

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

Отчёт
по лабораторной работе №1 «Метод ближайших соседей»
по дисциплине «Системы искусственного интеллекта»

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

Преподаватель _____ Л. В. Уткин

Содержание

1	Постановка задачи	3
1.1	Задание 1	3
1.2	Задание 2	3
1.3	Задание 3	3
1.4	Задание 4	3
2	Реализация	4
2.1	Задание 1	4
2.1.1	Исследование влияния объёма обучающей выборки на результат	4
2.1.2	Исследование влияния объёма тестовой выборки на результат	5
2.2	Задание 2	6
2.2.1	Исследование зависимости ошибки классификации от значения k и от типа ядра	6
2.2.2	Исследование зависимости ошибки классификации от метрики расстояния	6
2.2.3	Определение класса примера	7
2.3	Задание 3	8
2.4	Задание 4	9
3	Приложение	12

1 Постановка задачи

1.1 Задание 1

Исследовать, как объем обучающей выборки и количество тестовых данных, влияет на точность классификации или на вероятность ошибочной классификации в примере крестики-нолики и примере о спаме e-mail сообщений.

1.2 Задание 2

- Построить классификатор для обучающего множества Glass, данные которого характеризуются 10-ю признаками:

1. Id number: 1 to 214;
2. RI: показатель преломления;
3. Na: сода (процент содержания в соответствующем оксиде); 4. Mg;
5. Al;
6. Si;
7. K;
8. Ca;
9. Ba;
10. Fe.

Классы характеризуют тип стекла:

- (1) окна зданий, плавильная обработка
- (2) окна зданий, не плавильная обработка
- (3) автомобильные окна, плавильная обработка
- (4) автомобильные окна, не плавильная обработка (нет в базе)
- (5) контейнеры
- (6) посуда
- (7) фары

- Построить графики зависимости ошибки классификации от значения k и от типа ядра.
- Исследовать, как тип метрики расстояния (параметр distance) влияет на точность классификации.
- Определить, к какому типу стекла относится экземпляр с характеристиками $RI = 1.516$ $Na = 11.7$ $Mg = 1.01$ $Al = 1.19$ $Si = 72.59$ $K = 0.43$ $Ca = 11.44$ $Ba = 0.02$ $Fe = 0.1$
- Определить, какой из признаков оказывает наименьшее влияние на определение класса путем последовательного исключения каждого признака.

1.3 Задание 3

- Для построения классификатора использовать заранее сгенерированные обучающие и тестовые выборки, хранящиеся в файлах svmdata4.txt, svmdata4test.txt.
- Найти оптимальное значение k , обеспечивающее наименьшую ошибку классификации.
- Посмотреть, как выглядят данные на графике.

1.4 Задание 4

Разработать классификатор на основе метода ближайших соседей для данных Титаник (Titanic dataset) - <https://www.kaggle.com/c/titanic>

2 Реализация

2.1 Задание 1

2.1.1 Исследование влияния объёма обучающей выборки на результат

Для исследования влияния объёма обучающей выборки на точность классификации обучение производилось в несколько итераций, где в каждой итерации размер обучающей выборки менялся с 10% до 80% от размера всего датасета. Для получения усреднённого результата на каждой итерации было произведено 20 экспериментов, заключающихся в

1. формировании обучающей выборки, размер которой задан в начале итерации,
2. формировании тестовой выборки размера 100,
3. запуске процедуры обучения методом k ближайших соседей,
4. расчёте точности классификации, как отношения количества верных предсказаний к сумме всех предсказаний.

На рис. 1 приведён график зависимости точности предсказания от размера обучающей выборки для данных крестики-нолики.

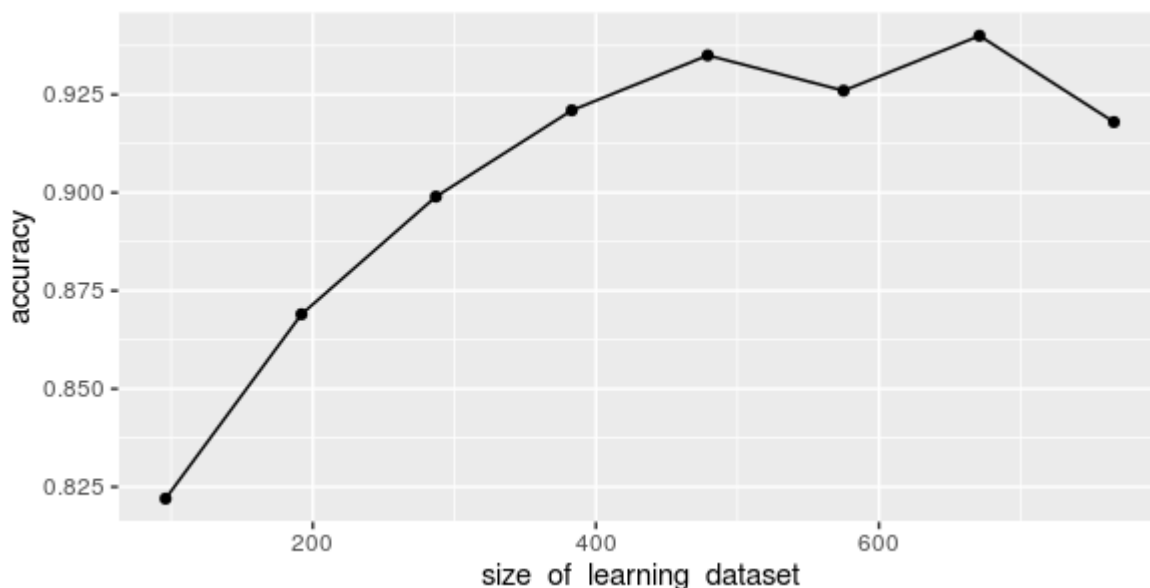


Рис. 1: Зависимость точности предсказания от размера обучающей выборки для данных крестики-нолики

По графику видно, что увеличение размера обучающей выборки до 500 даёт увеличение точности предсказания, дальнейшее увеличение может как улучшить, так и ухудшить результат, то есть достигнув достаточно больших размеров обучающей выборки не имеет смысла увеличивать его ещё больше. Также стоит заметить, что точность достигает практически 100% при больших размерах обучающей выборки, что скорее всего говорит о переобучении.

Рис. 2 представляет график зависимости точности предсказания от размера обучающей выборки для данных спам. Видно, что остались справедливыми выводы, сделанные для предыдущего примера.

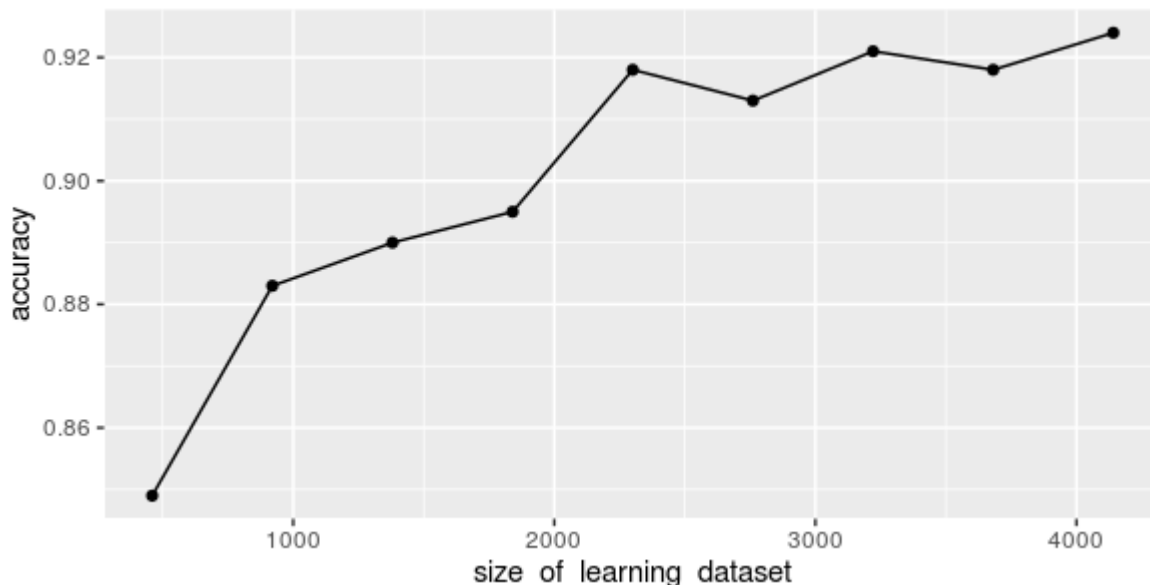


Рис. 2: Зависимость точности предсказания от размера обучающей выборки для данных спам

2.1.2 Исследование влияния объёма тестовой выборки на результат

Процедура проведения исследования влияния объёма тестовой выборки на результат аналогична описанной в предыдущем пункте, но в этом случае размер обучающей выборки равен 300, а размер тестовой выборки варьировался от 100 до 600. Результаты для примера крестики-нолики представлены на рис. 3.

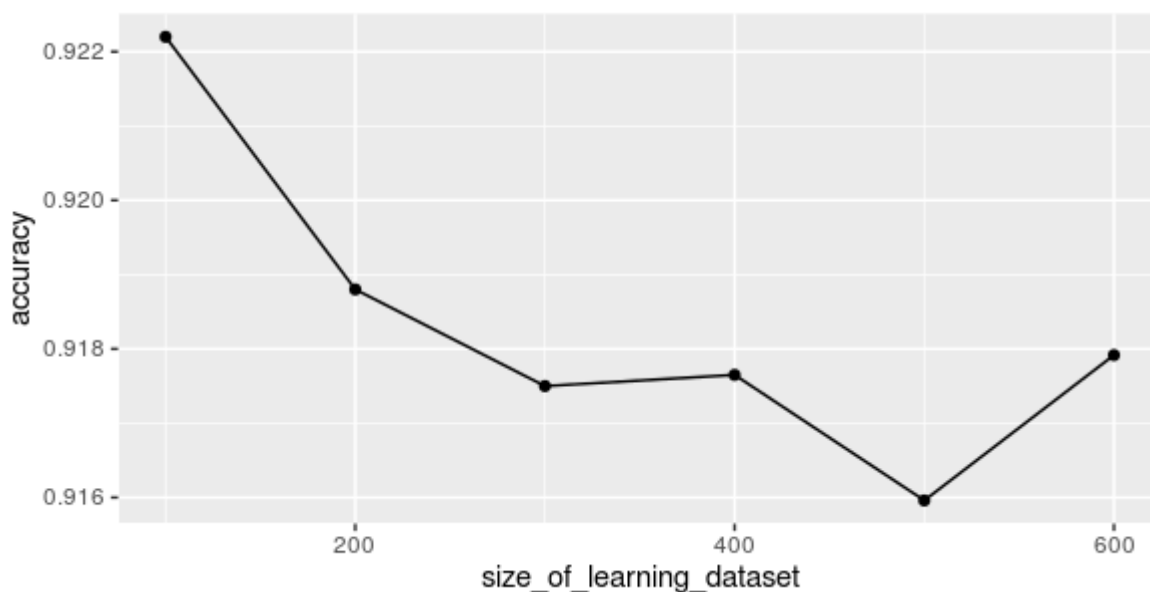


Рис. 3: Зависимость точности предсказания от размера тестовой выборки для данных крестики-нолики

Видно, что с увеличением размера тестовой выборки ухудшается точность результата. Это может быть связано с характеристикой данных, попавших в тестовую выборку: при её малом размере вероятность того, что попадутся примеры с выбросами или какими-то нетипичными параметрами, ниже. Такой случайностью выбора примеров в выборку можно объяснить и рост точности в конце при размере равном 600. Также стоит отметить небольшую величину разброса значений: разница между максимальной и минимальной точностью составила примерно 0.6%.

На рис. 4 представлен график зависимости точности предсказания от размера тестовой выборки для данных спам, где размер обучающей выборки составлял 2500.

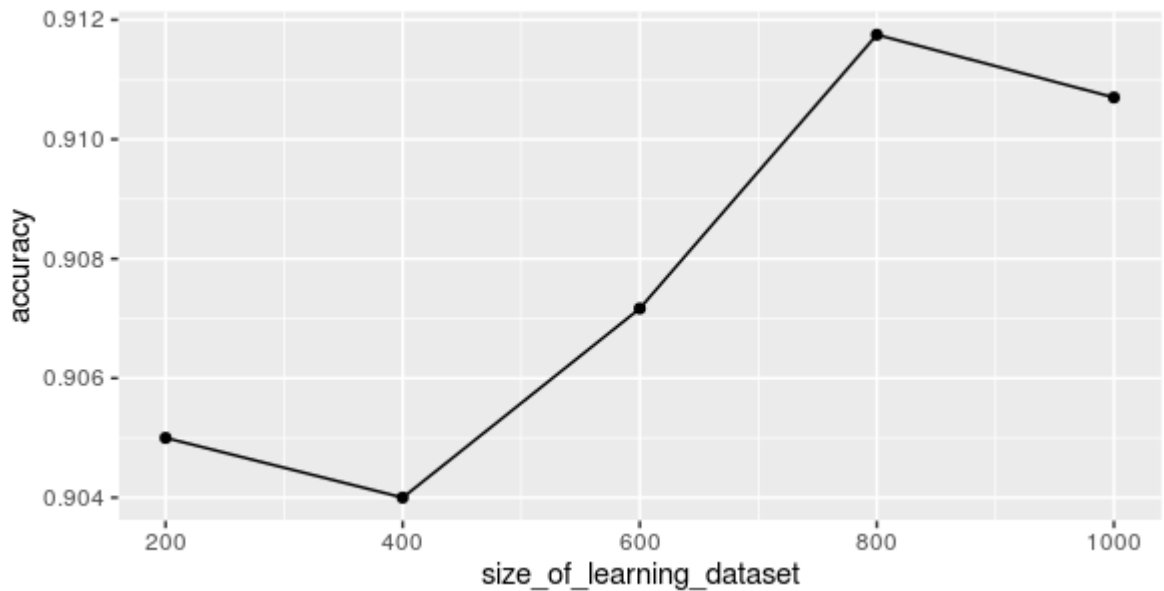


Рис. 4: Зависимость точности предсказания от размера тестовой выборки для данных спам

Здесь прослеживается обратная зависимость: чем больше размер тестовой выборки, тем больше точность. Это может объясняться тем, что в данных меньшего размера оказалось слишком много выбросов, а при увеличении выборки нормальные данные уменьшили влияние этих выбросов.

2.2 Задание 2

2.2.1 Исследование зависимости ошибки классификации от значения k и от типа ядра

На рис. 5 представлен график зависимости ошибки классификации от k и типа ядра. Лучшими параметрами получившегося классификатора, при которых ошибка равна 0.2616822, являются: $k = 8$, $\text{kernel} = \text{biweight}$.

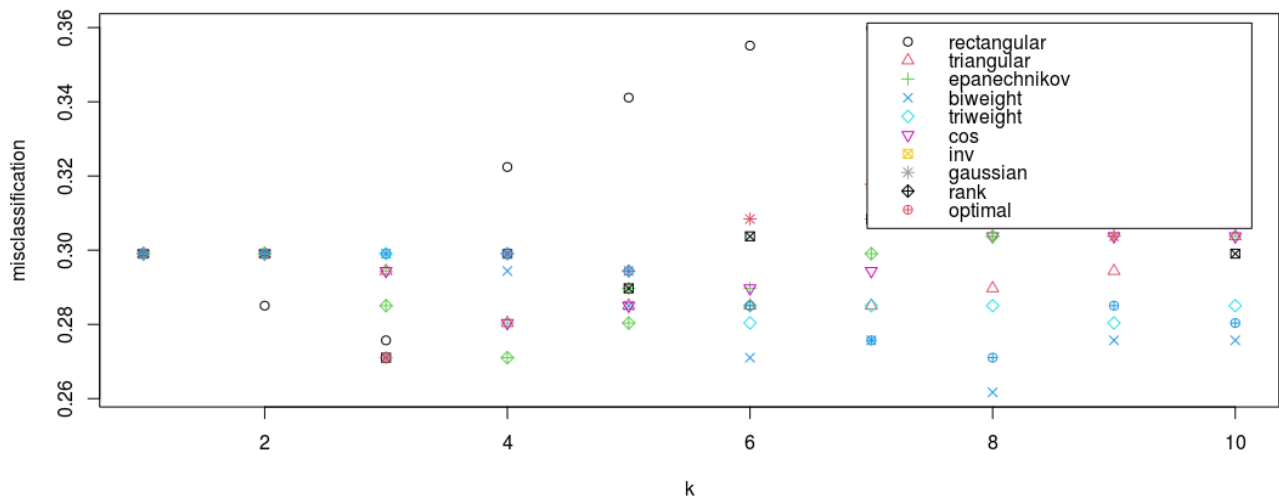


Рис. 5: Зависимость ошибки классификации от k и типа ядра

2.2.2 Исследование зависимости ошибки классификации от метрики расстояния

На рис. 6 представлен график зависимости ошибки классификации от метрики расстояния при параметрах $k = 8$ и $\text{kernel} = \text{biweight}$. Лучшим значением метрики является 1.

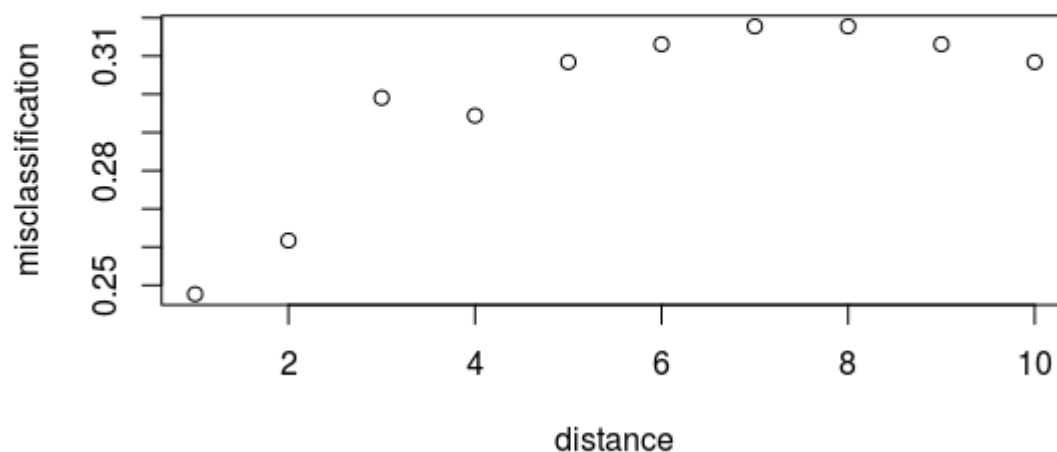


Рис. 6: Зависимость ошибки классификации от метрики расстояния

2.2.3 Определение класса примера

Экземпляр с данными характеристиками принадлежит к типу 5. При исключении одного из признаков получались следующие результаты классификации, выведенные с вероятностями определений данных классов:

- "RI"
5
0.9250552
- "Na"
5
0.5257786
- "Mg"
5
0.4049444
- "Al"
5
0.8371996
- "Si"
5
0.8727281
- "K"
5
0.7279954
- "Ca"
2
0.3762198
- "Ba"
5
0.7279954

- "Fe"
- 5
0.75672

Видно, что неправильно определяется класс при исключении признака Ca, что означает, что он сильно влияет на результат классификации. Также стоит обратить внимание на относительно низкую вероятность правильной классификации при удалении признака Mg, возможно он тоже влияет на результат, пусть и в меньшей степени.

2.3 Задание 3

Распределение данных на графике представлено на рис. 7.

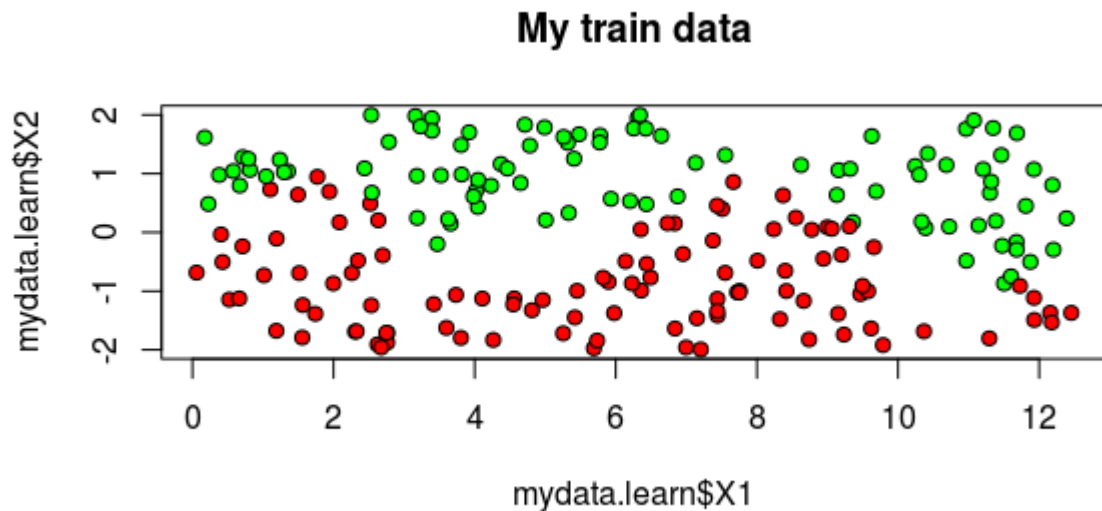


Рис. 7: Распределение данных

Обучение метода KNN методом LOO дало следующие параметры классификатора:

Minimal misclassification: 0.035

Best kernel: optimal

Best k: 8

На рис. 8 представлен график зависимости величины ошибки классификации от k.

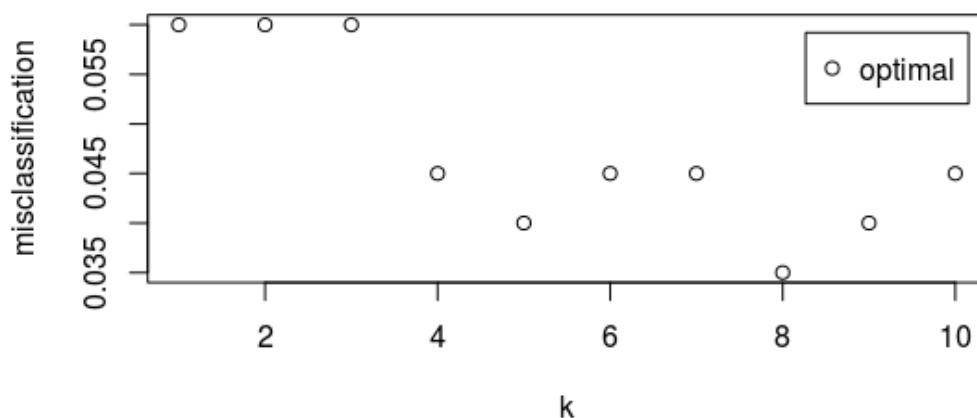


Рис. 8: Ошибка классификации при различных k

2.4 Задание 4

Данные содержат следующую информацию о пассажирах:

- Идентификатор,
- Спасся человек или нет,
- Класс билета,
- Имя,
- Пол,
- Возраст,
- Количество родственников на борту 2-го порядка,
- Количество родственников на борту 1-го порядка,
- Номер билета,
- Цена билета,
- Каюта,
- Порт посадки.

Логично предположить, что на вероятность быть спасённым не влияют такие факторы, как идентификатор, имя, номер билета, каюта и порт посадки. Тогда столбцы PassengerId, Name, Ticket, Cabin и Embarked можно удалить (в ходе дальнейших экспериментов было выявлено, что поле Fare также не улучшает точность предсказания модели). Влияние факторов класс, пол и наличие родственников на выживаемость наглядно продемонстрированы на рис. 9-12.

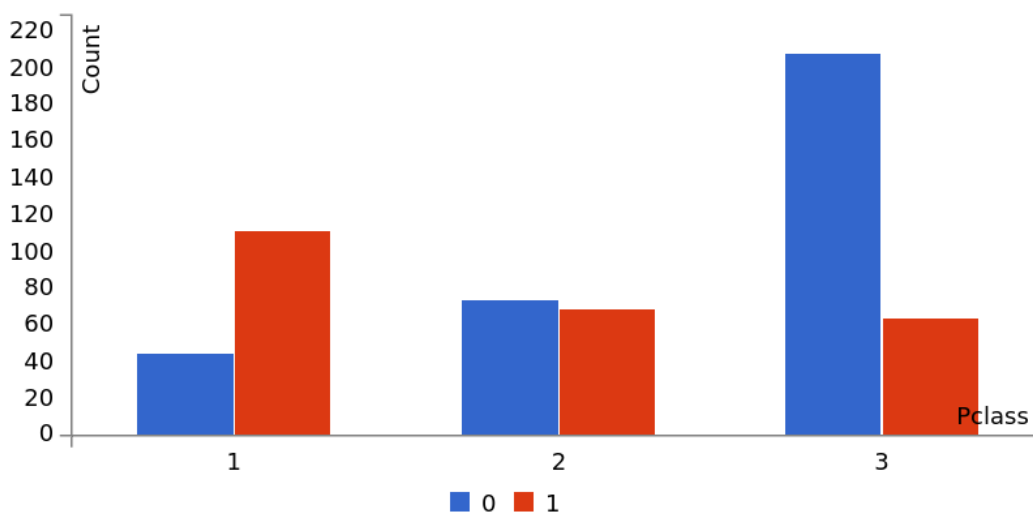


Рис. 9: График зависимости количества выживших от класса

Видно, что доля выживших преобладает только в категории пассажиров первого класса, а большая часть погибших были среди пассажиров третьего класса.

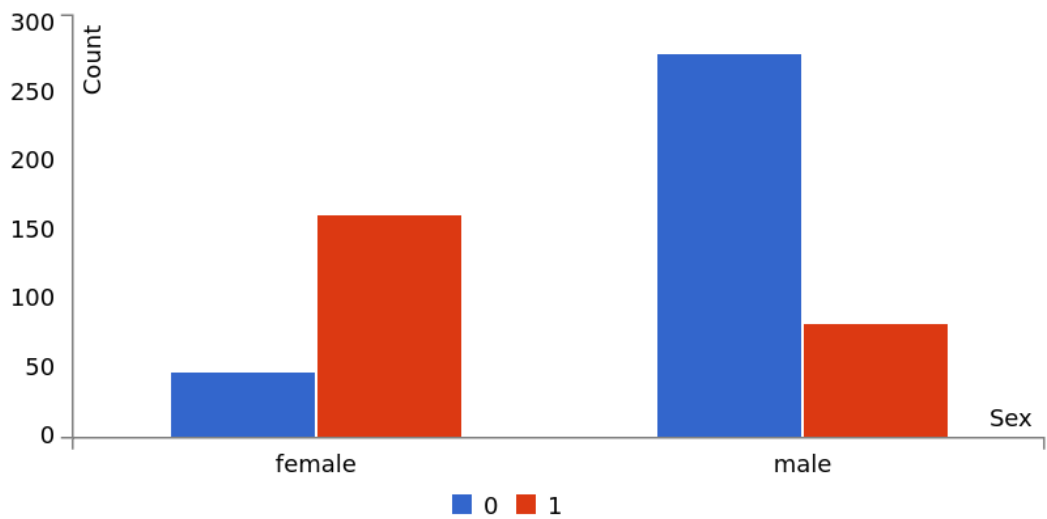


Рис. 10: График зависимости количества выживших от пола

Пол имеет большое значение для определения того, выживет пассажир или нет: большая часть женщин была спасена, а большая часть мужчин – нет.

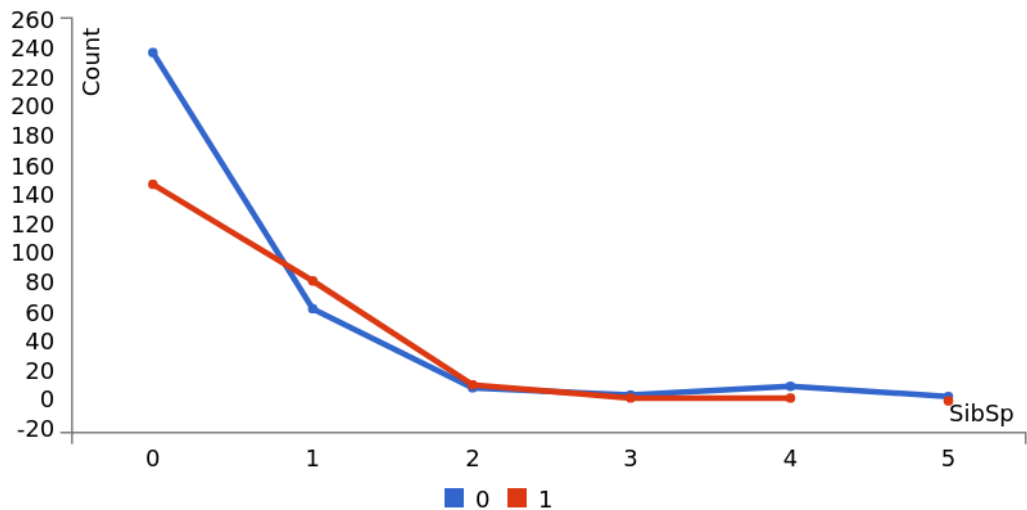


Рис. 11: График зависимости количества выживших от количества родственников 2-го порядка

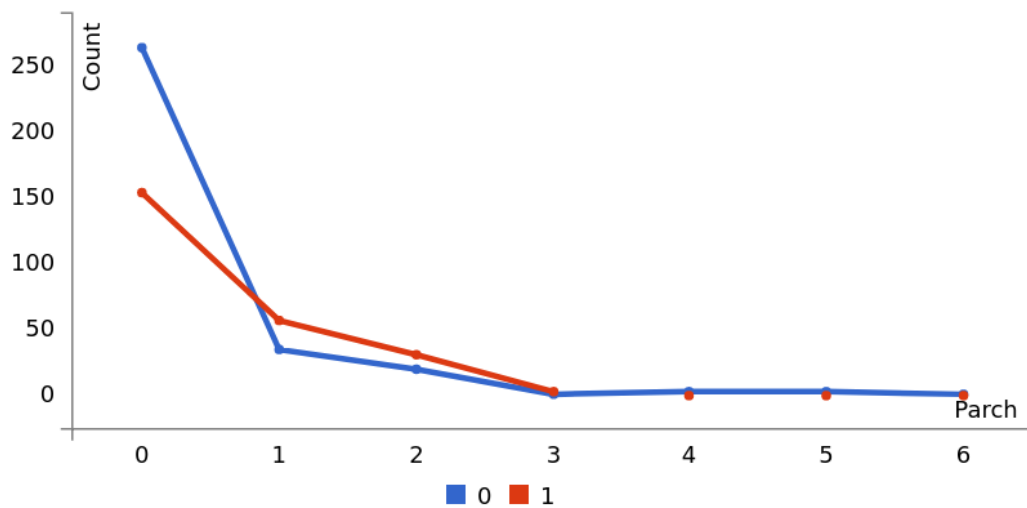


Рис. 12: График зависимости количества выживших от количества родственников 1-го порядка

Среди тех, кто имел большое количество родственников на борту, никто не выжил.

Данные содержат пропущенные значения (NA). В обучающих данных их всего 117, все они содержатся в поле Age. В данных для тестирования такие значения присутствуют в полях Age (86) и Fare (1). Записи с NA были удалены, т.к. попытка присвоить пропущенным значениям медианные значения всей выборки дала в итоге меньшую точность ответа.

В данных для тестирования отсутствовали метки Survived, поэтому была сформирована валидационная выборка для оценки точности предсказаний модели, составляющая 20% от обучающей выборки.

В итоге обучения на тестовых данных модель предсказала 195 погибших и 137 выживших пассажира. Точность на валидационной выборке составила 81%.

3 Приложение

Задание 1

```
1 library(kknn)
2 library(ggplot2)
3 library(dplyr)
4 setwd("/home/olga/MyProjects/Polikek/ML/KNN/datasets")
5
6 data <- read.delim("Tic_tac_toe.txt", sep=" ", head=FALSE)
7 data$V10 <- as.factor(data$V10)
8 m <- dim(data)[1]
9 mean_accuracy <- c()
10 #size_list <- seq(0.1, 0.8, by=0.1) * m
11 size_list <- c(100, 200, 300, 400, 500)
12 for (j in 1:length(size_list)){
13   sample_size <- size_list[j]
14   accuracy <- c()
15   for (k in 1:20){
16     val <- sample(1:m, size = 400)
17     data.learn <- data[val,]
18     data2 <- data[-val,]
19     val2 <- sample(m-400, size = sample_size)
20     data.valid <- data2[val2,]
21     data.kknn <- kknn(V10~., data.learn, data.valid)
22     accuracy_table <- table(data.kknn$fitted.values, select(data.valid, V10)[,1])
23     accuracy[k] <- sum(diag(accuracy_table)) / sum(accuracy_table)
24   }
25   mean_accuracy[j] <- mean(accuracy)
26 }
27
28
29 new_data <- data.frame(
30   size_of_learning_dataset=size_list,
31   accuracy=mean_accuracy
32 )
33 ggplot(new_data, aes(x=size_of_learning_dataset, y=accuracy)) +
34   geom_line() +
35   geom_point()
36
37
38 library(kernlab)
39 library(kknn)
40
41 data(spam)
42 m <- dim(spam)[1]
43 mean_accuracy <- c()
44 #size_list <- seq(0.1, 0.9, by=0.1) * m
45 size_list <- c(200, 400, 500, 600, 700, 800, 900, 1000)
46 for (j in 1:length(size_list)){
47   sample_size <- size_list[j]
48   accuracy <- c()
49   for (k in 1:100){
50     val <- sample(1:m, size = 2500)
51     spam.learn <- spam[val,]
52     data2 <- spam[-val,]
53     val2 <- sample(m-2500, size = sample_size)
54     spam.valid <- data2[val2,]
55     spam.kknn <- kknn(type~., spam.learn, spam.valid)
56     accuracy_table <- table(spam.kknn$fitted.values, select(spam.valid, type)[,1])
57     accuracy[k] <- sum(diag(accuracy_table)) / sum(accuracy_table)
58   }
59   mean_accuracy[j] <- mean(accuracy)
60 }#600
61
62 library(ggplot2)
```

```

63 new_data <- data.frame(
64   size_of_learning_dataset=size_list,
65   accuracy=mean_accuracy)
66
67 ggplot(new_data, aes(x=size_of_learning_dataset, y=accuracy)) +
68   geom_line() +
69   geom_point()

```

Задание 2

```

1  library(kknn)
2  library(mlbench)
3  library(ggplot2)
4  data(Glass)
5  kernels <- c("rectangular","triangular", "epanechnikov", "biweight",
6              "triweight", "cos", "inv", "gaussian", "rank", "optimal" )
7  result <- train.kknn(Type~., Glass, kmax=10,
8                      kernel=c("rectangular","triangular", "epanechnikov", "biweight",
9                              "triweight", "cos", "inv", "gaussian", "rank", "optimal" ))
10
11
12  plot(result)
13  misclassification <- c()
14  for (i in 1:10){
15    result <- train.kknn(Type~., Glass, distance=i, ks=c(8), kernel="biweight")
16    misclassification[i] <- result$MISCLASS
17  }
18  distance <- seq(1,10)
19  plot(x=distance, y=misclassification)
20
21  example <- list(RI =1.516, Na =11.7, Mg =1.01, Al =1.19, Si =72.59,
22                K=0.43, Ca =11.44, Ba =0.02, Fe =0.1 )
23
24  example_kknn <- kknn(Type~., Glass, example)
25  print(example_kknn$fitted.values, max.levels = 0)
26  print(example_kknn$prob[, "5"])
27  for (i in 1:length(example)){
28    print(names(example)[i])
29    new_ex <- example[-i]
30    new_Glass <- Glass[-i]
31    example_kknn <- kknn(Type~., new_Glass, new_ex)
32    #sprintf("class: %d", example_kknn$fitted.values, max.levels = 0)
33    #sprintf("probability: %f", example_kknn$prob[, "5"])
34    print(example_kknn$fitted.values, max.levels = 0)
35    print(example_kknn$prob[, example_kknn$fitted.values])
36    new_ex <- example
37    new_Glass <- Glass
38  }

```

Задание 3

```

1  library(kknn)
2  library(ggplot2)
3  setwd("/home/olga/MyProjects/Polikek/ML/KNN/datasets")
4  mydata.learn <- read.delim("svmdata4.txt", sep="\t")
5  mydata.learn$Colors <- as.factor(mydata.learn$Colors)
6  mydata.test <- read.delim("svmdata4test.txt", sep="\t")
7  mydata.test$Colors <- as.factor(mydata.test$Colors)
8  mydata.train <- train.kknn(Colors~., mydata.learn, kmax=10)
9  plot(mydata.train)
10 plot(mydata.learn$X1, mydata.learn$X2, pch=21,
11       bg=c("green", "red") [unclass(mydata.learn$Colors)], main="My train data")

```

Задание 4

```

1 library(kknn)
2 library(dplyr)
3 library(ggplot2)
4 library(rpivotTable)
5 library(dummies)
6 setwd("/home/olga/MyProjects/Polikek/ML/KNN/datasets")
7
8 data.train = read.csv("Titanic_train2.csv", sep=",", quote = "\"", stringsAsFactors = TRUE)
9 data.train = select(data.train, Survived, Pclass, Age, Sex, SibSp, Parch)
10 #data.train$Fare[data.train$Fare==0] <- NA
11 #data.train$Age[is.na(data.train$Age)] <- median(data.train$Age, na.rm = TRUE)
12 data.train = na.omit(data.train)
13 data.train$Survived = factor(data.train$Survived)
14 data.train$Pclass = factor(data.train$Pclass, ordered=TRUE, levels = c(3, 2, 1))
15 #colSums(is.na(data.train))
16
17
18 data.test = read.csv("Titanic_test2.csv", sep=",", quote = "\"", stringsAsFactors = TRUE)
19 data.test = select(data.test, Pclass, Age, Sex, SibSp, Parch)
20 data.test$Pclass = factor(data.test$Pclass, ordered=TRUE, levels = c(3, 2, 1))
21 #data.train$Age[is.na(data.train$Age)] <- median(data.train$Age, na.rm = TRUE)
22 data.test = na.omit(data.test)
23 #data.test$Age[is.na(data.test$Age)] <- mean(data.test$Age, na.rm = TRUE)
24 #data.test$Fare[is.na(data.test$Fare)] <- mean(data.test$Fare, na.rm = TRUE)
25 s = round(seq(dim(data.train)[1]*0.2))
26 data.valid = data.train[s,]
27 data.train = data.train[-s,]
28
29 data.knn = kknn(Survived~., data.train, data.test, k=8)
30 table(data.knn$fitted.values)
31 data.check = kknn(Survived~., data.train, data.valid, k=8)
32 t = table(data.check$fitted.values, data.valid$Survived)
33 accuracy = sum(diag(t)) / sum(t)
34
35
36 #data.train$Sex = factor(data.train$Sex)
37 #data.train$SibSp = factor(data.train$SibSp)
38 #data.train$Parch = factor(data.train$Parch)
39 rpivotTable(data.train, rows="Survived", cols="Sex")
40 rpivotTable(data.train, rows="Survived", cols="SibSp")
41 rpivotTable(data.train, rows="Survived", cols="Pclass")
42 rpivotTable(data.train, rows="Survived", cols="Age")
43 rpivotTable(data.train, rows="Survived", cols="Parch")

```
