

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

Отчёт
по курсовой работе
по дисциплине «Системы искусственного интеллекта»

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

Преподаватель _____ Л. В. Уткин

Санкт-Петербург,
2021г.

Содержание

1	Постановка задачи	3
2	Разработка классификаторов	3
2.1	Метод опорных векторов	4
2.2	Метод k ближайших соседей	4
2.3	Деревья решений	5
2.4	Результат	5
3	Кластеризация	6
4	Определение наиболее значимых признаков	6
	Заключение	7
	Список используемых источников	8
	Приложение	9

1 Постановка задачи

Данные для исследования в курсовом проекте взяты из базы данных «Wine recognition data». Эти данные являются результатами химического анализа вин, выращенных в одном регионе Италии, но полученных из трех разных сортов. Данные содержат следующие признаки:

1. Alcohol
2. Malic acid
3. Ash
4. Alcalinity of ash
5. Magnesium
6. Total phenols
7. Flavanoids
8. Nonflavanoid phenols
9. Proanthocyanins
10. Color intensity
11. Hue
12. OD280/OD315 of diluted wines
13. Proline

Курсовой проект заключается в разработке классификаторов для базы данных, визуализации данных, исследовании и настройке классификаторов.

В ходе работы необходимо выполнить следующие задания:

1. Разработать 3 классификатора и осуществить настройку их параметров для минимизации ошибки классификации на тестовых данных. Выполнить визуализацию данных при помощи метода t-SNE.
2. Сравнить классификаторы (по критерию вероятность ошибки классификации для тестовых данных) и обосновать выбор наилучшего из них.
3. Удалить их базы метки классов и осуществить кластеризацию данных. Построить дендограмму. Сравнить полученные результаты с реальными метками данных. Определить долю ошибочно кластеризованных данных.
4. Используя логистическую регрессию в рамках метода Лассо, определить наиболее значимые признаки, влияющие на отнесение объектов к определенному классу.

2 Разработка классификаторов

Визуализация данных, выполненная при помощи алгоритма t-SNE, представлена на рис.1

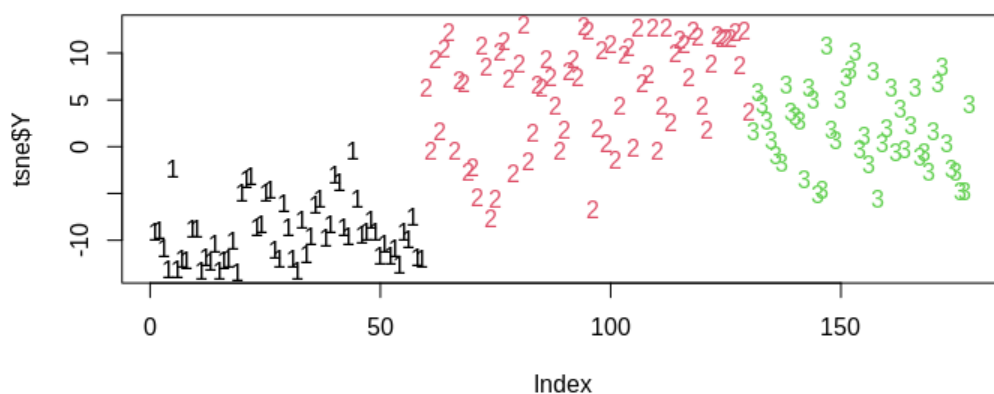


Рис. 1: Визуализация данных

Для классификации данных были обучены модели методом опорных векторов, k ближайших соседей и деревьев решений.

2.1 Метод опорных векторов

Изменение точности модели при различных параметрах, измеренное с применением скользящего контроля, представлено на рис. 2. Для модели с ядром типа «radial» и параметром C, равном 1, на тестовых данных получена точность 98%

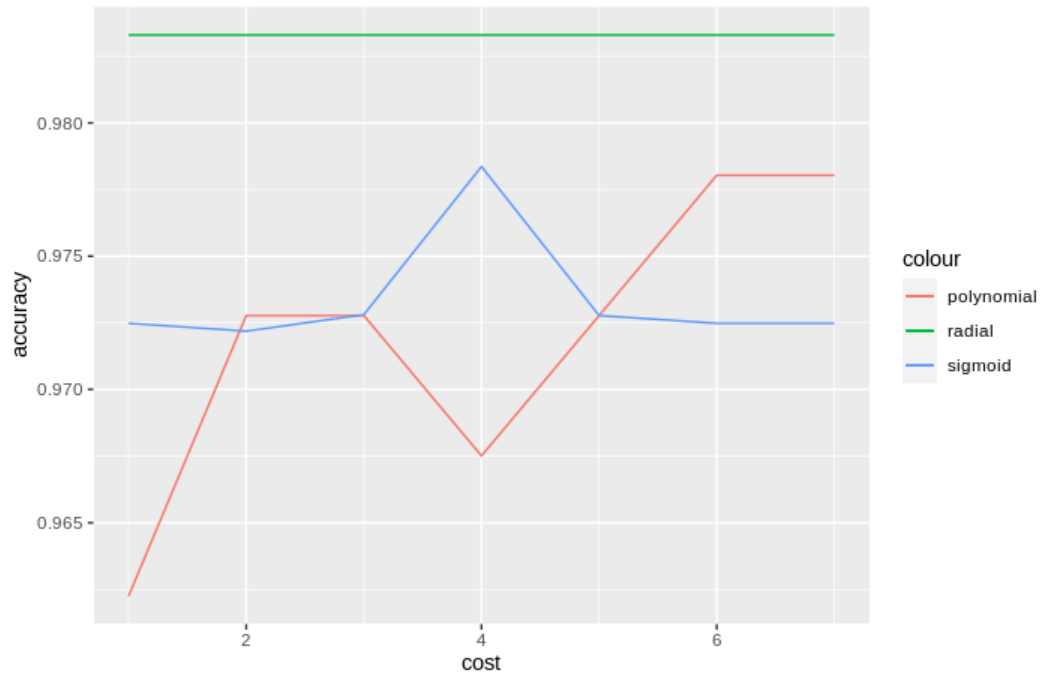


Рис. 2: График зависимости точности модели от типа ядра и параметра C

2.2 Метод k ближайших соседей

На рис. 3 и 4 представлены графики зависимости ошибки классификации от k и типа ядра при параметре distance = 1 и distance = 2 соответственно.

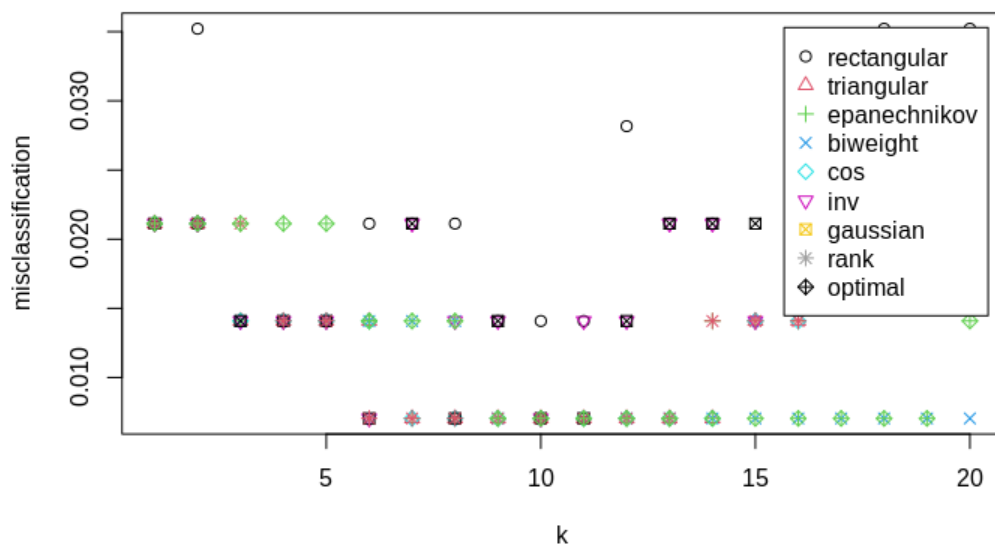


Рис. 3: distance = 1

Минимальная ошибка классификации: 0.007042254, лучшее ядро: triangular, лучшее k: 7.

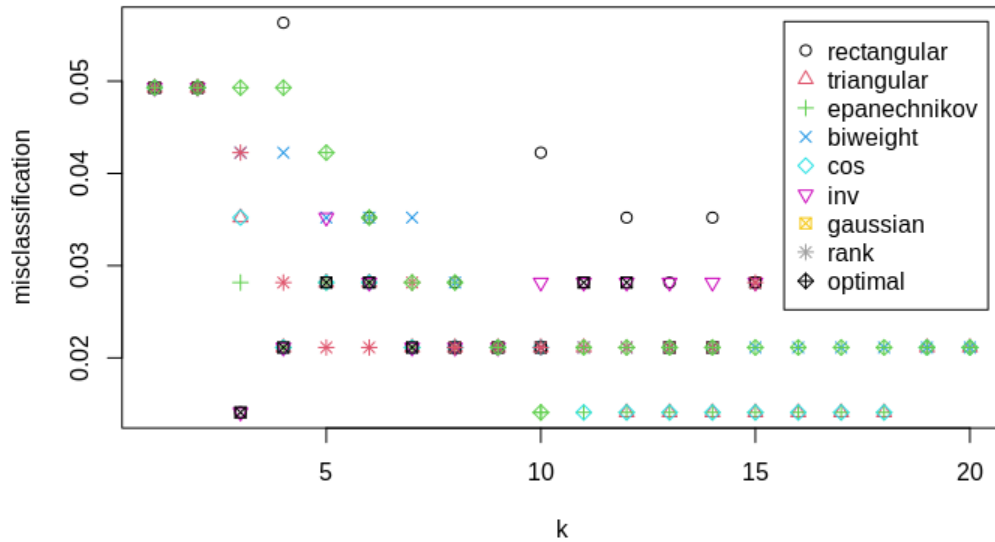


Рис. 4: distance = 2

Минимальная ошибка классификации: 0.01408451, лучшее ядро: rectangular, лучшее k: 3.

Модель с ядром triangular, k=7, distance=1 классифицировала тестовые данные с 97% точностью.

2.3 Деревья решений

На рис. 3 представлено одно из деревьев, полученных при обучении модели на обучающих данных. Средняя точность классификации, полученная при кросс-валидации на 10 разбиениях – 89%

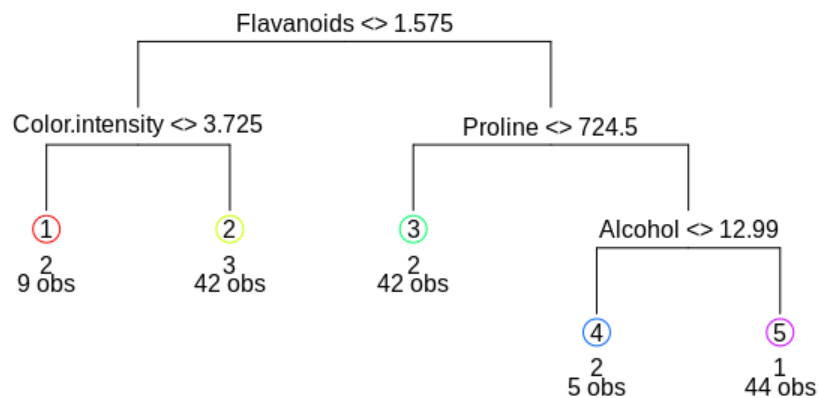


Рис. 5: Дерево решений

2.4 Результат

Были найдены оптимальные параметры для трёх моделей классификации: метод опорных векторов, метод k ближайших соседей и деревья решений. Точность классификации на тестовых данных составила:

- для метода опорных векторов – 98%;
- для k ближайших соседей – 97%
- для деревьев решений – 89%;

Таким образом, метод опорных является наиболее оптимальным.

3 Кластеризация

Кластеризация была проведена иерархическим методом. На рис. 5 представлена построенная дендограмма. Доля ошибочно кластеризованных данных – 0.522.

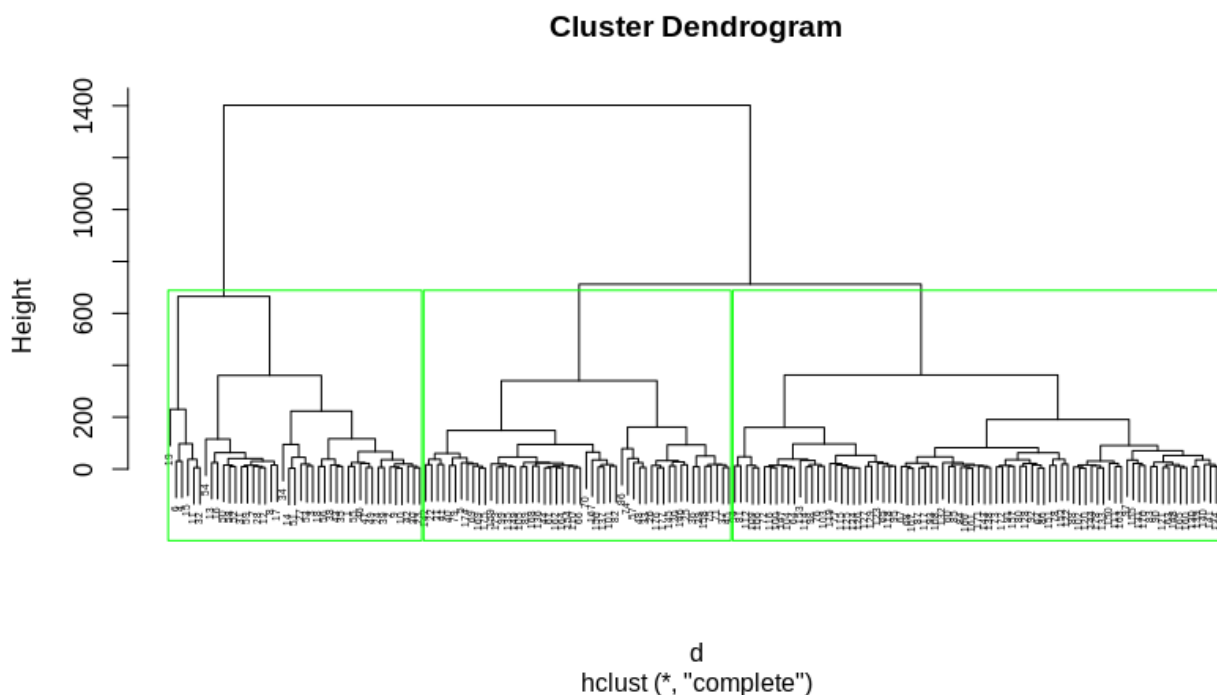


Рис. 6: Дендограмма

4 Определение наиболее значимых признаков

На рис. 6 – значения коэффициентов для всех признаков.

Наименее значимыми являются признаки Magnesium, Proline, Proanthocyanins. Наиболее значимыми – Flavanoids, Wines, Hue.

(Intercept)	4.245141633
Alcohol	-0.092500138
Malic.acid	0.018843655
Ash	-0.109275403
Alcalinity.of.ash	0.035653531
Magnesium	.
Total.phenols	0.018053043
Flavanoids	-0.293313197
Nonflavanoid.phenols	-0.062598400
Proanthocyanins	0.008081763
Color.intensity	0.067682052
Hue	-0.209990663
Wines	-0.247457272
Proline	-0.000691265

Рис. 7: Дендограмма

Заключение

В ходе курсовой работы было проделано следующее:

1. Выполнена визуализация данных при помощи метода t-SNE.
2. Разработано 3 классификатора (метод опорных векторов, метод k ближайших соседей и деревья решений), каждый из которых довольно точно классифицировал данные из тестовой выборки. Самым точным оказался метод опорных векторов, дающий результат, равный 98%.
3. Осуществлена кластеризацию данных, с довольно высокой долей ошибочно кластеризованных данных $\approx 52\%$.
4. Методом Лассо определены наиболее и наименее значимые признаки. Наименее значимые признаки: Magnesium, Proline, Proanthocyanins. Наиболее значимые: Flavanoids, Wines, Hue.

Список используемых источников

- [1] MachineLearning.ru Профессиональный информационно-аналитический ресурс, посвященный машинному обучению, распознаванию образов и интеллектуальному анализу данных
[Электронный ресурс]: <http://www.machinelearning.ru/wiki/index.php> (Дата обращения: 01.05.21)
- [2] Документация по языку R.

Приложение

```
1 library(e1071)
2 library(ggplot2)
3 library(kknn)
4 library(Rtsne)
5 library(tree)
6 library(maptree)
7 library(cluster)
8 library(factoextra)
9 library(dendextend)
10 library(glmnet)
11 library(pvclust)
12
13 setwd("/home/olga/MyProjects/Polikek/ML/Kurs")
14 data = read.delim("wine.data", sep = ",")
15
16 d = dim(data)[1]
17 target = data$class
18 data$class = as.factor(target)
19 #####
20 #t-sne
21 #####
22 data = subset(data, select = -c(class))
23 tsne = Rtsne(data, dims=1)
24 plot(tsne$Y, t='n')
25 text(tsne$Y, labels=target, col = target)
26
27 names(target) = seq(1:d)
28 data$class = as.factor(target)
29 #####
30 #KNN
31 #####
32 set.seed(1234)
33 s = sample(d, size = 0.8 * d)
34 train = data[s,]
35 test = data[-s,]
36
37 kernels <- c("rectangular", "triangular", "epanechnikov", "biweight", "cos", "inv", "gaussian", "rank", "optimal")
38 result <- train.kknn(class~, train, kmax=20, distance = 2,
39                     kernel=kernels)
40 plot(result)
41 knn = kknn(class~, train, test, k=3, kernel="optimal", distance = 2)
42 table_knn = table(test$class, knn$fitted.values)
43 table_knn
44 accuracy_knn = sum(diag(table_knn))/sum(table_knn)
45 accuracy_knn
46 prediction = predict(result, train[, -1])
47 table(prediction, as.numeric(train$class))
48
49 #####
50 #SVM
51 #####
52 s = sample(d)
53 data = data[s,]
54 m = round(d/10)
55 accuracy_svm = 0
56 kernel = "polynomial"
57 cost = 100
58 for (i in 1:10){
59   if (i == 1){
60     test = data[1:m,]
61     train = data[(m+1):d,]
62     test = na.omit(test)
63     svm = svm(class~, data=train)
```

```

64     prediction = predict(svm, test)
65     table_svm = table(test$class, prediction)
66     accuracy_svm = accuracy_svm + sum(diag(table_svm))/sum(table_svm)
67 }
68 else{
69     if (i == 10){
70         test = data[(m*10):d,]
71         train = data[1:(m*10-1),]
72         test = na.omit(test)
73         svm = svm(class~., data=train)
74         prediction = predict(svm, test)
75         table_svm = table(test$class, prediction)
76         accuracy_svm = accuracy_svm + sum(diag(table_svm))/sum(table_svm)
77     }
78     else{
79         test = data[(m*i):(m*i + m),]
80         train = rbind(data[1:(m*i - 1),], data[(m*i + m + 1):d,])
81         test = na.omit(test)
82         svm = svm(class~., data=train)
83         prediction = predict(svm, test)
84         table_svm = table(test$class, prediction)
85         accuracy_svm = accuracy_svm + sum(diag(table_svm))/sum(table_svm)
86     }
87 }
88 }
89 accuracy_svm = accuracy_svm/10
90 accuracy_svm
91 #####
92 #tree
93 #####
94
95 accuracy_tree = 0
96 for (i in 1:10){
97     if (i == 1){
98         test = data[1:m,]
99         train = data[(m+1):d,]
100        test = na.omit(test)
101        tree = tree(class~., train)
102        prediction = predict(tree, test)
103        error = 0
104        for (i in 1:length(test$class)){
105            if (prediction[i, test$class[i]] != 1){
106                error = error + 1
107            }
108        }
109        accuracy_tree = accuracy_tree + 1 - error/dim(test)[1]
110    }
111    else{
112        if (i == 10){
113            test = data[(m*10):d,]
114            train = data[1:(m*10-1),]
115            test = na.omit(test)
116            tree = tree(class~., train)
117            prediction = predict(tree, test)
118            error = 0
119            for (i in 1:length(test$class)){
120                if (prediction[i, test$class[i]] != 1){
121                    error = error + 1
122                }
123            }
124            accuracy_tree = accuracy_tree + 1 - error/dim(test)[1]
125        }
126        else{
127            test = data[(m*i):(m*i + m),]
128            train = rbind(data[1:(m*i - 1),], data[(m*i + m + 1):d,])
129            test = na.omit(test)

```

```

130     tree = tree(class~., train)
131     prediction = predict(tree, test)
132     error = 0
133     for (i in 1:length(test$class)){
134         if (prediction[i, test$class[i]] != 1){
135             error = error + 1
136         }
137     }
138     accuracy_tree = accuracy_tree + 1 - error/dim(test)[1]
139 }
140 }
141 }
142 accuracy_tree/10
143 draw.tree(tree)
144
145 tree = tree(class~., train)
146 prediction = predict(tree, test)
147 error = 0
148 for (i in 1:length(test$class)){
149     if (prediction[i, test$class[i]] != 1){
150         error = error + 1
151     }
152 }
153 1 - error/dim(test)[1]
154 #####
155 #cluster
156 #####
157 data = read.delim("wine.data", sep = ",")
158 target = data$class
159 data = subset(data, select = -c(class))
160 d = dist(data)
161 dend = hclust(d)
162 plot(dend, cex=0.4)
163 groups = cutree(dend, k=3)
164 rect.hclust(dend, k = 3, border = "green")
165 t = table(groups, target)
166 1 -(sum(diag(t)) / sum(t))
167
168 #####
169 #lasso
170 #####
171 x = as.matrix(data[, -1])
172 y = as.matrix(data[, 1])
173 lambda_seq <- 10^seq(2, -2, by = -.1)
174 cv_output <- cv.glmnet(x, y, alpha = 1,
175                        lambda = lambda_seq, nfolds = 5)
176 plot(cv_output)
177 best_lam <- cv_output$lambda.min
178 best_lam
179 lasso_best <- glmnet(x, y, alpha = 1, lambda = best_lam, label=TRUE)
180 coef(lasso_best)

```
