

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

Отчёт
по лабораторной работе №2 «Наивный Байесовский
классификатор»
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	4
2.1	Постановка задачи	4
2.2	Реализация	4
3	Задание 3	5
3.1	Постановка задачи	5
3.2	Реализация	6
	Приложение	7

1 Задание 1

1.1 Постановка задачи

Исследуйте, 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 Реализация

На рис. 1 и 2 представлены графики зависимости точности классификации от размера обучающей и тестовой выборки для данных крестики-нолики. При увеличении размера обучающей выборки, точность классификации на тестовых данных сначала повышается, так как новые данные позволяют модели лучше обучиться, но при достижении определённого уровня точность снижается из-за переобучения модели. При увеличении размера тестовой выборки точность также снижается, 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.

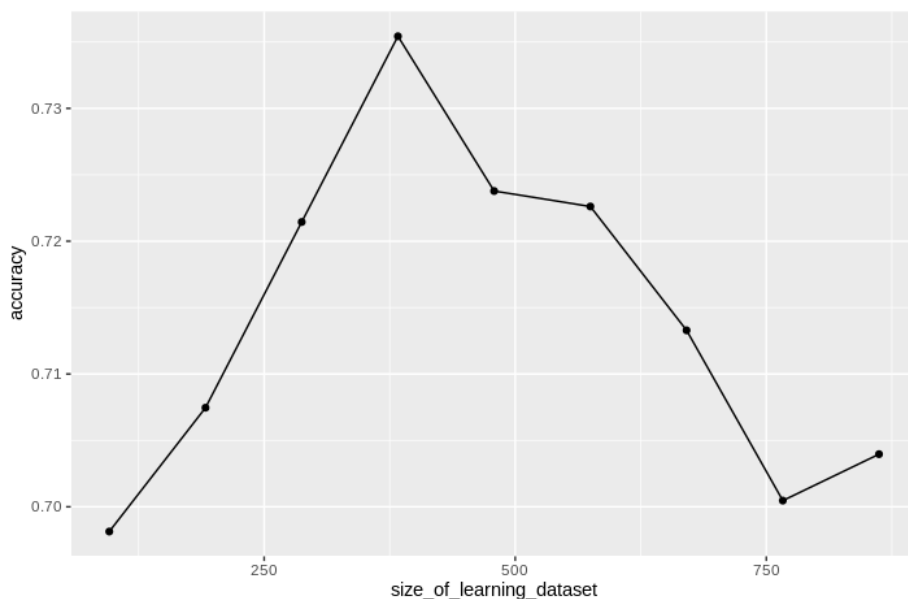


Рис. 1: Зависимость точности классификации от размера обучающей выборки для данных крестики-нолики

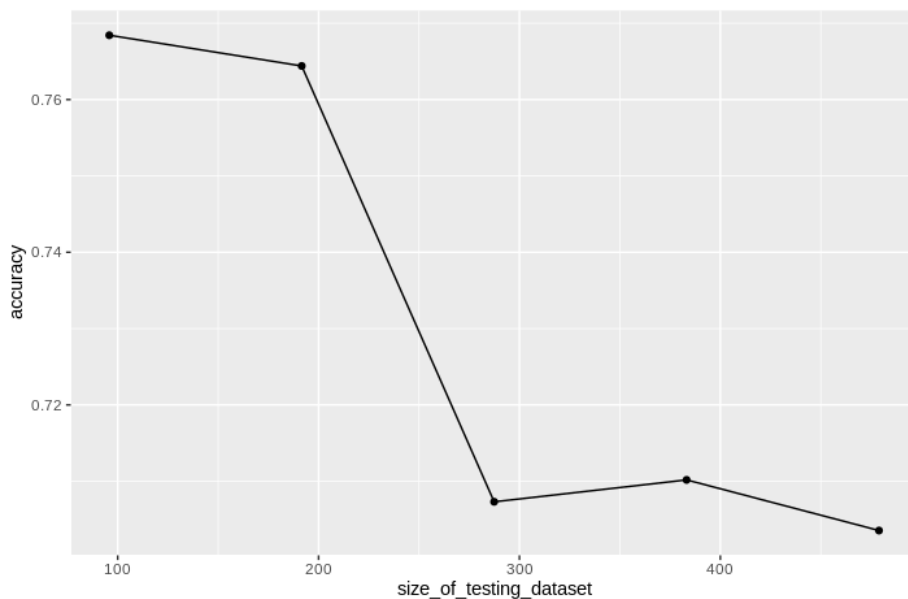


Рис. 2: Зависимость точности классификации от размера тестовой выборки для данных спам

На рис. 3 и 4 представлены графики зависимости точности классификации от размера обучающей и тестовой выборки для данных спам, подтверждающие выводы, сделанные выше.

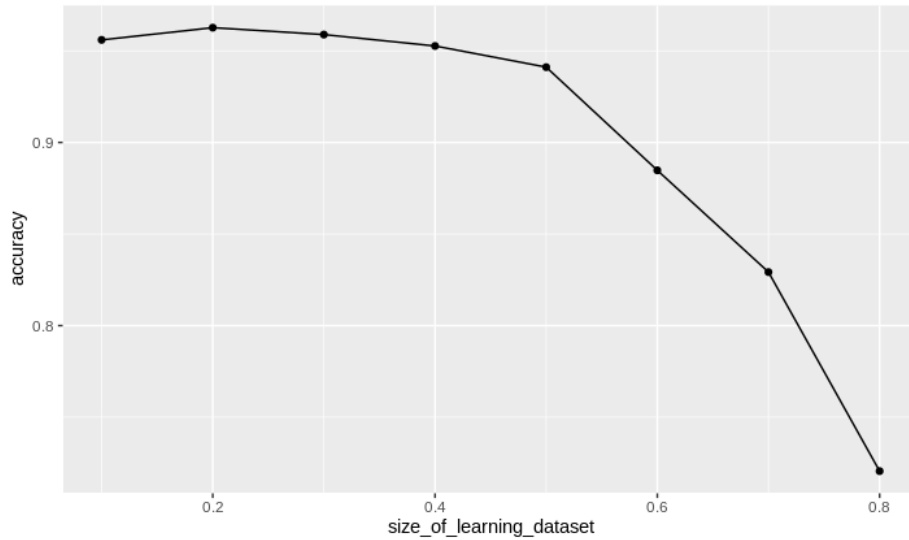


Рис. 3: Зависимость точности классификации от размера обучающей выборки для данных спам

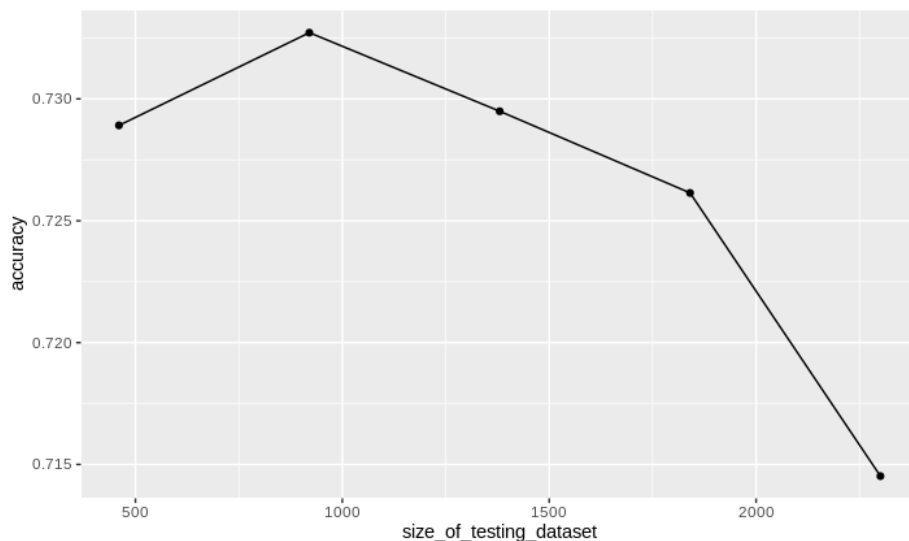


Рис. 4: Зависимость точности классификации от размера тестовой выборки для данных крестики-нолики

2 Задание 2

2.1 Постановка задачи

Сгенерируйте 100 точек с двумя признаками X_1 и X_2 в соответствии с нормальным распределением так, что первые 50 точек (class -1) имеют параметры: мат. ожидание X_1 равно 10, мат. ожидание X_2 равно 14, среднеквадратические отклонения для обеих переменных равны 4. Вторые 50 точек (class +1) имеют параметры: мат. ожидание X_1 равно 20, мат. ожидание X_2 равно 18, среднеквадратические отклонения для обеих переменных равны 3. Построить соответствующие диаграммы, иллюстрирующие данные. Построить байесовский классификатор и оценить качество классификации.

2.2 Реализация

Диаграмма рассеяния точек представлена на рис. 5. Рис. 6 и 7 показывают гистограммы распределения. Построенный байесовский классификатор предсказывает классы точек с точностью 95%.

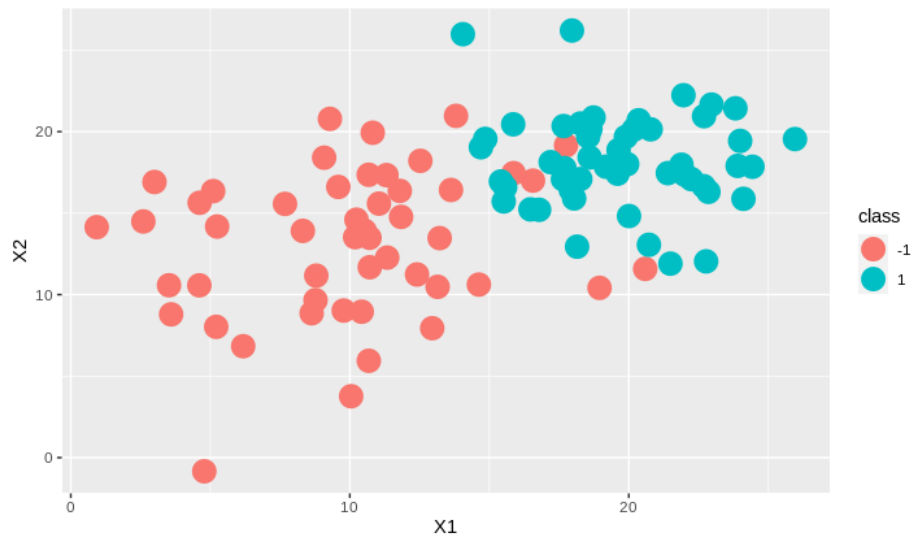


Рис. 5: Диаграмма рассеяния сгенерированных данных

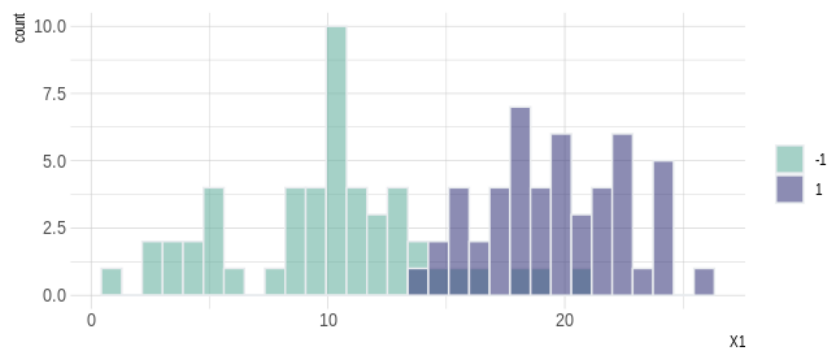


Рис. 6: Гистограмма распределения X1

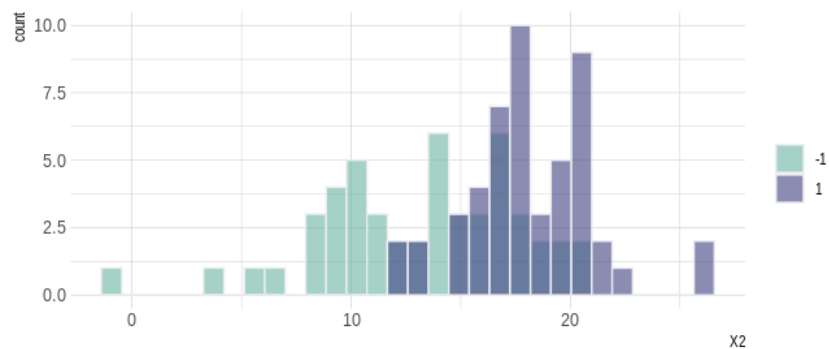


Рис. 7: Гистограмма распределения X2

3 Задание 3

3.1 Постановка задачи

Разработать байесовский классификатор для данных Титаник (Titanic dataset) - <https://www.kaggle.com/c/titanic>
 Исходные обучающие данные для классификации – в файле Titanic_train.csv Данные для тестирования – в файле Titanic_test.csv

3.2 Реализация

Были удалены, как не влияющие на выживаемость пассажира признаки столбцы PassengerId, Name, Ticket, Cabin и Embarked.

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.

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 160 погибших и 172 выживших пассажира. Точность на валидационной выборке составила 72%.

Приложение

Задание 1

```
1 library(e1071)
2 library(ggplot2)
3
4 setwd("/home/olga/MyProjects/Polikek/ML/Bayes/datasets")
5 A_raw <- read.table("Tic_tac_toe.txt", sep = ",", stringsAsFactors = TRUE)
6 m <- dim(A_raw)[1]
7 size <- seq(0.1, 0.5, by=0.1)
8 set.seed(12345)
9 mean_accuracy <- c()
10 for (j in 1:length(size)){
11   A_rand <- A_raw[ order(runif(m)), ]
12   nt <- as.integer(m*size[j])
13   accuracy <- c()
14   for (k in 1:10){
15     A_test <- A_rand[1:nt, ]
16     #A_test <- A_rand[(nt+1):n, ]
17     A_train <- A_rand[m-400:m, ]
18     A_classifier <- naiveBayes(V10 ~ ., data = A_train)
19     A_predicted <- predict(A_classifier, A_test)
20     accuracy_table <- table(A_predicted, A_test$V10)
21     accuracy[k] <- sum(diag(accuracy_table)) / sum(accuracy_table)
22   }
23   mean_accuracy[j] <- mean(accuracy)
24 }
25
26 new_data <- data.frame(
27   size_of_testing_dataset=size * m,
28   accuracy=mean_accuracy)
29 ggplot(new_data, aes(x=size_of_testing_dataset, y=accuracy)) +
30   geom_line() +
31   geom_point()
32
33
34 library(e1071)
35 library(ggplot2)
36 library(kernlab)
37
38 data(spam)
39 size <- seq(0.1, 0.5, by=0.1)
40 n = dim(spam)[1]
41 mean_accuracy <- c()
42 for (j in 1:length(size)){
43   accuracy <- c()
44   s = n * size[j]
45   for (k in 1:10){
46     idx <- sample(1:dim(spam)[1], s)
47     spamtrain <- spam[-idx, ][1:2000,]
48     spamtest <- spam[idx, ]
49     model <- naiveBayes(type ~ ., data = spamtrain)
50     accuracy_table <- table(predict(model, spamtest), spamtest$type)
51     #predict(model, spamtest, type = "raw")
52     accuracy[k] <- sum(diag(accuracy_table)) / sum(accuracy_table)
53   }
54   mean_accuracy[j] <- mean(accuracy)
55 }
56
57 new_data <- data.frame(
58   size_of_testing_dataset=size * n,
59   accuracy=mean_accuracy)
60 ggplot(new_data, aes(x=size_of_testing_dataset, y=accuracy)) +
61   geom_line() +
62   geom_point()
```

Задание 2

```
1 library(ggplot2)
2 library(ggExtra)
3 library(hrbrthemes)
4 library(e1071)
5
6 data = data.frame(X1 = c(rnorm(50, mean=10, 4), rnorm(50, mean=20, 3)),
7                    X2 = c(rnorm(50, mean=14, 4), rnorm(50, mean=18, 3)),
8                    class = c(rep("1",50), rep("1",50)))
9
10 ggplot(data, aes(x=X2, fill=class)) +
11   geom_histogram( color="#e9ecef", alpha=0.6, position = 'identity') +
12   scale_fill_manual(values=c("#69b3a2", "#404080")) +
13   theme_ipsum() +
14   labs(fill="")
15
16 ggplot(data, aes(x=X1, y=X2, color=class)) +
17   geom_point(size=6)
18
19 rows = sample(nrow(data))
20 data = data[rows, ]
21 train = data[1:60,]
22 test = data[61:100,]
23 bayes = naiveBayes(class~., train)
24 prediction = predict(bayes, test)
25 ac_tab = table(prediction, test$class)
26 correct = sum(diag(ac_tab))/sum(ac_tab)
```

Задание 3

```
1 library(dplyr)
2 library(ggplot2)
3 library(rpivotTable)
4 library(dummies)
5 setwd("/home/olga/MyProjects/Polikek/ML/Bayes/datasets")
6
7 data.train = read.csv("Titanic_train2.csv", sep=",", quote = "\"", stringsAsFactors = TRUE)
8 data.train = select(data.train, Survived, Pclass, Age, Sex, SibSp, Parch)
9 data.train = na.omit(data.train)
10 data.train$Survived = factor(data.train$Survived)
11 data.train$Pclass = factor(data.train$Pclass, ordered=TRUE, levels = c(3, 2, 1))
12
13 data.test = read.csv("Titanic_test2.csv", sep=",", quote = "\"", stringsAsFactors = TRUE)
14 data.test = select(data.test, Pclass, Age, Sex, SibSp, Parch)
15 data.test$Pclass = factor(data.test$Pclass, ordered=TRUE, levels = c(3, 2, 1))
16 data.test = na.omit(data.test)
17
18 bayes = naiveBayes(Survived ~ ., data = data.train)
19 prediction = predict(bayes, data.test)
20
21 s = round(seq(dim(data.train)[1]*0.2))
22 data.valid = data.train[s,]
23 data.train = data.train[-s,]
24 bayes = naiveBayes(Survived ~ ., data = data.train)
25 prediction = predict(bayes, data.valid)
26 t = table(prediction, data.valid$Survived)
27 accuracy = sum(diag(t)) / sum(t)
```
