

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

Отчёт
по лабораторной работе №4 «Кластеризация»
по дисциплине «Системы искусственного интеллекта»

Студентка гр. 3630201/70101 _____ О. В. Саксина

Преподаватель _____ Л. В. Уткин

Содержание

1	Задание 1	3
1.1	Постановка задачи	3
1.2	Реализация	3
2	Задание 2	3
2.1	Постановка задачи	3
2.2	Реализация	3
3	Задание 3	4
3.1	Постановка задачи	4
3.2	Реализация	4
4	Задание 4	5
4.1	Постановка задачи	5
4.2	Реализация	5
5	Задание 5	6
5.1	Постановка задачи	6
5.2	Реализация	6
	Приложение	8

1 Задание 1

1.1 Постановка задачи

Разбейте множество объектов из набора данных `pluton` в пакете «`cluster`» на 3 кластера методом центров тяжести (`kmeans`). Сравните качество разбиения в зависимости от максимального числа итераций алгоритма.

1.2 Реализация

Была обучена модель методов центров тяжести при максимальном числе итераций равном 3, 5, 10, 20. В результате каждая модель давала разбиение с отношением межкластерной дисперсии к внутрикластерной равной 92,1 %, не зависящим от данного параметра, т.к. алгоритм разбиения останавливался после второй итерации.

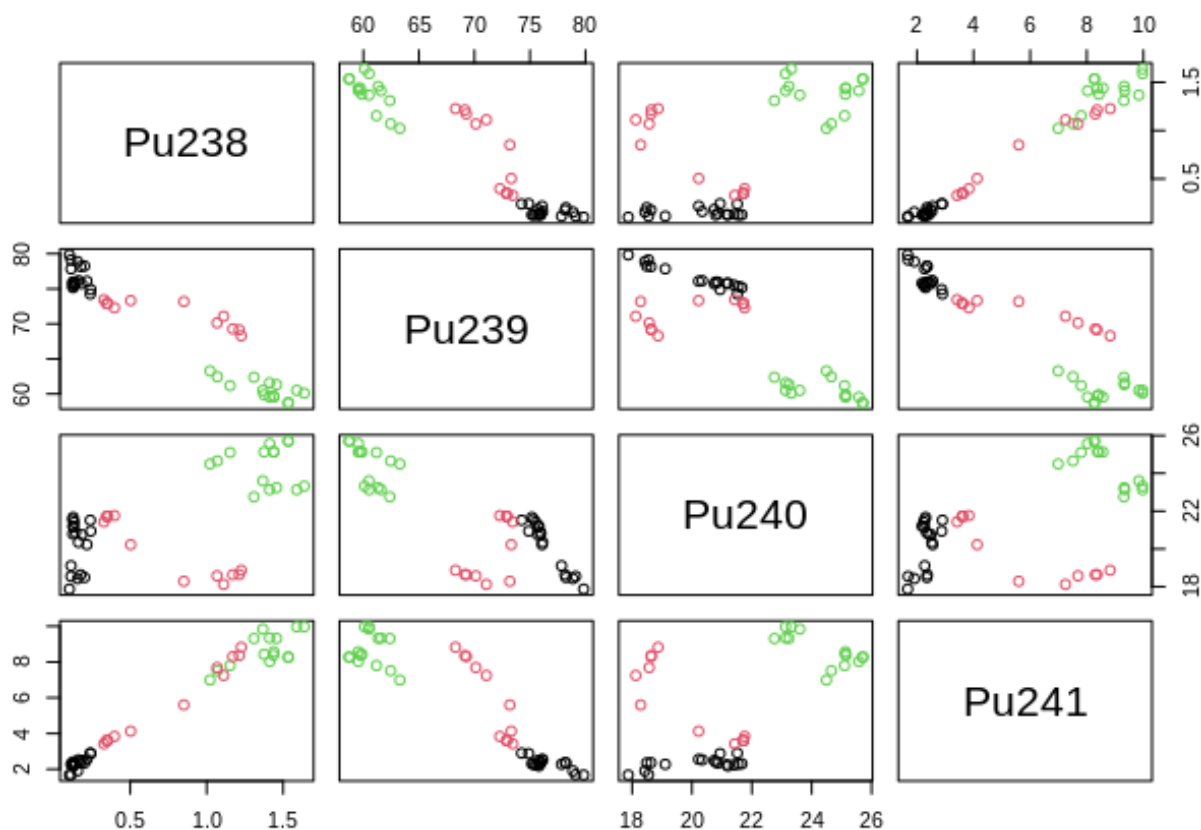


Рис. 1: Разбиение множества объектов на кластеры

2 Задание 2

2.1 Постановка задачи

Сгенерируйте набор данных в двумерном пространстве, состоящий из 3 кластеров, каждый из которых сильно “вытянут” вдоль одной из осей. Исследуйте качество кластеризации методом `slaga` в зависимости от 1) использования стандартизации; 2) типа метрики. Объясните полученные результаты.

2.2 Реализация

Был сгенерирован набор данных, представленный на рис. 2. Качество кластеризации при разных параметрах модели, определяемое как ширина силуэта, показано в таблице 1. Значение силуэта показывает, насколько объект похож на свой кластер по сравнению с другими кластерами. Для всей кластерной структуры чем ближе данная оценка к 1, тем лучше. Видно, что евклидово и манхэттонское расстояния по качеству примерно одинаковы и оценки не зависят от использования стандартизации (последнее можно объяснить характером

распределения: положение точек задавалось по нормальному распределению с одинаковыми параметрами для всех трёх кластеров), жаккардово расстояние показало худшие результаты.

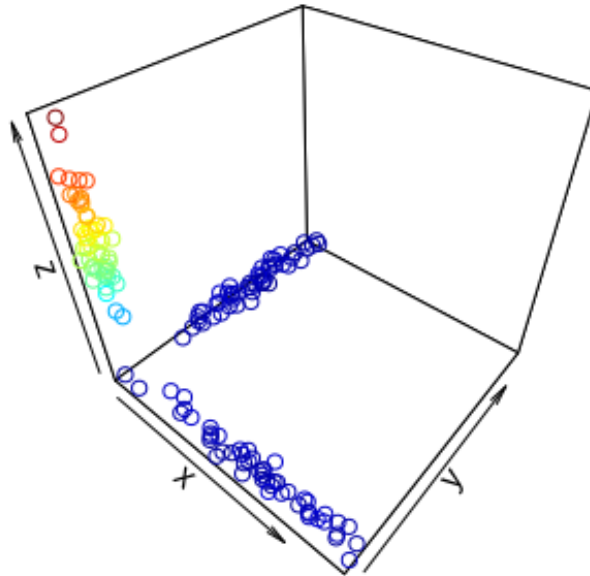


Рис. 2: Распределение данных

Использование стандартизации	Тип метрики		
	euclidean	manhattan	jaccard
TRUE	0.73	0.77	0.52
FALSE	0.73	0.77	0.01

Таблица 1: Зависимость ширины силуэта от типа метрики и использования стандартизации

3 Задание 3

3.1 Постановка задачи

Постройте дендрограмму для набора данных `votes.repub` в пакете «cluster» (число голосов, поданных за республиканцев на выборах с 1856 по 1976 год). Строки представляют 50 штатов, а столбцы - годы выборов (31). Проинтерпретируйте полученный результат.

3.2 Реализация

Построенная дендрограмма представлена на рис. 3. Дендрограмма показывает степень близости отдельных объектов и кластеров, а также наглядно демонстрирует в графическом виде последовательность их объединения или разделения. В данном случае первое разделение всего множества объектов выделяет в левое поддерево группу штатов, где было отдано меньше всего голосов (в штате Миссисипи в среднем за все годы было отдано 22.33407 голосов, что является наименьшим значением), а в правое – штаты, где больше голосов (в штате Вермонт в среднем за все годы было отдано 66.31677 голосов, что является наибольшим значением).

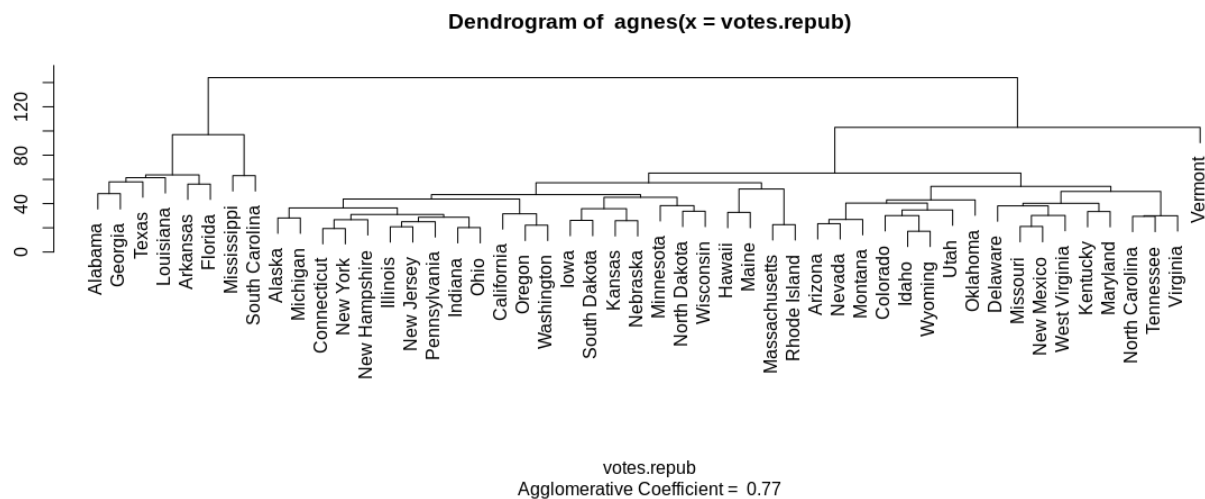


Рис. 3: Дендрограмма для набора данных votes.repub

4 Задание 4

4.1 Постановка задачи

Постройте дендрограмму для набора данных animals в пакете «cluster». Данные содержат 6 двоичных признаков для 20 животных. Переменные

- [, 1] wag теплокровные;
- [, 2] fly летающие;
- [, 3] ver позвоночные;
- [, 4] end вымирающие;
- [, 5] gro живущие в группе;
- [, 6] hai имеющие волосяной покров.

Проинтерпретируйте полученный результат.

4.2 Реализация

Построенная дендрограмма представлена на рис. 4. Каждое поддерево объединяет максимально похожих друг на друга животных. Например на втором разбиении в один кластер попали bee (пчела) и fly (муха), а в другой – cpl (гусеница) и spi (паук). Можно рассматривать эту дендрограмму, как подобие филогенетического дерева.

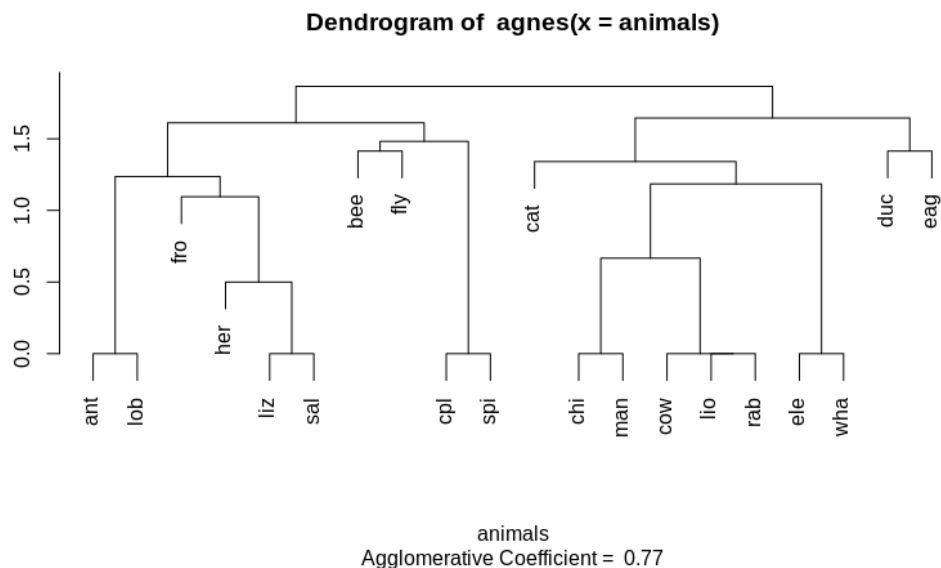


Рис. 4: Дендрограмма для набора данных animals

5 Задание 5

5.1 Постановка задачи

Рассмотрите данные из файла `seeds_dataset.txt`, который содержит описание зерен трех сортов пшеницы: `Kama`, `Rosa` and `Canadian`. Признаки:

1. область A ,
2. периметр P ,
3. компактность $C = 4 * \pi * A / P^2$,
4. длина зерна,
5. ширина зерна,
6. коэффициент ассиметрии,
7. длина колоска.

5.2 Реализация

Посмотрев на матрицу рассеяния для данных, представленную на рис. 5, можно увидеть, что есть параметры, которые лучше других объясняют разделение зёрен на сорта. Применив метод главных компонент, получим Компоненту 1, которая объясняет 71% дисперсии исходных данных, и Компоненту 2, объясняющую 17% дисперсии исходных данных. Иллюстрация этого приведена на рис. 6.

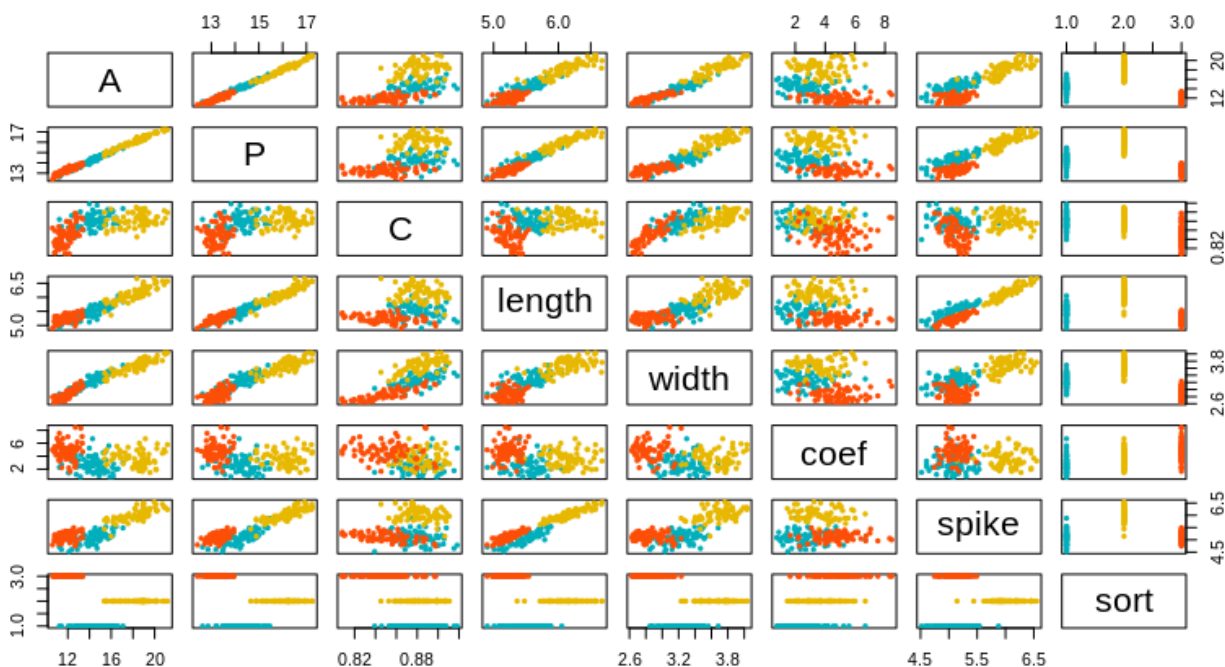


Рис. 5: Матрица рассеяния для данных `seeds_dataset`

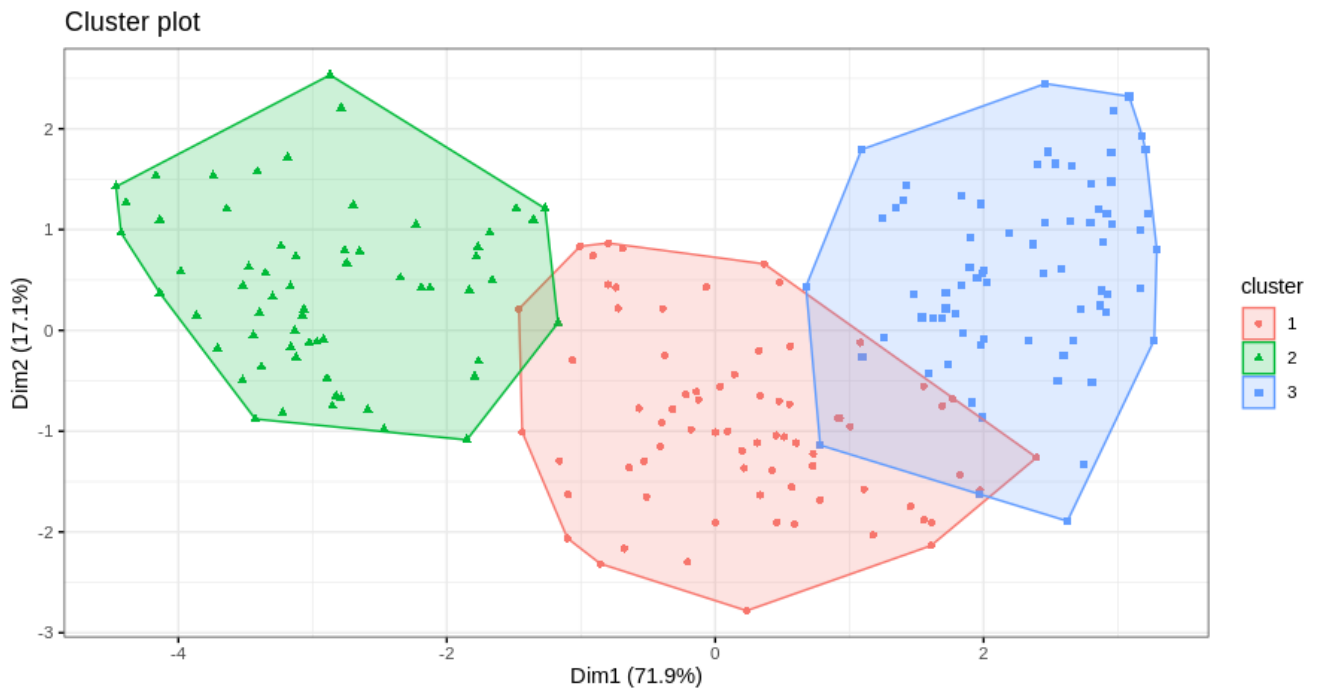


Рис. 6: Кластеризация по главным компонентам

Приложение

Задание 1

```
1 library(cluster)
2 library(factoextra)
3
4 c1 = kmeans(pluton, centers = 3, iter.max = 1)
5 c2 = kmeans(pluton, centers = 3, iter.max = 2)
6 c3 = kmeans(pluton, centers = 3, iter.max = 3,
7             algorithm = "Hartigan-Wong", trace=TRUE )
8 summary(c3)
9 fviz_cluster(c2, data = pluton,
10             geom = "point",
11             ggtheme = theme_bw()
12 )
13 plot(pluton, col = c2$cluster)
```

Задание 2

```
1 library(cluster)
2 library(scatterplot3d)
3 library("plot3D")
4 library(factoextra)
5
6 n = 50
7 mmax = 150
8 mmin = 0
9 sdmax = 50
10 sdmin = 5
11
12 data = data.frame(x = c(rnorm(n, mmax, sdmax), rnorm(n, mmin, sdmin), rnorm(n, mmin, sdmin)),
13                   y = c(rnorm(n, mmin, sdmin), rnorm(n, mmax, sdmax), rnorm(n, mmin, sdmin)),
14                   z = c(rnorm(n, mmin, sdmin), rnorm(n, mmin, sdmin), rnorm(n, mmax, sdmax)))
15 scatterplot3d(x = data$x, y = data$y, z = data$z)
16 scatter3D(x = data$x, y = data$y, z = data$z)
17
18 cl_ef = clara(data, k = 3, metric = "euclidean", stand = FALSE)
19 cl_mf = clara(data, k = 3, metric = "manhattan", stand = FALSE)
20 cl_et = clara(data, k = 3, metric = "euclidean", stand = TRUE)
21 cl_mt = clara(data, k = 3, metric = "manhattan", stand = TRUE)
22 cl_jf = clara(data, k = 3, metric = "jaccard", stand = FALSE)
23 cl_jt = clara(data, k = 3, metric = "jaccard", stand = TRUE)
24 fviz_cluster(cl_ef, data = data,
25             geom = "point",
26             ggtheme = theme_bw()
27 )
28 plot(cl_ef)
29 plot(cl_mf)
30 plot(cl_et)
31 plot(cl_mt)
32 plot(cl_jf)
33 plot(cl_jt)
34 plot(data, col=cl_ef$clustering)
```

Задание 3

```
1 library(cluster)
2
3 data = votes.repub
4 plot(agnes(votes.repub))
5 states = row.names(data)
6 means = c()
```



```

7 for (i in 1:length(states)){
8   means[i] = mean(unlist(data[states[i],]), na.rm=TRUE)
9 }
10 rank = mapply(c, states, means , SIMPLIFY = FALSE )
11 df = data.frame(states=states, means=means)
12 df = df[order(df$means),]

```

Задание 4

```

1 library(cluster)
2
3 plot(agnes(animals))
4 mona(animals)

```

Задание 5

```

1 library(cluster)
2 library(factoextra)
3
4 setwd("/home/olga/MyProjects/Polikek/ML/Cluster/datasets")
5
6 data = read.delim("seeds_dataset.txt")
7 my_cols <- c("#00AFBB", "#E7B800", "#FC4E07")
8 pairs(data[,1:8], pch = 19, cex = 0.5,
9        col = my_cols[data$sort])
10
11 sorts = data[,ncol(data)]
12 data = data[, -ncol(data)]
13
14 cl = clara(data, k = 3)
15 fviz_cluster(cl, data = data,
16              geom = "point",
17              ggtheme = theme_bw()
18 )
19 plot(cl)
20
21 set.seed(12345)
22 data2 = subset(data, select = -c(A,P))
23 k2 = kmeans(data2, centers = 3)
24 k2$betweenss/k2$totss * 100
25 tab = table(k2$cluster, sorts)
26 tab
27 cl = clara(data, k = 3)
28 fviz_cluster(cl, data = data,
29              geom = "point",
30              ggtheme = theme_bw()
31 )
32 pca <- prcomp(data, center = TRUE, scale = TRUE)
33 summary(pca)
34 pca$rotation
35 biplot(pca)

```
