	<p>AveExpectation= -0,03 maxExpectation= 0,00 minExpectation= -0,04 AveStdDeviation= 8698787,61 maxStdDeviation= 13130396,47 minStdDeviation= 1260447,72 AveRange= 76,50 maxRange= 86,39 minRange= 56,84 AveCorrelation= -0,02 maxCorrelation= 0,04 minCorrelation= -0,10</p>	
	20	<p>AveExpectation= 0,05 maxExpectation= 0,16 minExpectation= 0,00 AveStdDeviation= 13304925,36 maxStdDeviation= 19083532,59 minStdDeviation= 3392606,23 AveRange= 95,59 maxRange= 96,02 minRange= 95,00 AveCorrelation= -0,03 maxCorrelation= 0,00 minCorrelation= -0,08</p>	

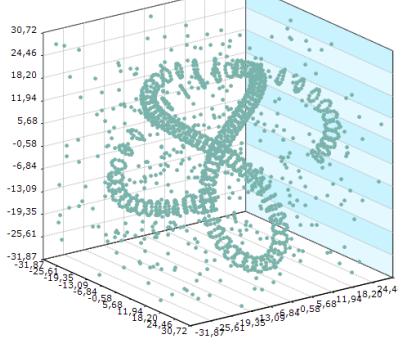
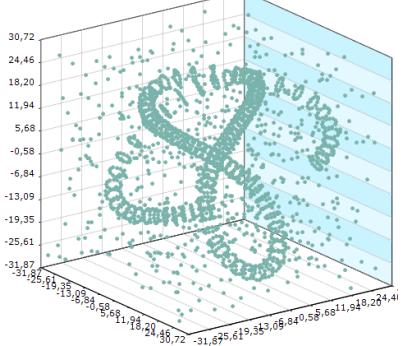
Продовження таблиці А.1.

	40	<p>AveExpectation= -0,05 maxExpectation= 0,00 minExpectation= -0,16 AveStdDeviation= 16912056,08 maxStdDeviation= 24188415,21 minStdDeviation= 4333312,74 AveRange= 95,59 maxRange= 96,02 minRange= 95,00 AveCorrelation= -0,02 maxCorrelation= 0,00 minCorrelation= -0,06</p>	
	60	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 20458703,13 maxStdDeviation= 29141035,22 minStdDeviation= 5281287,93 AveRange= 95,59 maxRange= 96,02 minRange= 95,00 AveCorrelation= -0,02 maxCorrelation= 0,00 minCorrelation= -0,05</p>	

Продовження таблиці А.1.

1294	0	<p>AveExpectation= 0,00 maxExpectation= 0,01 minExpectation= 0,00 AveStdDeviation= 14292162,05 maxStdDeviation= 17868392,54 minStdDeviation= 9762846,59 AveRange= 52,45 maxRange= 60,01 minRange= 41,37 AveCorrelation= 0,00 maxCorrelation= 0,01 minCorrelation= 0,00</p>	
	20	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00 AveStdDeviation= 18999204,54 maxStdDeviation= 23385494,06 minStdDeviation= 13156549,88 AveRange= 60,86 maxRange= 61,72 minRange= 60,00 AveCorrelation= 0,00 maxCorrelation= 0,00 minCorrelation= 0,00</p>	

Кінець таблиці А.1.

	40	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00</p> <p>AveStdDeviation= 23367447,26 maxStdDeviation= 28523867,45 minStdDeviation= 16244155,10</p> <p>AveRange= 60,86 maxRange= 61,72 minRange= 60,00</p> <p>AveCorrelation= 0,00 maxCorrelation= 0,00 minCorrelation= 0,00</p>	
	60	<p>AveExpectation= 0,00 maxExpectation= 0,00 minExpectation= 0,00</p> <p>AveStdDeviation= 27727745,75 maxStdDeviation= 33648341,17 minStdDeviation= 19328719,57</p> <p>AveRange= 60,86 maxRange= 61,72 minRange= 60,00</p> <p>AveCorrelation= 0,00 maxCorrelation= 0,00 minCorrelation= 0,00</p>	

Додаток Б

Слайди презентації

Актуальність дослідів

- ▲ Необхідність синхронізувати і організовувати на єдиній критеріальній основі вибір метода класифікації на основі даних аналізуємої вибірки
- ▲ Потребність уніфіцирувати алгоритми класифікації і за рахунок цього зменшити час вибору метода
- ▲ Необхідність забезпечення користувачів якісним рішенням завдання аналізу при різних дослідюваних даних
- ▲ Постійне зростання обсягу надходженої інформації і різноманітність цієї інформації вимагає розвитку технологій аналізу цих даних
- ▲ Необхідність створення універсального метода класифікації. окремі методи добре працюють на певних даних, але не є універсальними
- ▲ Необхідність аналізу складних наборів даних, з пересекаючими та накладаючими класами

Цель роботи

Розробка універсальної математичної моделі залежності вибору комбінації алгоритмів на різних етапах ієрархічного алгоритму Хамелеон від початкових характеристик аналізуємого набору даних з метою підвищення якості класифікації.

2

Постановка завдання дослідів

Представимо постановку завдання як розробку такої математичної моделі вибору алгоритмів A_r для різних етапів алгоритма класифікації Хамелеон для розв'язання R завдання вибору оптимальної комбінації алгоритмів $A = \sum A_r$, що складається з множини $A_r, r = 1, R$. Вибір будь-якої комбінації сводиться до вибору певного підмножества, розв'язання завдання, яке використовує оптимальну критерію якості та статистичні характеристики об'єкта. Представимо модель наступним чином:

$$\{D, H_D\} \xrightarrow{A_r \in A} \{E^r\}$$

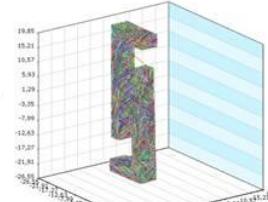
де A – множина алгоритмів; D – дослідувана выборка даних H_D – множина характеристик дослідуваної выборки даних; E – множина показників якості розв'язаної завдання за вибірку алгоритмів для різних етапів алгоритму класифікації Хамелеон, оптимального для даної выборки.

3

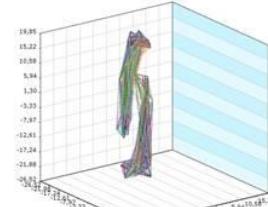
Модификация этапа огрубления графа

В процессе стадии огрубления строится последовательность меньших графов, каждый с меньшим количеством узлов.

- Случайное паросочетание (RM).
- Паросочетание из тяжелых ребер (HEM).
- Модифицированный алгоритм паросочетания тяжелых ребер (Modified Heavy Edge Matching - HEM*).
- Паросочетание из наиболее тяжелых ребер (heaviest-edge matching).
- Модифицированное паросочетание из наиболее тяжелых ребер HEM*+.
- Паросочетание легких ребер (LEM).
- Паросочетание из тяжелых клик (HCM).
- Сочетание тяжелых треугольников. Heavy-triangle matching (HTM).
- Сочетание тяжелых схем (Heaviest Schema Matching (HSM)).
- Сочетание гиперребер (Hyperedge Coarsening HEC).
- Видоизмененое сочетание гиперребер (Modified Hyperedge Coarsening MHEC).
- Сочетание лучшего (первого) выбора (First Choice Coarsening FCC).



6



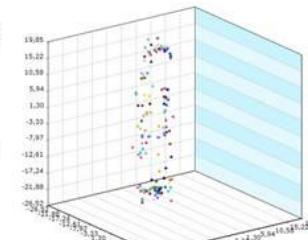
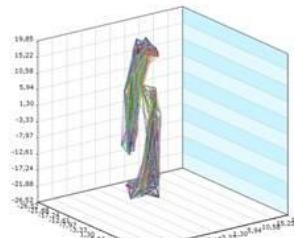
6

Модификация этапа начального разделения графа в рамках алгоритма Хамелеон (Initial partitioning)

Для решения задачи разбиения графа можно рекурсивно применить метод бинарного деления, при котором на первой итерации граф разделяется на две равные части, далее на втором шаге каждая из полученных частей также разбивается на две части и т.д.

В случае, когда требуемое количество разбиений n не является степенью двойки, каждое деление пополам необходимо осуществлять в соответствующем соотношении.

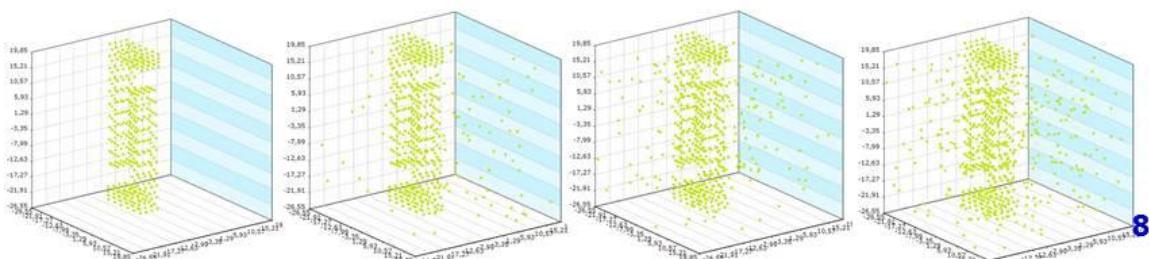
- Покоординатное разбиение (Coordinate NestedDissection (CND)).
- Деление сети с использованием кривых, заполняющих пространство (Space-filling Curv Techniques (SFCT)).
- Рекурсивное разделение графа (Recursive Graph Bisection (RGB)).
- Алгоритм возрастающего графа (Graph Growth Partitioning (GGP)).
- Алгоритм возрастающего графа с учетом выгод (Greedy Graph Growing Algorithm (GGGP)).
- Уровневое ячеичное разбиение(Levelized Nested Dissection (LND)).
- Seed-Growth bisection (SGB).



7

Описание реальных и экспериментальных выборок

- ➡ Реальные выборки данных были получены с профильных ресурсов. Выбрано 46 выборок для анализа.
- ➡ С ресурсов так же получены экспериментальные выборки используемые для проверки кластеризации. Выбрано 35 выборок с различными характеристиками и количеством объектов в выборке.
- ➡ Предложен метод генерации трехмерных экспериментальных выборок на основе экспорттированных данных построенной модели в системе 3Ds Max. Сгенерировано 33 выборки. Каждая из выборок анализировалась в 4 вариантах:
 - ➡ без добавления шума;
 - ➡ с добавлением 20% шума;
 - ➡ с добавлением 40% шума;
 - ➡ с добавлением 60% шума.



8

Математическая модель зависимости выбора параметра k при построении k -пп графа от исходных характеристик выборки

В рамках работы было проведено 3 эксперимента для выбора управляемых параметров:

- ➡ В первом эксперименте анализировались такие характеристики как количество объектов в выборке, минимальные и максимальные значения матожидания, дисперсии и разброса. Зависимости между данными параметрами и значением k не выявлено;
- ➡ Во втором эксперименте в качестве управляемого параметра были выбраны длина наибольшего остовного ребра полностью связного графа и среднее значение длины всех остальных ребер остова. Данные характеристики показывают зависимость, но использование данного подхода не является целесообразным в связи с трудоемкостью построения остова полностью связного графа;
- ➡ В третьем эксперименте в качестве характеристики использовались количество компонент связности и максимальное расстояние между компонентами связности и количество элементов в компоненте связности. Вторая характеристика вычисляется следующим образом:

$$SetDist = \max \left(\frac{dist(avComponent_i, avComponent_j)}{\max \left(\frac{ComponentOstovEdge_{ij}}{ComponentVertexNum_{ij}} \right)} \right)$$

где $avComponent$ -центроид компоненты связности, $ComponentOstovEdge$ -ребро соединяющее вершины принадлежащие одной компоненте, $ComponentVertexNum$ -количество вершин в компоненте. Данные характеристики нетрудоемки в расчете и существует зависимость между ними и значением k .

9

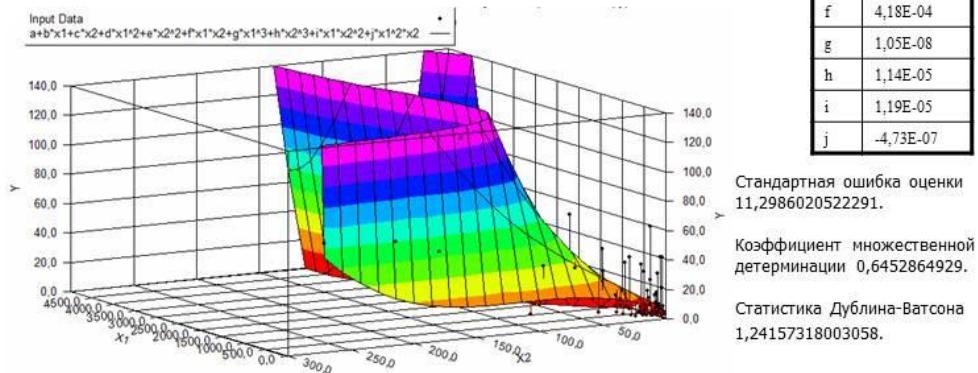
Математическая модель зависимости выбора параметра k для асимметричного графа

В результате исследования была получена следующая мат. модель:

$$k = a + b \cdot x_1 + c \cdot x_2 + d \cdot x_1^2 + e \cdot x_2^2 + f \cdot x_1 \cdot x_2 + g \cdot x_1^3 + h \cdot x_2^3 + i \cdot x_1 \cdot x_2^2 + j \cdot x_1^2 \cdot x_2$$

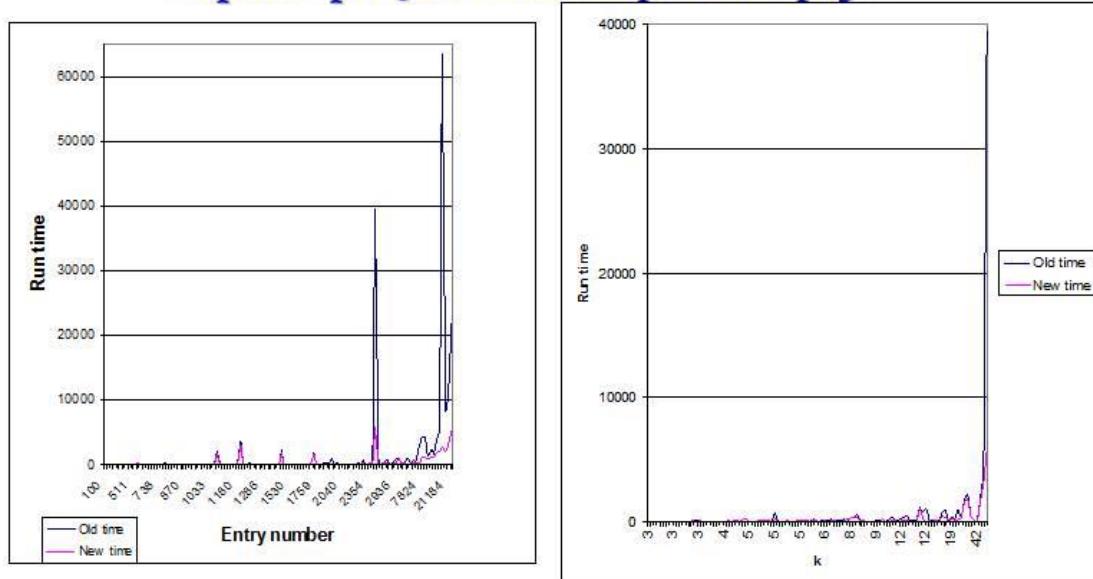
где x_1 – коэффициент расстояния, x_2 – количество компонент связности.

a	4,963024
b	2,33E-02
c	0,42939
d	-4,45E-05
e	-3,86E-03
f	4,18E-04
g	1,05E-08
h	1,14E-05
i	1,19E-05
j	-4,73E-07



10

Применение математической модели зависимости выбора параметра k для асимметричного графа



Применение данной модели улучшили время выполнения этапа построения графа в 62.45% случаев. В 37.55% случаев время выполнения ухудшилось. Отрицательный результат применения модели получен в 7.71% случаев. В среднем время выполнения улучшилось на 161%.

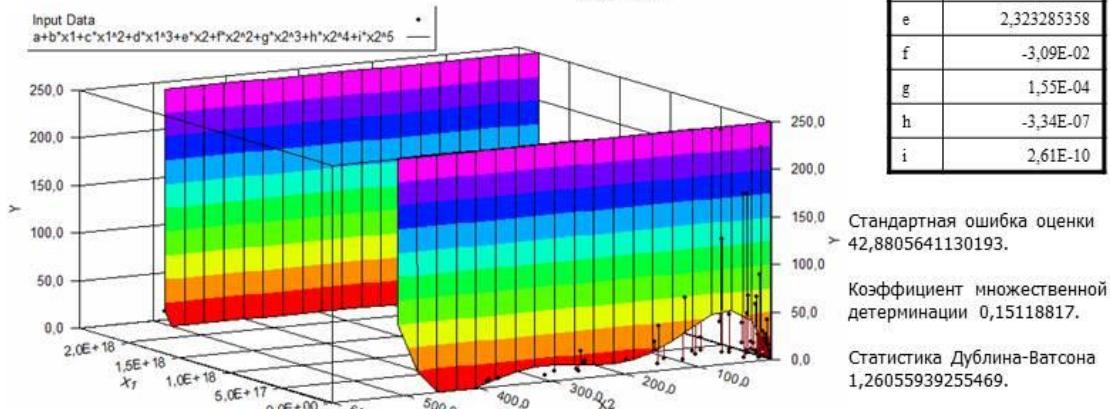
11

Математическая модель зависимости выбора параметра k для симметричного графа

В результате исследования была получена следующая мат. модель:

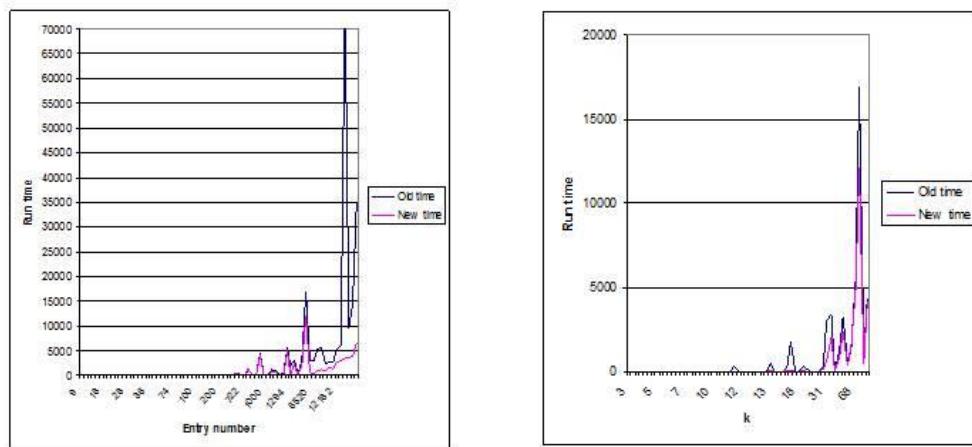
$$k = a + b \cdot x_1 + c \cdot x_1^2 + d \cdot x_1^3 + e \cdot x_2 + f \cdot x_2^2 + g \cdot x_2^3 + h \cdot x_2^4 + i \cdot x_2^5$$

где x_1 - коэффициент расстояния, x_2 - количество компонент связности.



12

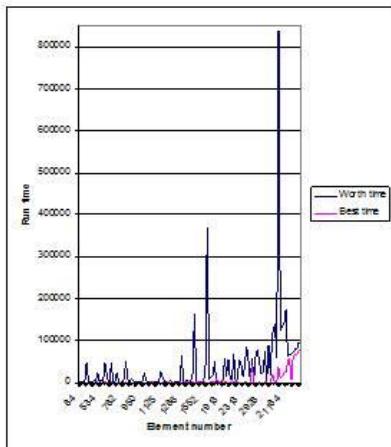
Применение математической модели зависимости выбора параметра k для симметричного графа



Применение данной модели улучшили время выполнения этапа построения графа в 69.23% случаев. В 20.51% случаев время выполнения ухудшилось. Отрицательный результат применения модели получен в 5.12% случаев. В среднем время выполнения улучшилось на 169%

13

Выводы по экспериментальным результатам работы модифицированного алгоритма Хамелеон



• Для сложных выборок иногда есть смысл использовать алгоритм KL, для остальных в большинстве случаев целесообразно использовать более быстрый алгоритм FM и более простые граничные алгоритмы KL и FM.

• Для выборок со сложной структурой данных качество разделения наиболее точно оценивается с помощью индекса Данна.

• Применение алгоритма построения симметричного графа нецелесообразно из-за низкой скорости работы.

• Качество работы алгоритмов на этапе восстановления и улучшения графа падает с увеличением количества классов.

14

ВЫВОДЫ

- Впервые построена модель зависимости качества кластеризации сложных линейноразделимых защумленных данных различного объема от характеристик данных и алгоритмов динамической кластеризации.
- Дальнейшее развитие получил метод Хамелеон, который в отличие от существующего использует разные алгоритмы на разных этапах кластеризации, что позволяет учитывать разницу в данных и использовать методы, которые работают лучше на конкретных данных.
- Усовершенствован критерий сравнения алгоритмов и анализа исходных данных применяемых в рамках алгоритма Хамелеон
- Усовершенствован метод построения графа, который позволяет ускорить процесс построения графа посредством выбора оптимального k , основываясь на характеристиках анализируемых данных. Созданы модели выбора k для алгоритма k -ближайших соседей используемого в алгоритме Хамелеон, основанная на характеристиках данных.

16