

## **Кластерный анализ данных о рейтинге лучших вузов России с помощью программного пакета визуального программирования Orange**

*Голубева Евгения Павловна*

*Приамурский государственный университет имени Шолом-Алейхема  
Студент*

### **Аннотация**

Цель данной статьи – выполнить кластерный анализ данных о рейтинге лучших вузов России за 2021 г. Для кластерного анализа был использован программный пакет визуального программирования на основе компонентов для визуализации данных Orange и данные о рейтинге лучших вузов России. С помощью средств визуализации Orange выполнили кластерный анализ данных о рейтинге лучших вузов России и получили итоговую схему.

**Ключевые слова:** Orange, виджет, визуализация данных, кластерный анализ.

## **Cluster analysis of data on the ranking of the best universities in Russia using the Orange visual programming software package**

*Golubeva Evgeniya Pavlovna*

*Sholom-Aleichem Priamursky State University  
Student*

### **Abstract**

The purpose of this article is to perform a cluster analysis of data on the ranking of the best universities in Russia for 2021. For cluster analysis, a visual programming software package based on Orange data visualization components and data on the ranking of the best universities in Russia were used. With the help of Orange visualization tools, we performed a cluster analysis of data on the ranking of the best universities in Russia and obtained the final scheme.

**Keywords:** Orange, widget, data visualization, cluster analysis.

## **1 Введение**

### **1.1 Актуальность**

Современный мир стал невероятно зависим от данных, которые создаются и хранятся в большом количестве в различных сферах деятельности. При этом, для их обработки и анализа требуются соответствующие методы и инструменты. Одним из наиболее популярных методов анализа данных является кластерный анализ.

Преимущества кластерного анализа включают возможность автоматического выделения групп объектов на основе их характеристик, что позволяет сократить время и усилия, затрачиваемые на анализ данных

вручную. Кроме того, кластерный анализ может помочь выявить скрытые связи между объектами, которые могут быть невидимы на первый взгляд.

Программа Orange благодаря своим функциям позволяет провести кластерный анализ данных.

### **1.2 Обзор исследований**

А. В. Леонов в статье рассматривал основные алгоритмы кластеризации категориальных данных применительно к различным типам пользовательских интерфейсов, определяются их достоинства и недостатки. [1]. В данной статье проанализировали основные аспекты кластерного анализа больших объемов данных при помощи различных методов, их сравнения и выделения наиболее эффективного А. В. Клименко, И. С. Слащев [2]. Д. В. Гринченков, Ф. Х. Нгуен, Т. Т. Нгуен, Д. А. Горбушин выполнили краткий обзор и сравнительный анализ возможностей алгоритмов, используемых для интеллектуального анализа данных [3]. В статье рассматривали исследование программного обеспечения Data Mining Ю. С. Кривенко, А.Т. Минасян и А.О. Разиньков [4]. С.В. Пальмов и А.А. Жуйкова в статье описали функционал аналитического пакета Orange, предназначенный для поиска часто встречающихся наборов элементов и ассоциативных правил [5]. В статье рассмотрел использование методов кластеризации в программе Orange на основе реальной базы данных. Н. Юсупов [6]. С. С. Мастевой и А.Н. Петрова статью посвятили краткому обзору методов Data Mining в период информационной эры [7].

### **1.3 Цель исследования**

Цель исследования - выполнить кластерный анализ данных с помощью программы Orange.

## **2 Материалы и методы**

Для кластерного анализа используется программа Orange. Работа будет происходить на готовых данных рейтинг вузов.xlsx, скачать которые можно по ссылке:

[https://docs.google.com/spreadsheets/d/1GsA6zDjMFrzWnXgcXubjvJ0jdu\\_CT\\_MC/edit?usp=sharing&ouid=104272149632818699735&rtpof=true&sd=true](https://docs.google.com/spreadsheets/d/1GsA6zDjMFrzWnXgcXubjvJ0jdu_CT_MC/edit?usp=sharing&ouid=104272149632818699735&rtpof=true&sd=true)

## **3 Результаты и обсуждения**

Перед началом работы требуется установить Orange с официального сайта и установить.

- 1) Создадим новый файл (рис.1).

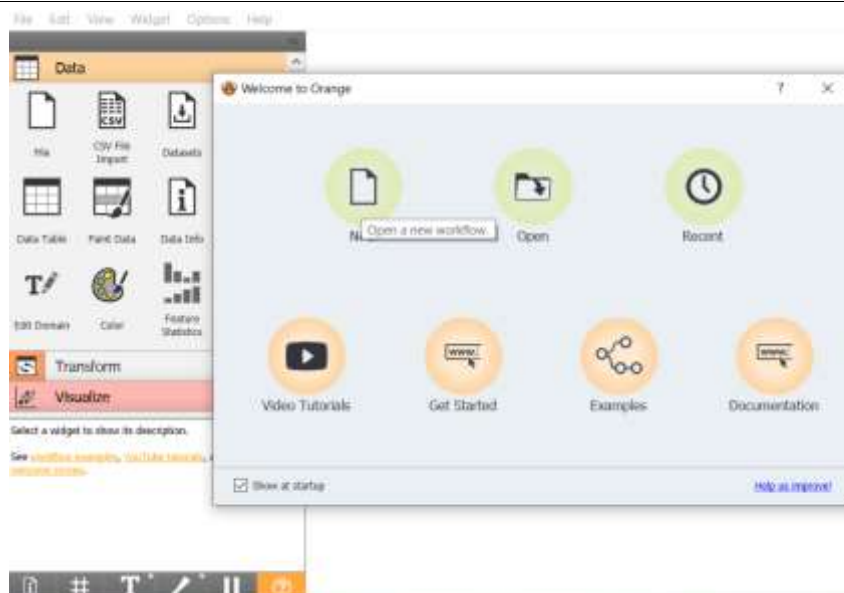


Рис.1. Создание нового файла

2) Добавляем виджет File на холст (Рис.2).

Виджет File считывает файл входных данных (таблицу данных с экземплярами данных) и отправляет набор данных в выходной канал. История последних открытых файлов хранится в виджете. Виджет также включает каталог с образцами наборов данных, которые поставляются с предустановленным с Orange.

Виджет считывает данные из Excel (.xlsx), простых файлов с разделителями табуляции (.txt), файлов, разделенных запятыми (.csv) или URL-адресов.



Рис.2. Добавление виджета File на холст

Есть 3 способа добавления виджета на холст:

1. Дважды щелкните на виджет.
2. Перетащите виджет на холст.
3. Щелкните правой кнопкой мыши на холсте для меню виджета.

3) Чтобы добавить файл необходимо открыть виджет file на холсте (Рис.3).

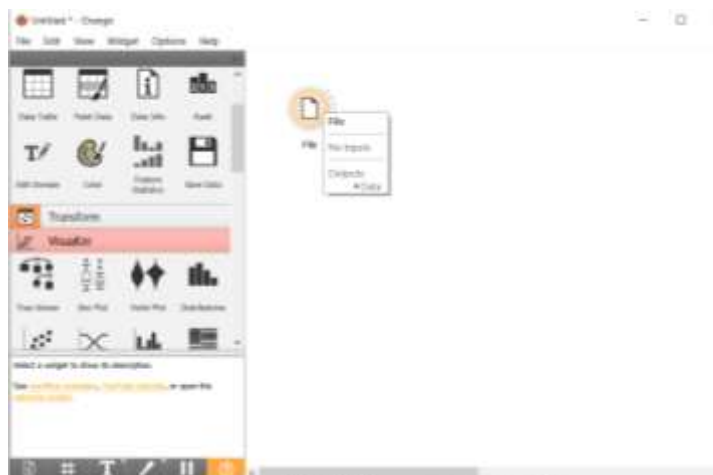


Рис.3. Открытие виджета File

4) Открылось диалоговое окно File (Рис.4).

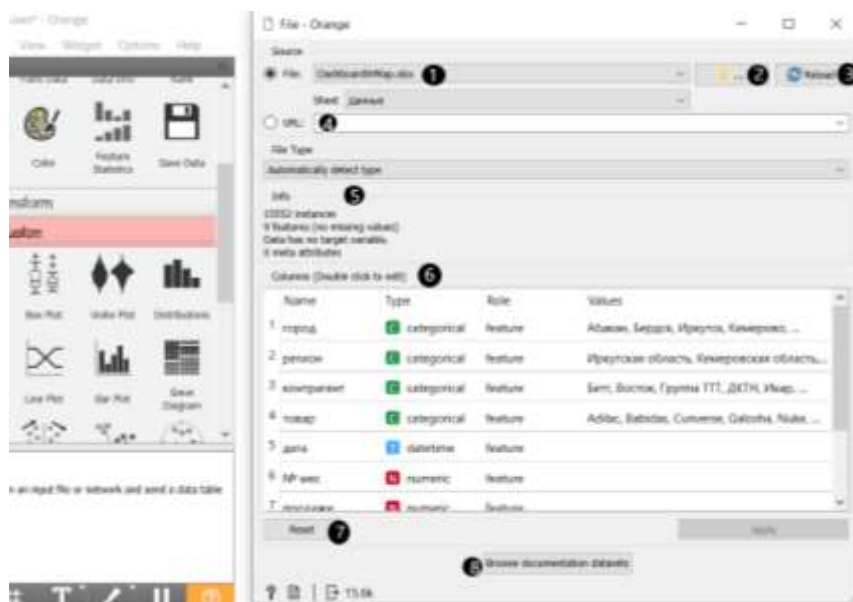


Рис.4. Диалоговое окно File

Описание диалогового окна File (Рис.4).

1. Просмотр ранее открытых файлов данных или загрузка любой из образцов.
2. Найти файл данных.
3. Перегрузка выбранного в данный момент файл данных.
4. Вставка данных из URL-адресов, включая данные из Google Таблиц.
5. Информация о загруженном наборе данных: размер набора данных, количество и типы объектов данных.
6. Дополнительные сведения о функциях в наборе данных. Объекты можно редактировать, дважды щелкнув по ним. Пользователь может изменить имена атрибутов, выбрать тип переменной для каждого атрибута

(Continuous, Nominal, String, Datetime) и выбрать способ дальнейшего определения атрибутов (как Features, Targets или Meta). Пользователь также может проигнорировать атрибут.

7. Сброс.
8. Просмотр наборов данных документации.

5) Добавляем файл рейтинг вузов.xlsx (Рис.5).

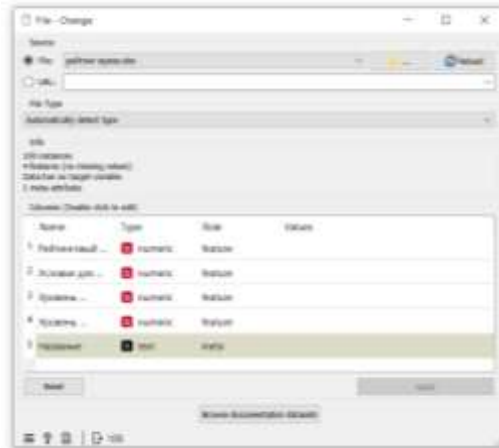


Рис.5. Добавление файла

6) Добавляем виджет Data Table на холст (Рис.6).

Data Table - получает один или несколько наборов данных в виде входных данных и представляет их в виде электронной таблицы. Экземпляры данных могут быть отсортированы по значениям атрибутов. Виджет также поддерживает ручной выбор экземпляров данных.



Рис.6. Добавление виджета Data Table

7) Чтобы посмотреть данные таблицы необходимо соединить два виджета на холсте File и Data Table (Рис.7).



Рис.7. Соединение виджетов

8) Открываем виджет Data Table что бы просмотреть данные загруженной таблицы (Рис.8).

Instance	Идентификатор	Масса (кг)	Средняя температура	Средняя влажность	Средняя скорость ветра
1	Минеральный	4.66451	5	11	7
2	Минеральный	4.61078	2	9	2
3	Минеральный	4.5179	3	3	6
4	Сланцевый	4.48281	3	11	9
5	Минеральный	3.47629	6	4	12
6	Минеральный	3.40297	8	2	13
7	Минеральный	3.25449	4	7	34
8	Минеральный	4.23225	7	14	11
9	Сланцевый	4.14146	12	12	8
10	Песчаный	3.95152	14	8	25
11	Минеральный	3.94883	11	28	7
12	Песчаный	3.93036	18	18	14
13	Минеральный	3.87487	15	3	25
14	Минеральный	3.86883	9	40	5
15	Песчаный	3.85269	13	8	36
16	Минеральный	3.8268	21	13	10
17	Минеральный	3.79472	18	52	3
18	Песчаный	3.72441	17	14	15
19	Песчаный	3.67385	16	24	17
20	Сланцевый	3.65391	18	28	4
21	Минеральный	3.64463	23	15	21
22	Песчаный	3.63075	22	21	28
23	Сланцевый	3.59118	18	18	19
24	Минеральный	3.388	25	17	20
25	Доломитовый	3.37847	27	22	22
26	Песчаный	3.37286	21	16	48
27	Песчаный	3.31832	25	40	38
28	Сланцевый	3.31281	27	20	16
29	Минеральный	2.99621	24	32	100
30	Минеральный	2.98642	19	52	89
31	Минеральный	2.98312	25	70	18
32	Минеральный	2.98182	36	25	52
33	Сланцевый	2.96282	26	81	23
34	Минеральный	2.91627	41	18	24
35	Минеральный	2.88884	24	21	65
36	Сланцевый	2.87513	33	34	63
37	Минеральный	2.85847	38	46	54

Рис.8. Диалоговое окно Data Table

Описание диалогового окна Data Table (Рис.8).

1. Сведения о текущем размере набора данных, количестве и типах атрибутов
2. Значения непрерывных атрибутов могут быть визуализированы с помощью баров; знача может быть отнесен к различным классам.
3. Экземпляры данных (строки) могут быть выбраны и отправлены на выход виджета канал.
4. Используется кнопка Restore Original Order для изменения порядка экземпляров данных после сортировки на основе атрибутов.

9) Для того чтобы провести кластерный анализ необходимо использовать виджет K-Means. Кластеризация K-средних (K-Means) – простой метод

разделения множества данных на  $K$  различных непересекающихся кластеров. Для выполнения кластеризации сначала нужно определить желаемое число кластеров  $K$ , затем алгоритм  $K$  средних будет относить каждое наблюдение в точности к одному из  $K$  кластеров.

10) Из раздела *unsupervised* выбираем виджет *k-Means*, добавляем на холст и соединяем с виджетом *Data Table* (рис.9).



Рис.9. Добавление виджета *k-Means*

11) Открываем виджет *k-Means*, выбираем количество кластеров 4 (рис.10).



Рис.10. Настройка виджета *k-Means*

12) Для того чтобы просмотреть данные виджета *k-Means* добавляем на холст виджет *Data Table* и соединяем с *k-Means* (рис.11).



Рис.11. Соединение виджетов

13) Открываем виджет Data Table. В таблице появилось 2 новых столбца Cluster и Silhouette. Столбец Cluster содержит информацию к какому кластеру относится вуз (рис.12).

вуз	вуз	Cluster	size	Результат функции (D)	Поправка при сохранении качества обучения, size	принадлежность к кластеру (silhouette)	отношение к кластеру (silhouette)
39	Уральский федеральный университет	01	2,11201	2,11201	27	29	96
81	Новосибирский государственный университет	04	2,23466	2,23466	64	77	59
83	Новосибирский государственный политехн.	02	2,32846	2,32846	95	70	81
11	Уральский федеральный университет	01	1,94971	1,94971	79	100	55
12	Уральский федеральный университет	01	3,91607	3,91607	15	3	29
15	Уральский федеральный университет	01	2,4279	2,4279	65	57	81
18	Уральский федеральный университет	01	2,91636	2,91636	18	10	74
21	Уральский федеральный университет	01	2,65346	2,65346	74	86	102
84	Уральский государственный медицинский ун-верситет (УГМУ)	02	1,80663	1,80663	4	40	3
88	Технический институтский университет	03	2,21806	2,21806	118	46	101
70	Технический государственный университет	03	2,11872	2,11872	75	76	71
90	Технический государственный инженерный у-	03	1,80827	1,80827	93	81	80
96	Технический государственный университет	03	2,27978	2,27978	88	80	50
98	Технический государственный университет	03	2,70998	2,70998	98	100	713
99	Технический государственный университет	03	1,80848	1,80848	71	100	128
79	Технический государственный университет	03	2,70976	2,70976	87	72	62
80	Технический государственный университет	03	2,30909	2,30909	90	74	44
13	Образовательный государственный аграрный	02	1,39118	1,39118	28	19	99
14	Образовательный государственный аграрный ун-	02	1,40371	1,40371	44	62	80
16	Образовательный государственный аграрный ун-	02	1,74636	1,74636	65	67	74
86	Образовательный государственный аграрный ун-	02	2,20371	2,20371	58	86	81
86	Образовательный государственный аграрный ун-	02	2,23464	2,23464	44	50	83
86	Образовательный государственный аграрный ун-	02	2,21718	2,21718	58	64	104
81	Саратовский национальный исследовательский	02	2,20876	2,20876	68	117	30
8	Санкт-Петербургский политехнический универ-	01	4,14146	4,14146	12	12	0
24	Санкт-Петербургский государственный уни-	04	2,50891	2,50891	32	79	41
24	Санкт-Петербургский государственный уни-	04	2,87313	2,87313	35	30	82
4	Санкт-Петербургский государственный уни-	01	4,40201	4,40201	3	71	9
69	Санкт-Петербургский государственный уни-	02	2,21938	2,21938	52	90	100
81	Санкт-Петербургский государственный уни-	04	2,45425	2,45425	67	26	76
81	Санкт-Петербургский государственный уни-	04	2,86988	2,86988	28	81	83
46	Самарский государственный технический уни-	02	2,80793	2,80793	73	81	23
47	Самарский государственный технический уни-	02	2,1899	2,1899	121	40	47
80	Самарский государственный технический уни-	04	2,45486	2,45486	40	74	73
78	Рязанский государственный университет	03	2,0081	2,0081	65	44	122
22	Рязанский государственный университет	03	1,89988	1,89988	13	8	30
41	Рязанский государственный университет	04	2,3181	2,3181	46	92	87

Рис.12. Данные виджета k-Means

14) Для того чтобы визуально посмотреть распределение вузов по схожим признакам (баллов), на холст добавляем виджет scarlett plot и соединяем с виджетом k-Means.

15) Открываем Scatter Plot (Рис.13).



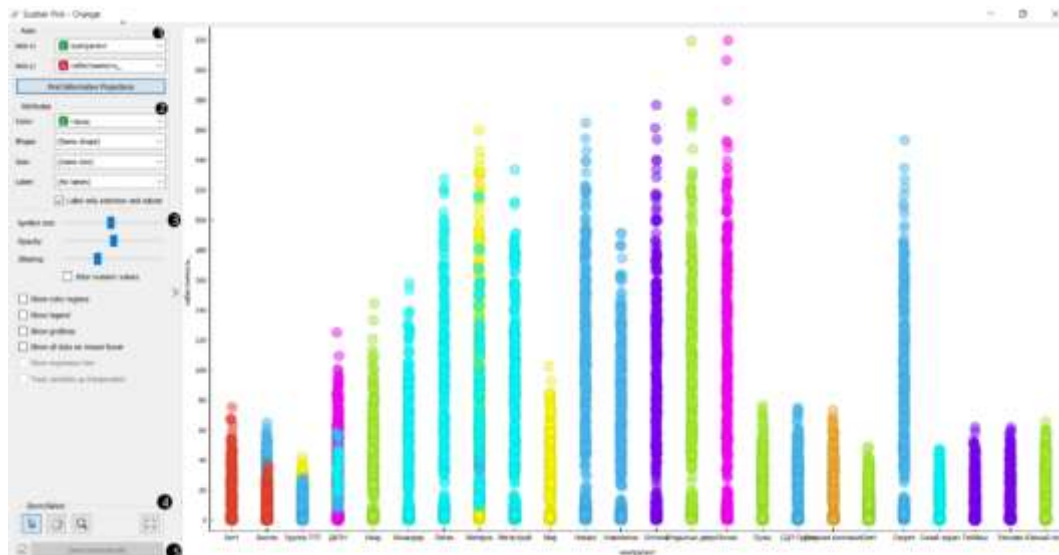


Рис.13. Диалоговое окно Scatter Plot

1. Выбор атрибутов  $x$  и  $y$ . Оптимизируйте свою проекцию с помощью функции поиска информативных проекций. Эта функция оценивает пары атрибутов по средней точности классификации и возвращает пару с наибольшим баллом с одновременным обновлением визуализации.

2. Атрибуты: Задайте цвет отображаемых точек (вы получите цвета для категориальных значений и сине-зелено-желтые точки для числовых). Задайте метку, форму и размер, чтобы различать точки. Метка только выбранных точек позволяет выбирать отдельные экземпляры данных и помечать только их.

3. Задайте размер и непрозрачность символа для всех точек данных. Установите дрожание, чтобы предотвратить перекрытие точек. Дрожание будет случайным образом рассеивать точки только вокруг категориальных значений. Если установлен флажок Числовые значения Джиттера, точки также рассеиваются вокруг их фактических числовых значений.

- Отображение цветовых областей цвета графика по классам.
- Показать легенду отображает легенду справа. Щелкните и перетащите легенду, чтобы переместить ее.
- Показать линии сетки отображает сетку за графиком.
- Показывать все данные при наведении указателя мыши позволяет создавать информационные пузырьки, если курсор помещен на точку.
- Показать линию регрессии рисует линию регрессии для пары числовых атрибутов. Если для раскраски графика выбрана категориальная переменная, будут отображаться отдельные линии регрессии для каждого значения класса.
- Отношение к переменным как к независимым соответствиям линии регрессии к группе точек (минимизация расстояния от точек), а не как к функции  $x$  (минимизация вертикальных расстояний).

1. Выбор, масштабирование, панорамирование и масштабирование по размеру - это варианты изучения графика. Ручной выбор экземпляров данных работает как инструмент углового/квадратного выделения. Дважды

щелкните, чтобы переместить проекцию. Прокрутка вверх или вниз для масштабирования.

2. Если установлен флажок отправить автоматически, изменения сообщаются автоматически. Также можно нажать кнопку Отправить.

16) Используя Scarlett plot проведем кластерный анализ данных Рейтингового функционала вуза, для этого выберем атрибуты в axis x «Cluster», а в axis y «Рейтинговый функционал вуза». С помощью графика можно увидеть, что группа кластера 1 (C1) имеет наивысший балл по рейтинговому функционалу (Рис.14).

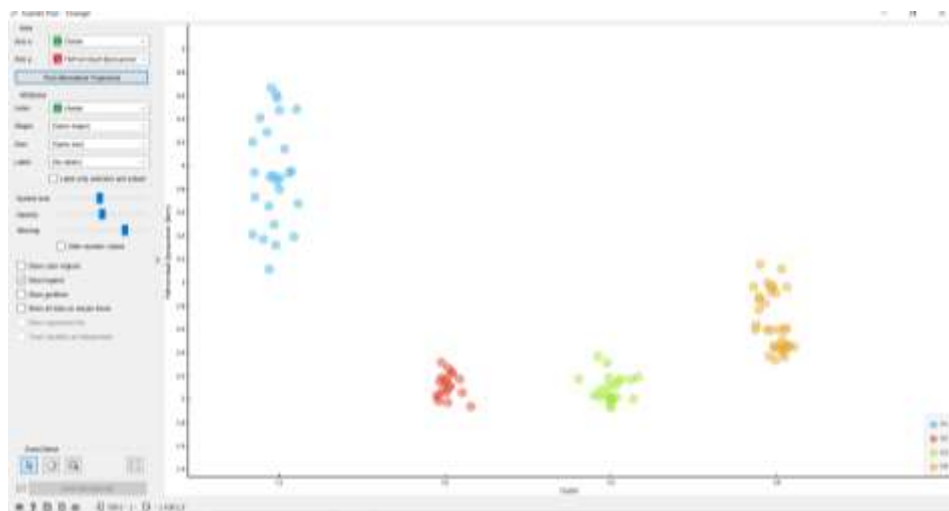


Рис.14. Кластерный анализ данных Рейтингового функционала вуза

17) По такому же принципу проведем кластерный анализ данных условий для получения качественного образования вуза, для этого выберем атрибуты в axis x «Cluster», а в axis y «Условия для получения качественного образования». С помощью графика можно увидеть, что группа кластера 3 (C3) имеет наивысший балл по получению качественного образования вуза (Рис.15).

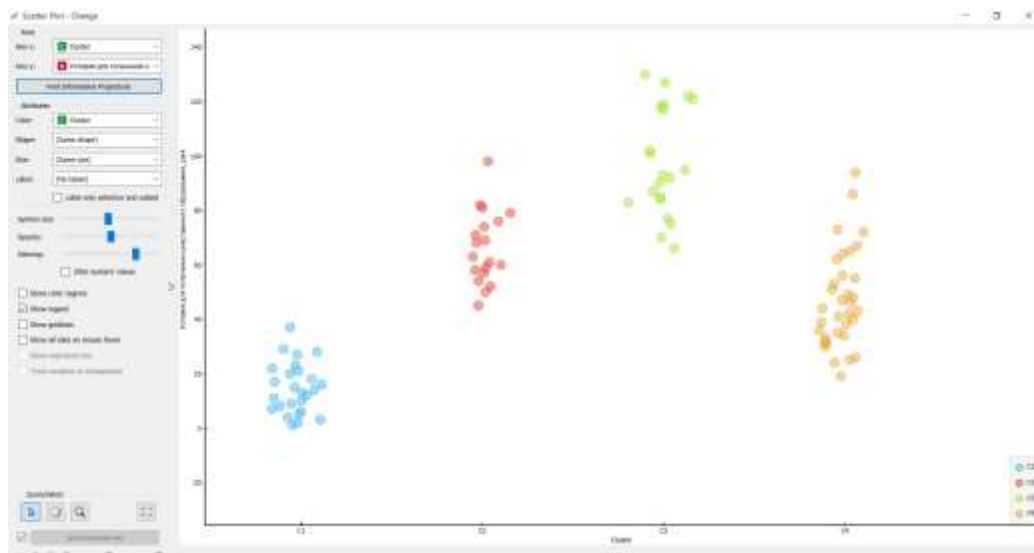


Рис.15. Кластерный анализ данных условий для получения качественного образования вуза

18) Также проведем кластерный анализ уровня востребованности выпускников среди работодателей, выберем атрибуты в axis x «Cluster», а в axis y «Уровень востребованности выпускников работодателями». На графике можно увидеть, что группа кластера 2 (C2) имеет наивысший балл среди других университетов по востребованности среди работодателей (Рис.16).

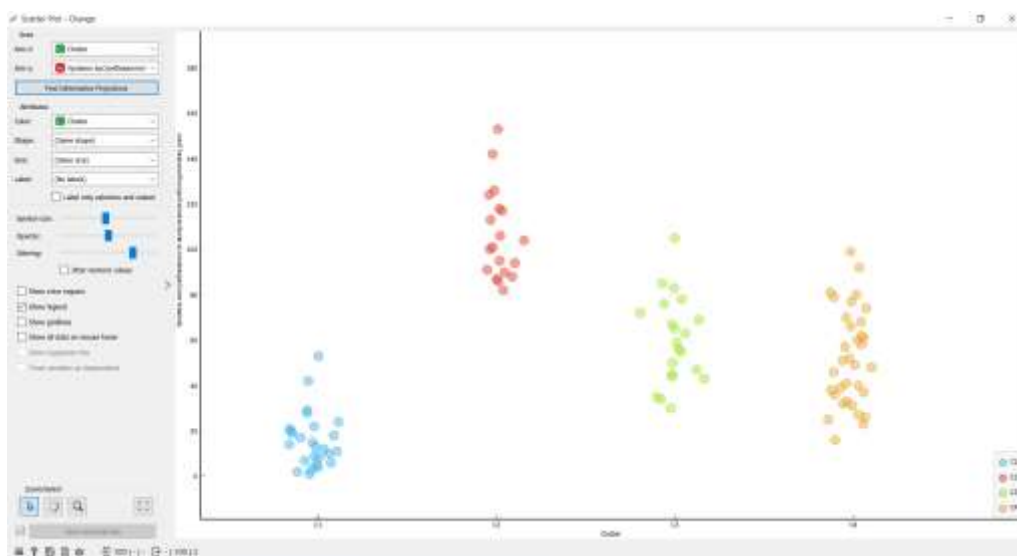


Рис.16. Кластерный анализ уровня востребованности выпускников среди работодателей

19) Далее проведем еще один кластерный анализ уровня научно-исследовательской деятельности, в атрибуте axis y выберем «а в axis y «Уровень востребованности выпускников работодателями». На графике можно заметить, что группа кластера 2 (C2) имеет наивысший балл среди других групп кластера (рис.17).

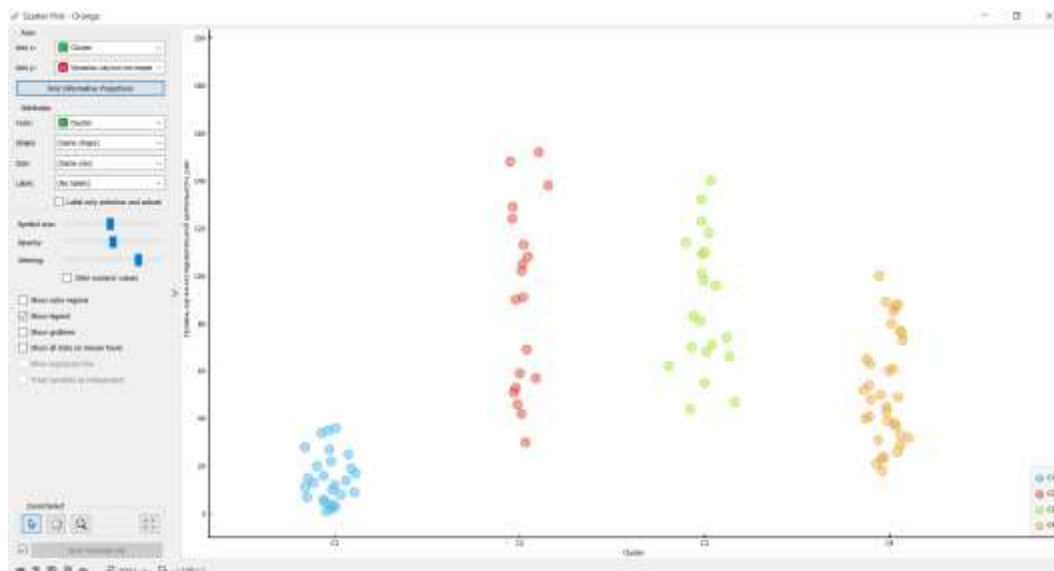


Рис.17. Кластерный анализ уровня научно-исследовательской деятельности

20) В итоге получилась готовая схема разведочного анализа (Рис.18).

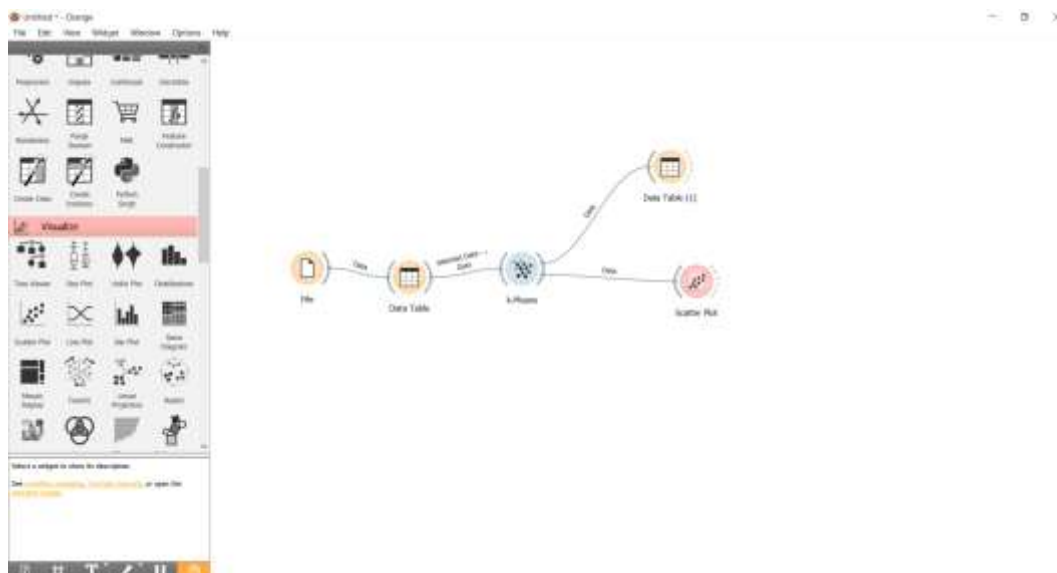


Рис.18. Итоговая схема

### Выводы

В данной работе был выполнен кластерный анализ данных с помощью программного пакета визуального программирования на основе компонентов для визуализации данных Orange. С помощью виджетов File, Data Table, k-Means, Scatter Plot и выполнили кластерный анализ данных о рейтинге лучших вузов России по разным критериям.

### Библиографический список

1. Мастевой С. С., Петрова А. Н. Data mining: обзор методов и области их применения // Наука, инновации и технологии: от идей к внедрению. 2022. С. 38-40.

2. Клименко А. В., Слащев И. С. Кластерный анализ данных //Вестник науки. 2019. Т. 1. №. 1. С. 159-163.
3. Гринченков Д. В. и др. Сравнительный анализ алгоритмов интеллектуального анализа данных //Моделирование. Теория, методы и средства. 2016. С. 263-266.
4. Маматкасымова А.Т., Кульматова Н.А. Orange: Использование системы визуального программирования при обработке больших данных// Материаловедение. 2022. №1(36). С. 2232.
5. Кривенко Ю. С., Минасян А. Т., Разиньков А. О. Исследование технологий интеллектуального анализа данных (Data Mining) // Актуальные проблемы управления в электронной экономике. 2018. С. 182-184.
6. Юсупов Н. Исследование методов кластеризации в программе Orange //Молодежная школа-семинар по проблемам управления в технических системах имени АА Вавилова. 2020. Т. 1. С. 35-37.
7. Токарев А. И., Брякин А. Н. Разведочный анализ данных и data mining // Перспективные направления развития отечественных информационных технологий. 2017. С. 201-203