

Решение задачи кластеризации набора данных различных слов с помощью программного пакета визуального программирования Orange

Голубева Евгения Павловна

Приамурский государственный университет имени Шолом-Алейхема

Студент

Аннотация

Цель данной статьи – решить задачу кластеризации набора данных состоящий из различных слов. Для решения задачи кластеризации был использован программный пакет визуального программирования на основе компонентов для визуализации данных Orange и набор данных различных слов. С помощью средств визуализации Orange решили задачу кластеризации набора данных состоящий из различных слов.

Ключевые слова: Orange, виджет, слова, кластеризация.

Solving the problem of clustering a dataset of different words using the Orange visual programming software package

Golubeva Evgeniya Pavlovna

Sholom-Aleichem Priamursky State University

Student

Abstract

The purpose of this article is to solve the problem of clustering a dataset of different words. To solve the clustering problem, a visual programming software package based on Orange data visualization components and a dataset of various words were used. Using Orange visualization tools, we solved the problem of clustering a dataset of different words.

Keywords: Orange, widget, words, clustering.

1 Введение

1.1 Актуальность

Кластеризация является важным методом анализа данных, позволяющим группировать объекты на основе их сходства. Применение кластеризации к набору данных различных слов может привести к появлению новых знаний и пониманию сходств и различий между этими словами.

Программный пакет визуального программирования Orange предоставляет удобный и интуитивно понятный интерфейс для решения задач анализа данных, включая кластеризацию. Благодаря этому пакету исследователи и специалисты в области анализа данных могут легко применять методы кластеризации к наборам данных, в том числе и к наборам из различных слов.

1.2 Обзор исследований

Е.А. Долгих, Т.А. Першина, Л.А. Давлетшина показали методы многомерного анализа данных: кластерный анализ, дерево решений и случайный лес, анализ изображений [1]. Рассмотрел программные комплексы, включающие в свою функциональность проведение какого-либо вида кластерного анализа, изучил их особенности, перечислил преимущества и недостатки А.В. Гладилин [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 рассмотрел использование методов кластеризации в программе Orange на основе реальной базы данных. Н. Юсупов [4]. А. В. Леонов в статье рассматривал основные алгоритмы кластеризации категориальных данных применительно к различным типам пользовательских интерфейсов, определяются их достоинства и недостатки [5].

1.3 Цель исследования

Цель исследования - решить задачу кластеризации набора данных состоящий из различных слов.

2 Материалы и методы

Для решения задачи кластеризации используется программа Orange. Работа будет происходить на готовом наборе данных состоящий из различных слов, скачать которые можно по ссылке:

<https://view.officeapps.live.com/op/view.aspx?src=https%3A%2F%2Ffile.biolab.si%2Fdatasets%2Fwords.xlsx&wdOrigin=BROWSELINK>

3 Результаты и обсуждения

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

Создадим новый файл (см.рис.1).



Рисунок 1- Создание нового файла

Для решения задачи кластеризации необходимо установить дополнение Text. Для того, чтобы скачать дополнение, необходимо перейти в Options, далее в Add-ons, в появившемся окне выбираем Text (см.рис.2).

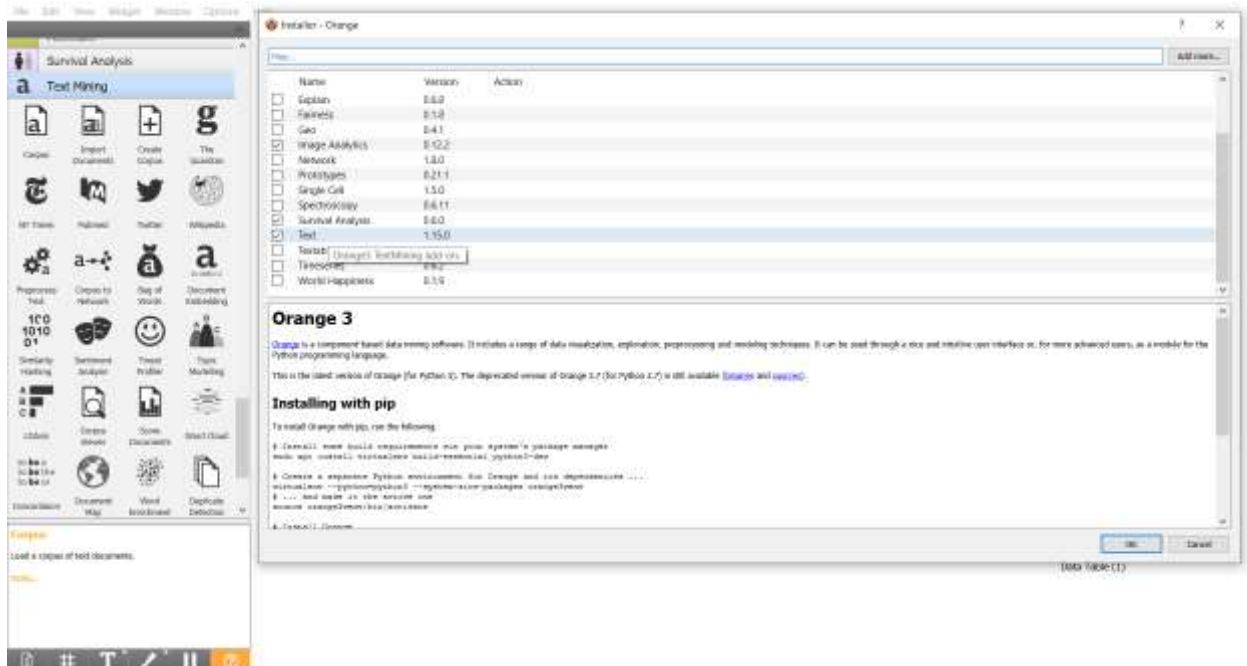


Рисунок 2 - Установка дополнения Text

Для того, чтобы загрузить набор данных состоящий из различных слов, необходимо из раздела Text Mining выбрать виджет Corpus и перенести его на холст (см.рис.3).



Рисунок 3 - Добавление виджета Corpus на холст

Открываем виджет Corpus и добавляем набор данных words.xlsx (см.рис.4).

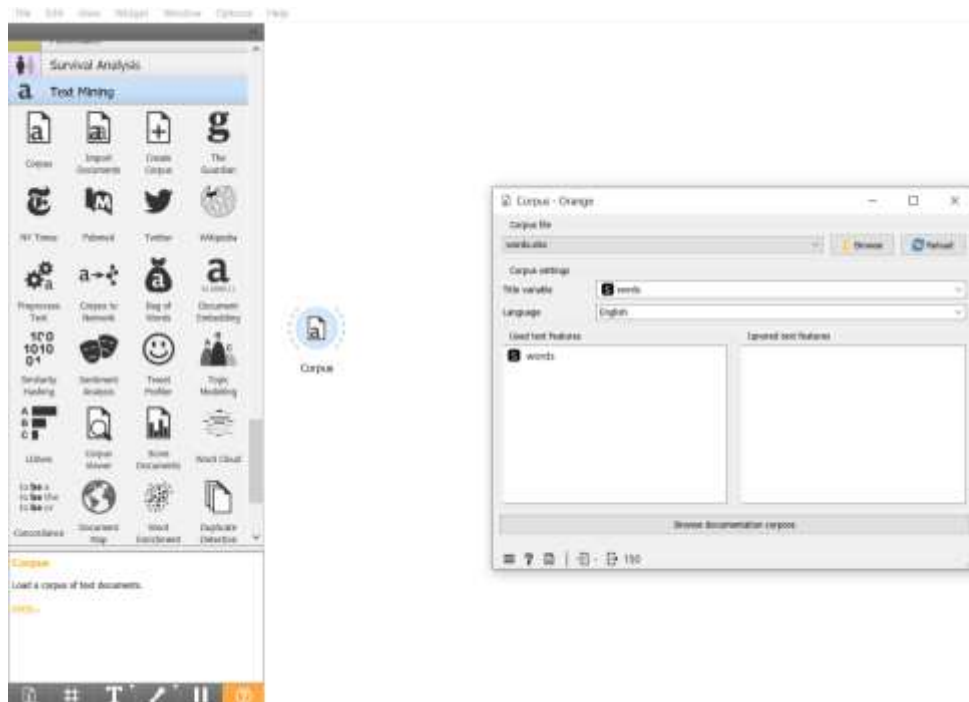


Рисунок 4 - Добавление набора данных words.xlsx

Далее добавляем виджет Corpus Viewer на холст, и соединяем с виджетом Corpus (см.рис.5).

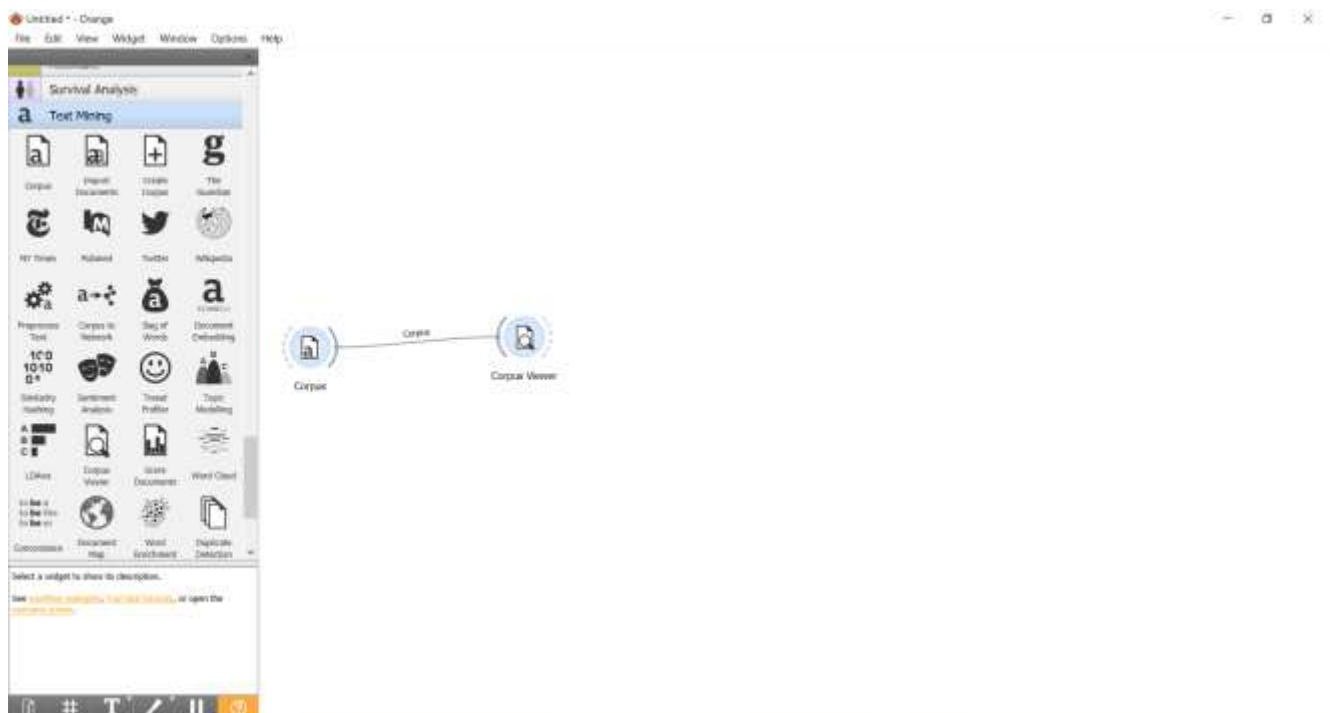


Рисунок 5 - Добавление виджета Corpus Viewer на холст

Открываем окно виджета Corpus Viewer. В открывшемся окне можно увидеть, что набор данных содержит 150 различных слов (см.рис.6)

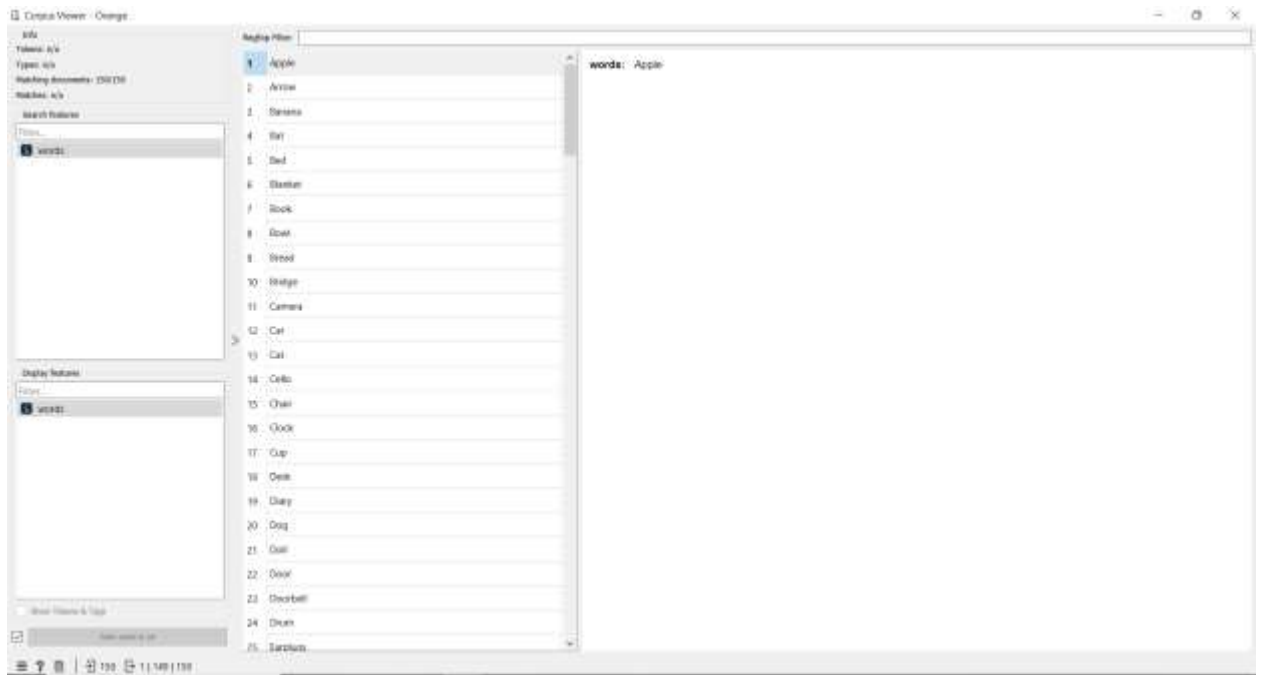


Рисунок 6 - Просмотр набора данных

Далее добавляем виджет Document Embedding на холст, и соединяем с виджетом Corpus. Виджет Document Embedding представляет слова в многомерном пространстве таким образом, что слова со схожими значениями имеют сходное вложение. Это означает, что каждое слово сопоставляется с вектором вещественных чисел, представляющих слово (см.рис.7).

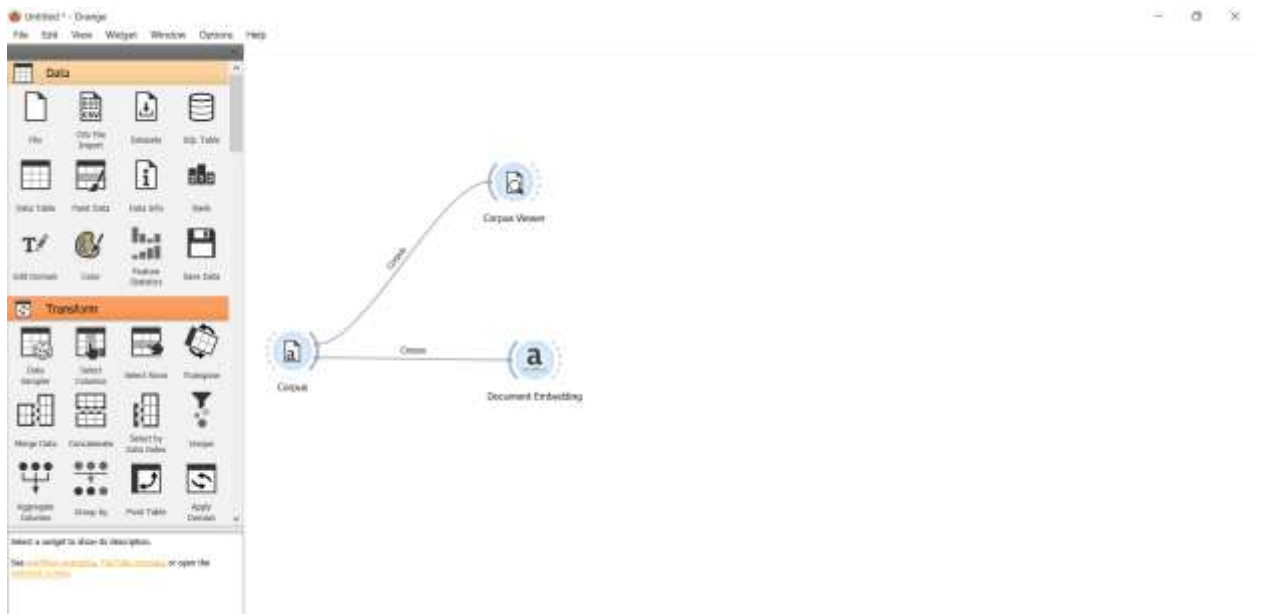


Рисунок 7 - Добавление виджета Document Embedding на холст

Открываем виджет Document Embedding, и в появившемся окне выбираем fastText (см.рис.8).

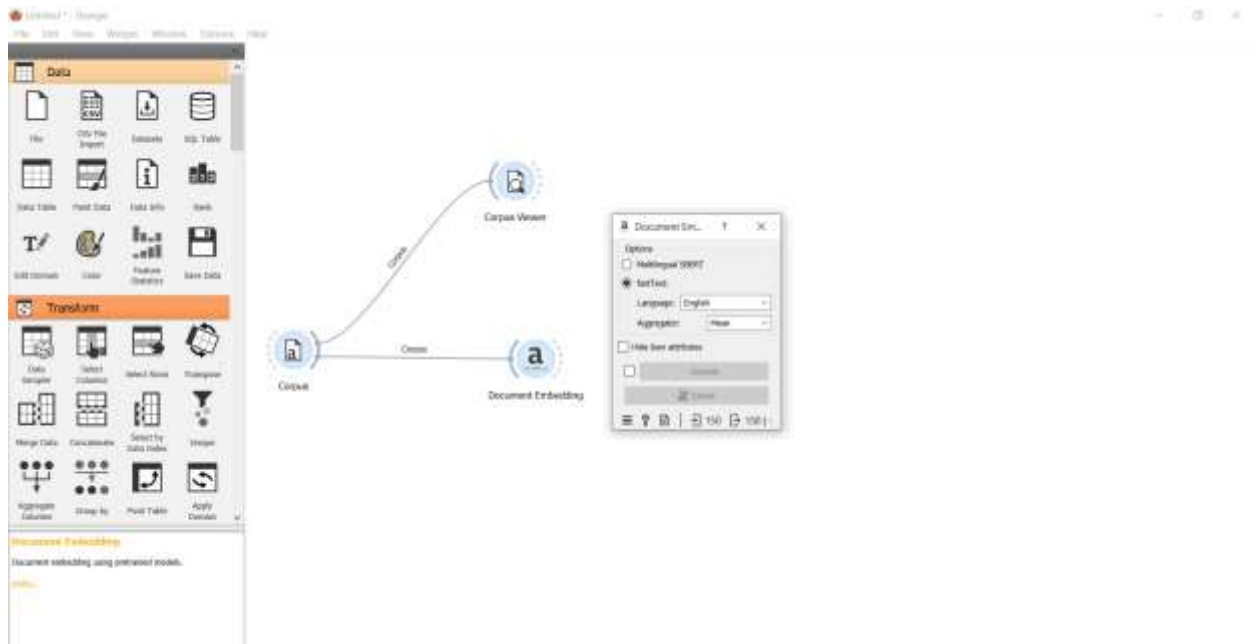


Рисунок 8 - Изменение настроек виджета Document Embedding

Добавим виджет Data Table на холст, и соединим с виджетом Document Embedding, для того чтобы посмотреть данные виджета Document Embedding (см.рис.9).

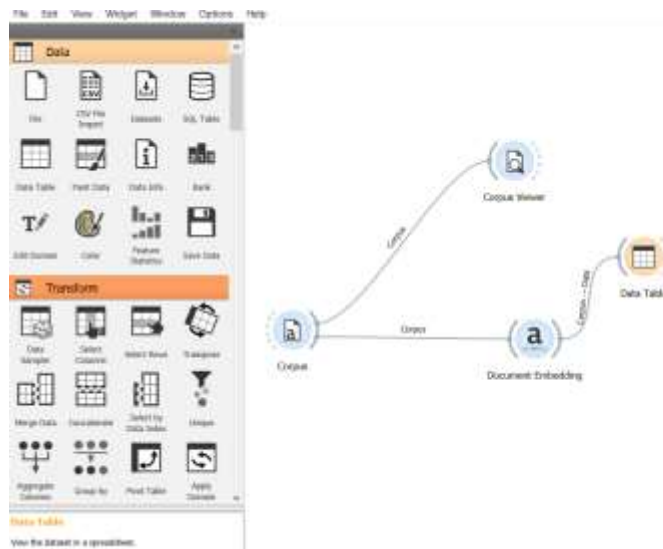


Рисунок 9 - Добавление виджета Data Table

Открываем виджет Data Table. С помощью таблицы можем увидеть, что для каждого слова добавились 300 дополнительных признаков (см.рис.10).

	embeddings	Features	Doc1 True Value	Doc2 True Value	Doc3 True Value	Doc4 True Value	Doc5 True Value	Doc6 True Value	Doc7 True Value	Doc8 True Value	Doc9 True Value	Doc10 True Value	Doc11 True Value	Doc12 True Value
1	Apple		0.0144061	-0.0103407	0.0770189	0.0962327	0.0611029	-0.0605826	0.125942	-0.0456897	-0.0733104	0.0503177	-0.008206	-0.008206
2	Apple		-0.0248477	-0.0406891	0.106488	0.092464	0.0782177	-0.0882894	0.075297	-0.0237103	-0.0439557	-0.0684937	-0.0833334	-0.0833334
3	Apple		0.125815	-0.047761	-0.0138652	0.0954233	0.094663	-0.0848803	0.104205	0.0317418	0.113215	0.0038646	0.00101323	0.00101323
4	Apple		0.082825	-0.152479	-0.0719205	0.230678	-0.272549	-0.405306	0.187599	-0.234859	0.0002458	-0.0940839	-0.274227	-0.274227
5	Apple		-0.0308096	0.0876754	-0.0153359	0.148502	0.152188	0.014891	0.133136	0.110841	-0.0211907	-0.00889377	-0.302628	-0.302628
6	Apple		-0.0129948	0.0113701	0.0477443	0.068803	0.0739672	-0.065791	0.083432	-0.0266496	0.057287	-0.019549	-0.0626419	-0.0626419
7	Apple		-0.121906	0.105171	0.0548215	0.117933	0.112326	-0.084728	0.0279075	0.0466285	-0.0748174	-0.0357998	-0.0862117	-0.0862117
8	Apple		-0.142988	-0.0528668	-0.0487126	0.115545	-0.0919048	0.0351815	0.187689	0.0832881	0.108001	-0.0718201	-0.0718201	-0.0718201
9	Apple		0.026856	0.040337	-0.0487977	0.0729793	-0.0080728	0.127115	0.133223	-0.0208486	-0.034121	0.0248094	-0.003978	-0.003978
10	Apple		0.0429815	-0.030748	-0.000139989	0.0733813	-0.0121149	-0.0401393	0.0319885	-0.0114801	-0.0474602	-0.0737488	-0.007328	-0.007328
11	Apple		0.107201	0.044283	0.0974046	0.0824754	0.00003882	-0.0905546	0.0764077	-0.0102116	0.08118647	-0.0207251	-0.00124884	-0.00124884
12	Apple		0.127486	0.173652	0.0155174	0.230301	-0.20487	0.107426	0.261468	0.00272304	0.0488513	-0.290248	-0.0021758	-0.0021758
13	Apple		0.0610558	-0.028323	-0.0322681	0.205153	0.138584	-0.107725	0.101803	0.0196826	-0.103044	0.0431342	-0.148338	-0.148338
14	Apple		0.0648897	-0.0348333	-0.0117679	0.0569561	0.0643195	-0.0070254	0.0281171	0.072821	0.162947	-0.002754	-0.041153	-0.041153
15	Apple		0.121511	0.0858865	-0.00539952	0.0474481	0.0133746	-0.091892	0.133039	0.0638031	0.0246134	0.0728636	-0.114619	-0.114619
16	Apple		0.152432	0.028642	0.0185232	0.0807913	-0.00747076	0.14428	0.00841824	0.018423	0.062521	0.082476	-0.097724	-0.097724
17	Apple		0.0351856	0.01320	-0.248407	0.122801	-0.0480127	0.144851	0.040148	0.0150823	0.05121	-0.188586	-0.003308	-0.003308
18	Apple		0.087	0.0994515	0.0530782	0.0388728	0.0075269	-0.133431	0.245062	0.111498	-0.0788	0.0680158	-0.120571	-0.120571
19	Apple		0.090208	0.103172	-0.0250191	0.0631701	0.0949189	0.0281949	0.0018234	0.05559	0.08254785	0.0071836	-0.0018848	-0.0018848
20	Apple		0.167972	-0.00128181	0.0181908	0.277907	-0.106211	-0.095717	-0.204739	0.0641571	-0.043053	-0.232542	-0.232542	-0.232542
21	Apple		-0.00118447	0.0529009	-0.00389904	0.183462	0.0166256	-0.2827313	0.151684	-0.0517531	0.0352273	0.074872	-0.060806	-0.060806
22	Apple		0.00154471	-0.0324318	0.0862192	0.124277	-0.052089	-0.00196252	0.115688	-0.07541	0.0238197	0.125152	-0.093184	-0.093184
23	Apple		0.0197156	-0.0294415	0.0117346	0.0411715	-0.090829	0.0525437	0.018951	0.0349462	0.0515931	0.0994639	0.002434	0.002434
24	Apple		-0.011626	0.0040809	0.0313933	0.244262	-0.0135643	0.124802	0.0579891	0.0830944	0.039423	-0.136988	-0.178714	-0.178714
25	Apple		-0.031207	-0.0403963	0.0198711	0.0400449	0.0509133	-0.0785797	0.0242889	0.0470754	0.151177	-0.0352081	0.0392589	0.0392589
26	Apple		0.0137874	-0.0259603	0.0184685	0.123841	0.0113934	-0.0421838	-0.0300023	-0.0068001	0.0788007	0.0710894	-0.703773	-0.703773
27	Apple		0.0647491	-0.011807	0.0087755	0.0421117	0.0385032	0.0390403	0.0238713	0.0177578	0.0458807	-0.030571	0.087	0.087
28	Apple		0.270674	0.044218	-0.238871	0.262304	0.0282625	-0.003222	0.186259	0.0908802	0.0332263	-0.208801	0.0761294	0.0761294
29	Apple		0.0408821	-0.054889	0.0180282	0.0781888	0.0137115	-0.0809102	0.031073	0.0154881	-0.079188	-0.030304	-0.067368	-0.067368
30	Apple		-0.141486	0.06351712	0.036009	0.0874453	0.157036	-0.019349	0.094471	0.094772	0.10986	0.0580181	-0.061804	-0.061804
31	Apple		-0.047505	-0.078829	-0.0207131	0.048863	0.071979	-0.086558	0.0081365	0.050272	-0.105266	-0.112526	-0.109409	-0.109409
32	Apple		-0.151585	-0.0719667	-0.0497386	0.0413607	0.0388627	-0.078682	0.106187	-0.077128	0.0088113	0.008453	-0.0417053	-0.0417053
33	Apple		0.0270312	-0.018212	0.0215625	0.101708	-0.0717159	-0.0478881	0.0268774	0.0199228	-0.0271112	0.0035885	-0.0811732	-0.0811732
34	Apple		0.00745472	0.06201772	0.0155884	0.0404372	0.0400682	-0.0737818	0.0001184	0.0248201	0.113882	-0.0078739	-0.0222231	-0.0222231

Рисунок 10 - Просмотр данных виджета Document Embedding

Далее добавляем виджет Distances на холст, и соединим с виджетом Document Embedding. Откроем виджет Distances, и выберем метрику расстояния cosine (см.рис.11).

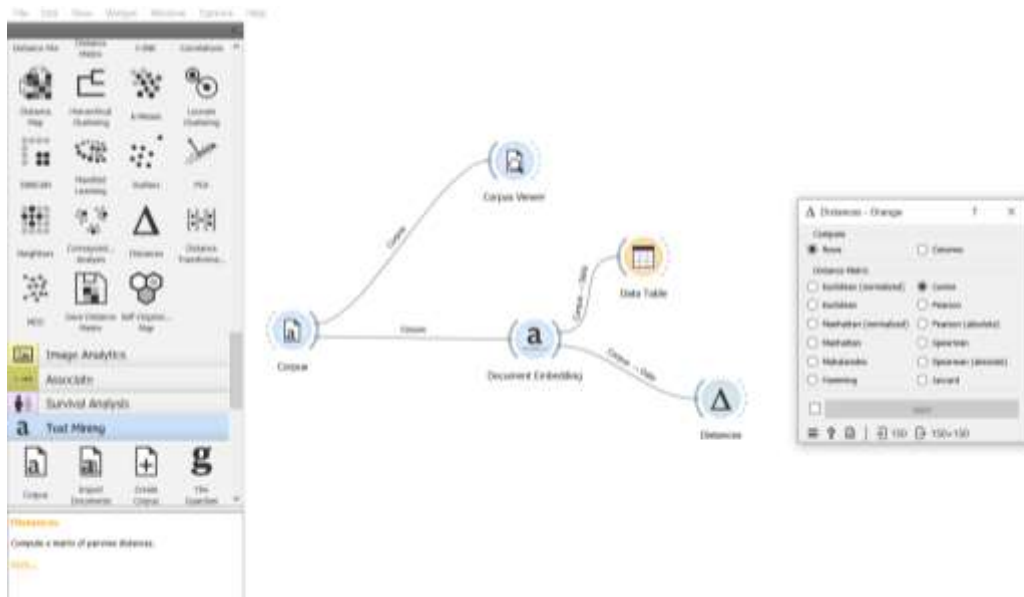


Рисунок 11 - Добавление виджета Distances на холст

Для того чтобы посмотреть данные результата виджета Distances, добавим виджет Distance Matrix и соединим с виджетом Distances (см.рис.12)

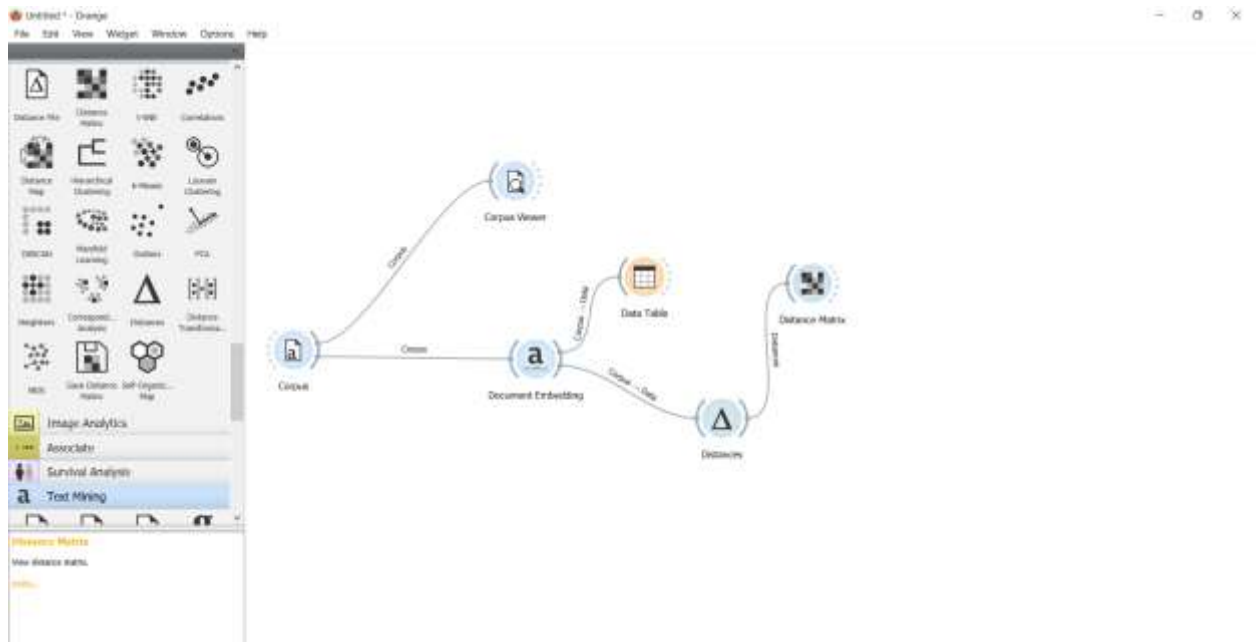


Рисунок 12 - Добавление виджета Distance Matrix на холст

На таблице матрица расстояния можем определить, какое расстояние между словами. По матрице можно определить, чем меньше значение между словами, тем больше они схоже по признакам (см.рис.13).

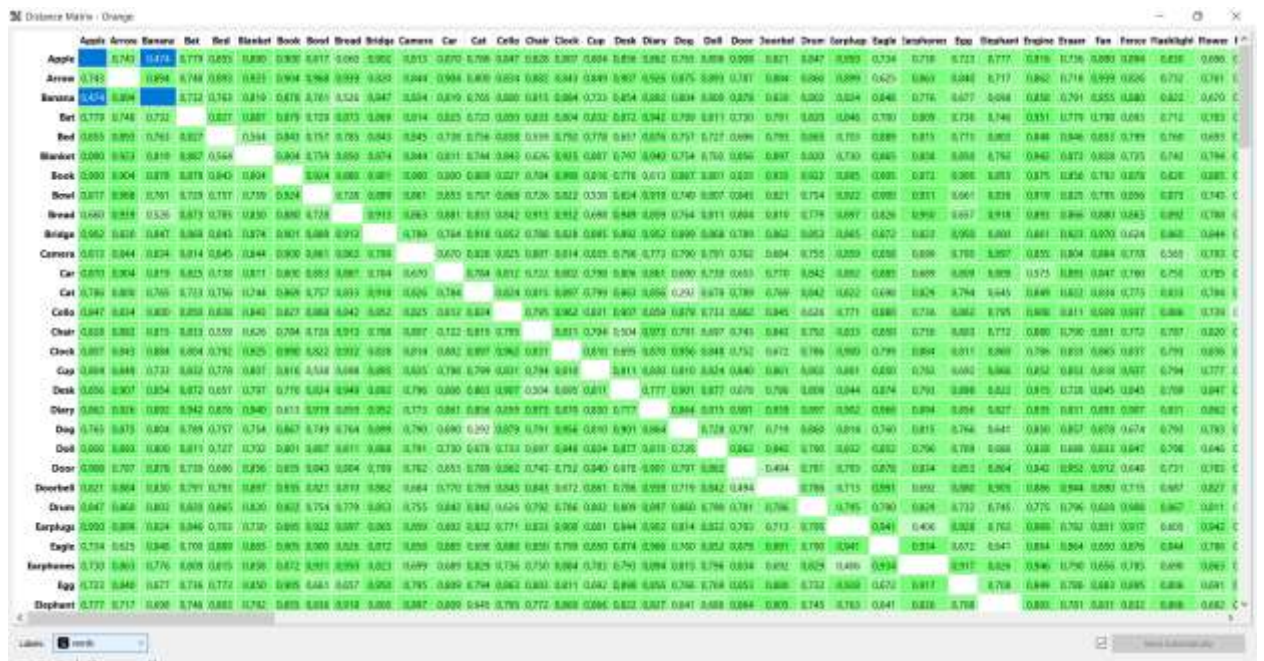


Рисунок 13 - Матрица расстояния

Далее проведем кластеризацию слов. Для этого добавим виджет Hierarchical Clustering на холст и соединим с виджетом Distances (см.рис.14).

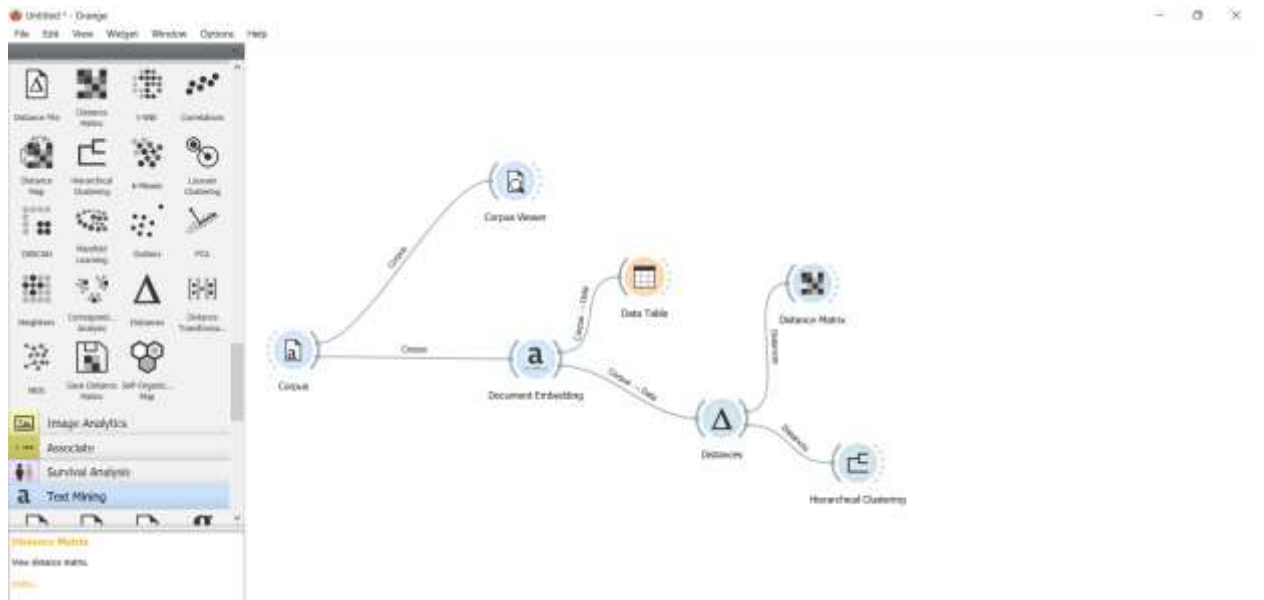


Рисунок 14 - Добавление виджета Hierarchical Clustering на холст

Открываем виджет Hierarchical Clustering. В появившемся окне можно увидеть, как слова распределились на 5 кластеров. Например, кластер c1 можно отнести к группе музыкальные инструменты, кластер c2 можно отнести к группе вещи (см.рис.15).

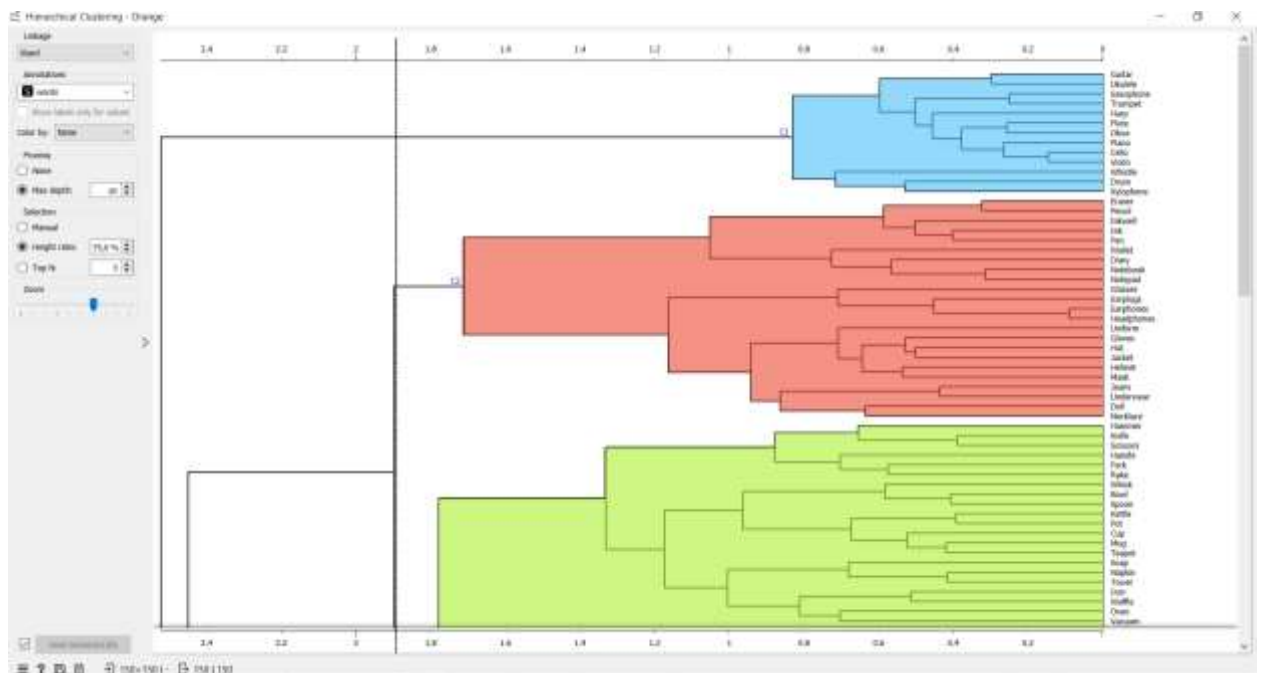


Рисунок 15 - Распределение слов по кластерам

Также можно визуализировать кластерный анализ. Для этого добавляем виджет t-SNE на холст, и соединяем с виджетами Document Embedding и Hierarchical Clustering (см.рис.16).

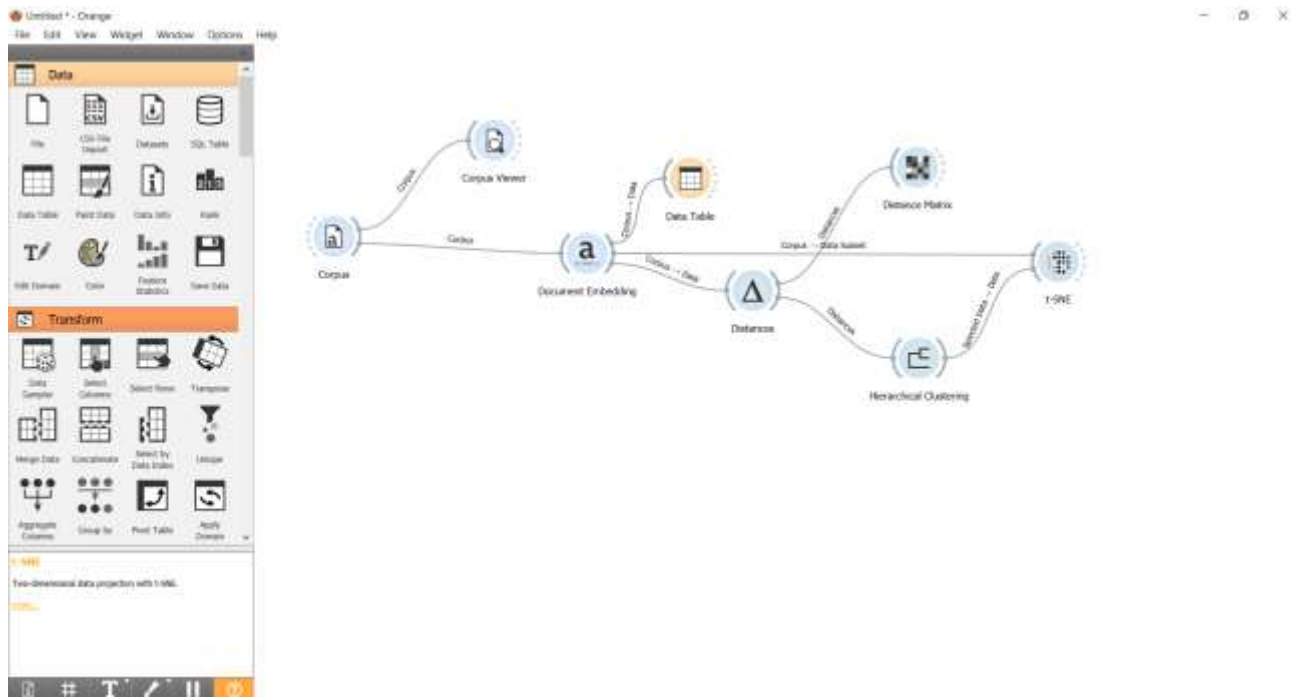


Рисунок 16 - Добавление виджета t-SNE на холст

Открываем виджет t-SNE. Виджет покажет 2D-карту, где образцы с похожими профилями экспрессии генов будут расположены близко друг к другу (см.рис.17).

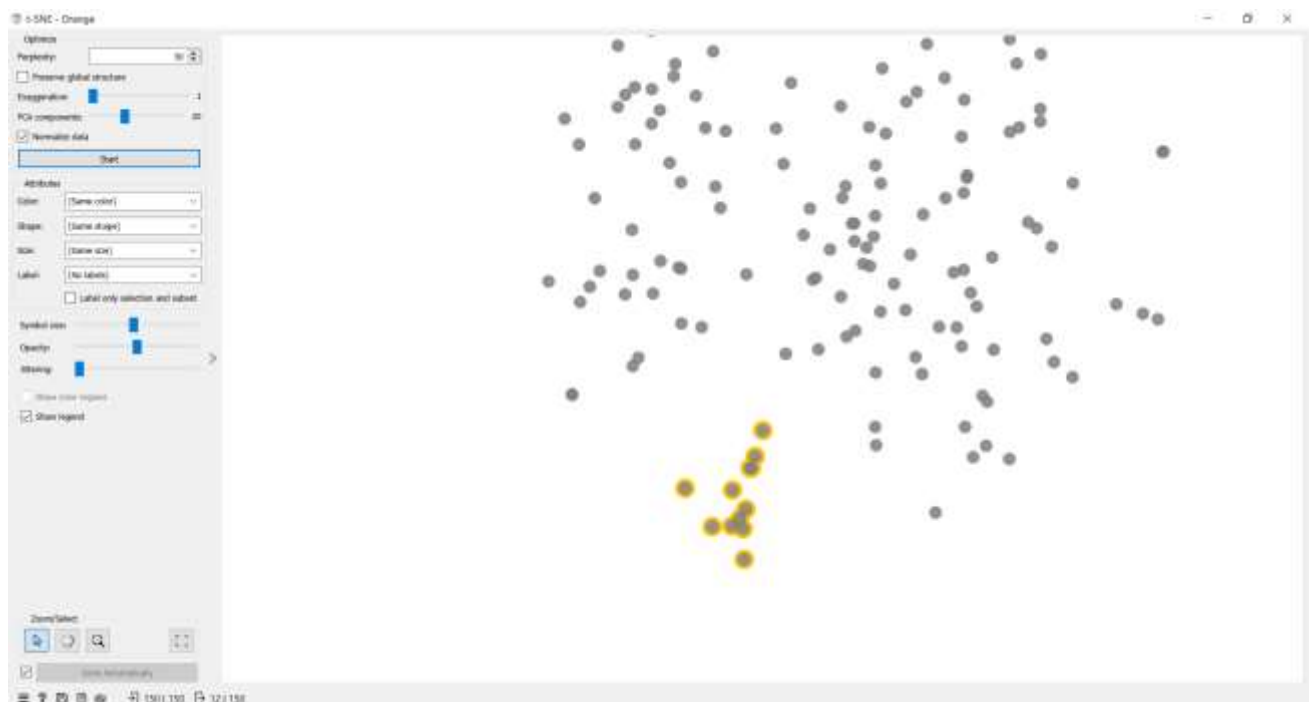


Рисунок 17 - Диалоговое окно t-SNE

Для того чтобы убедиться, что расположенные рядом точки схожи по значению добавляем на холст виджет Corpus Viewer и соединяем с виджет t-SNE (см.рис.18).

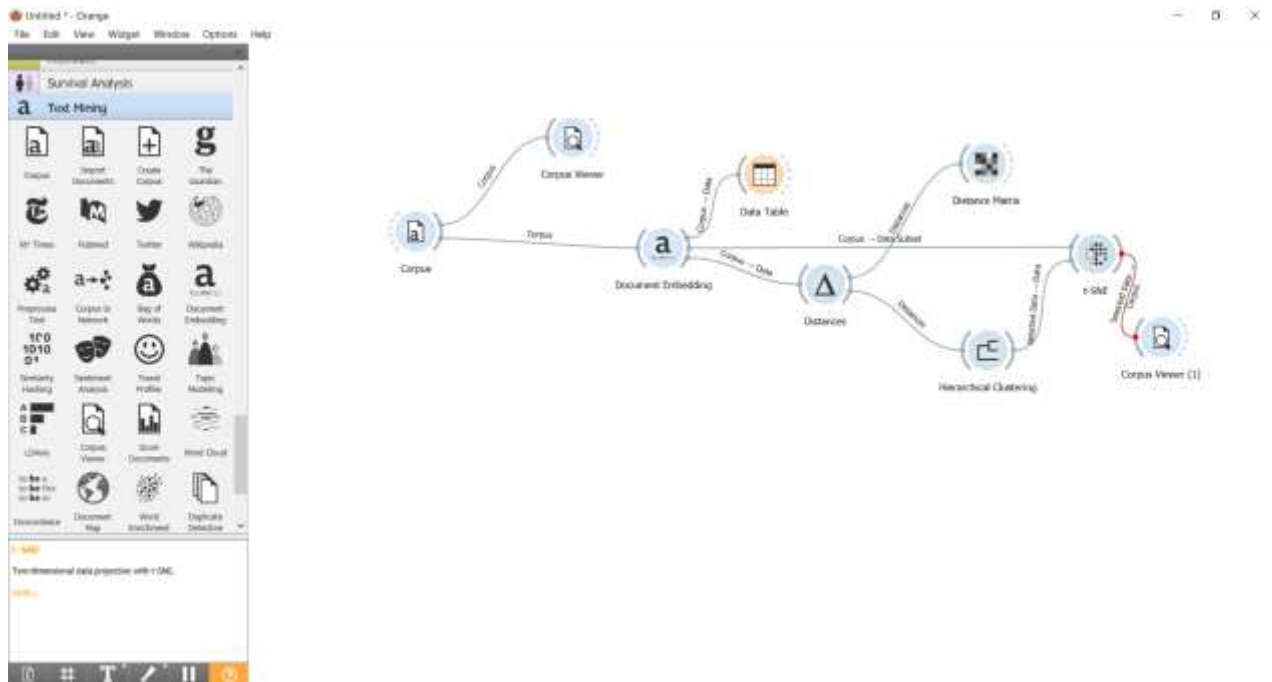


Рисунок 18 - Добавление виджета Corpus Viewer на холст

Открываем виджеты t-SNE и Corpus Viewer. Для того чтобы проверить правильность значений в виджете t-SNE выделяем группу точек, которые расположены близко к друг другу. В окне Corpus Viewer появились значение выделенных точек. Можем сделать вывод, что точки расположены близко к друг другу относятся к категории музыкальные инструменты (см.рис.19).

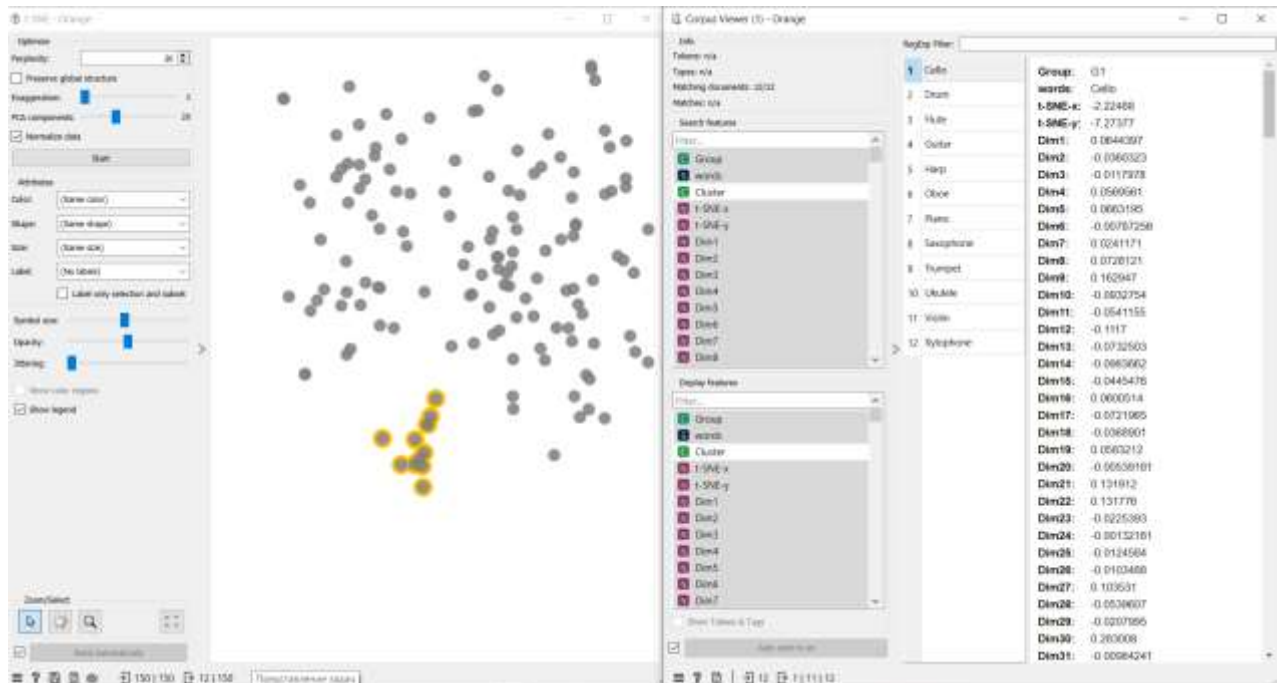


Рисунок 19 - Проверка выделенных точек

В итоге получилась готовая схема, с помощью которой можно решить задачу кластеризации различных слов (см.рис.20).

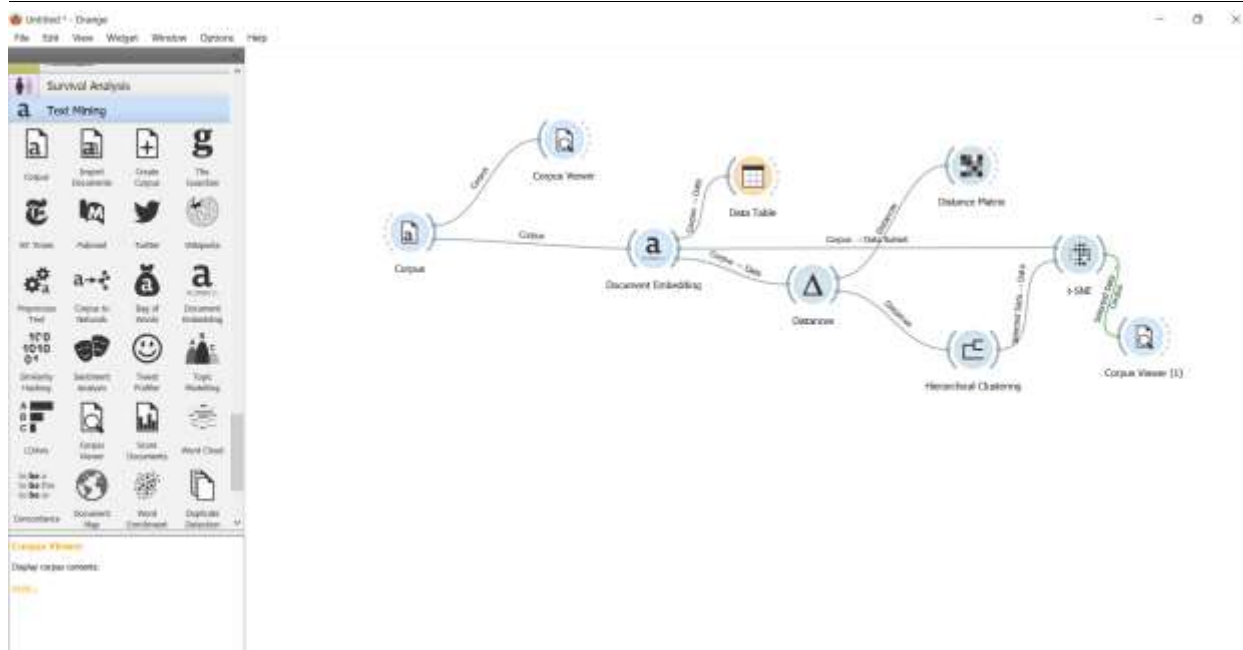


Рисунок 20 - Итоговая схема

Выводы

В данной работе была выполнена задача кластеризации набора данных различных слов с помощью программного пакета визуального программирования на основе компонентов для визуализации данных Orange. С помощью виджетов Corpus, Corpus Viewer, Document Embedding, Data Table, Distances, Hierarchical Clustering, t-SNE, Distance Matrix выполнили кластеризацию различных слов и получили итоговую схему.

Библиографический список

1. Гладилин А. В. Обзор существующих программных комплексов для проведения кластерного анализа // Новые информационные технологии и системы. 2017. С. 224-226.
2. Моисеенко Г. А. и др. Классификация и распознавание изображений живой и неживой природы // Оптический журнал. 2015. Т. 82. №. 10. С. 53-64.
3. 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.
4. Клименко А. В., Слащев И. С. кластерный анализ данных // Вестник науки. 2019. Т. 1. №. 1. С. 159-163.
5. Гринченков Д. В. и др. Сравнительный анализ алгоритмов интеллектуального анализа данных // Моделирование. Теория, методы и средства. 2016. С. 263-266.
6. 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.