

Компьютерный практикум по статистическому анализу данных

Лабораторная работа № 7. Введение в работу с данными

Демидова Екатерина Алексеевна

Содержание

1	Введение	4
2	Теоретическое введение	5
3	Выполнение лабораторной работы	6
4	Выводы	20
	Список литературы	21

Список иллюстраций

3.1	Примеры	6
3.2	Примеры	7
3.3	Примеры	8
3.4	Примеры	9
3.5	Примеры	10
3.6	Примеры	11
3.7	Примеры	12
3.8	Примеры	13
3.9	Задание 1	14
3.10	Задание 1	15
3.11	Задания 2	15
3.12	Задание 2	16
3.13	Задание 3	17
3.14	Задание 3	18
3.15	Задание 3	19

1 Введение

Цель работы

Основной целью работы является освоение специализированных пакетов Julia для обработки данных.

Задачи

1. Используя Jupyter Lab, повторите примеры.
2. Выполните задания для самостоятельной работы.

2 Теоретическое введение

Julia — высокоуровневый свободный язык программирования с динамической типизацией, созданный для математических вычислений[1]. Эффективен также и для написания программ общего назначения. Синтаксис языка схож с синтаксисом других математических языков, однако имеет некоторые существенные отличия.

Для выполнения заданий была использована официальная документация Julia[2].

3 Выполнение лабораторной работы

Выполним примеры из лабораторной работы для изучения циклов и функций (рис. 3.1 - 3.8)

```
> using CSV, DataFrames, DelimitedFiles, Plots, Statistics, GLM
[198]

Считывание данных

# Считывание данных и их запись в структуру:
P = CSV.File("programminglanguages.csv") |> DataFrame

# Функция определения по названию языка программирования года его создания:
function language_created_year(P, language::String)
    loc = findfirst(P[:,2].==language)
    return P[loc,1]
end

# Пример вызова функции и определение даты создания языка Python:
language_created_year(P, "Python")
# Пример вызова функции и определение даты создания языка Julia:
language_created_year(P, "Julia")

# Функция определения по названию языка программирования
# года его создания (без учёта регистра):
function language_created_year_v2(P, language::String)
    loc = findfirst(lowercase.(P[:,2]).==lowercase.(language))
    return P[loc,1]
end

language_created_year_v2(P, "julia")

# Построчное считывание данных с указанием разделителя:
Tx = readlm("programminglanguages.csv", ',')

Запись данных в файл

# Запись данных в CSV-файл:
CSV.write("programming_languages_data2.csv", P)
Можно задать тип файла и разделитель данных:
# Пример записи данных в текстовый файл с разделителем ',':
writedlm("programming_languages_data.txt", Tx, ',')
# Пример записи данных в текстовый файл с разделителем '-':
writedlm("programming_languages_data2.txt", Tx, '-')

# Построчное считывание данных с указанием разделителя:
P_new_delim = readlm("programming_languages_data2.txt", '-')

```

Рис. 3.1: Примеры

```

Словари

# Инициализация словаря:
dict = Dict{Integer,Vector{String}}{()}

# Инициализация словаря:
dict2 = Dict{Integer,Vector{String}}{()}

# Заполнение словаря данными:
for i = 1:size(P,1)
    year,lang = P[i,:]
    if year in keys(dict)
        dict[year] = push!(dict[year],lang)
    else
        dict[year] = [lang]
    end
end

# Пример определения в словаре языков программирования, созданных в 2003 году:
dict[2003]

DataFrames

# Подгружаем пакет DataFrames:
using DataFrames
# Задаём переменную со структурой DataFrame:
df = DataFrame(year = P[:,1], language = P[:,2])

# Вывод всех значения столбца year:
df[:,year]

# Получение статистических сведений о фрейме:
describe(df)

```

Рис. 3.2: Примеры

```

Работа с переменными отсутствующего типа (Missing Values)

# Отсутствующий тип:
a = missing
typeof(a)

# Пример операции с переменной отсутствующего типа:
a + 1

# Определение перечня продуктов:
foods = ["apple", "cucumber", "tomato", "banana"]
# Определение калорий:
calories = [missing, 47, 22, 105]

# Определение типа переменной:
typeof(calories)

# Подключаем пакет Statistics:
using Statistics
# Определение среднего значения:
mean(calories)

# Определение среднего значения без значений с отсутствующим типом:
mean(skipmissing(calories))

# Задание сведений о ценах:
prices = [0.85, 1.6, 0.8, 0.6]
# Формирование данных о калориях:
dataframe_calories = DataFrame(item=foods, calories=calories)
# Формирование данных о ценах:
dataframe_prices = DataFrame(item=foods, price=prices)
# Объединение данных о калориях и ценах:
DF = innerjoin(dataframe_calories, dataframe_prices, on=:item)

```

Рис. 3.3: Примеры


```

# Загрузка данных:
houses = CSV.File("houses.csv") |> DataFrame

# Построение графика:
using Plots
plot(size=(500,500),leg=false)
x = houses[!,:sq_ft]
y = houses[!,:price]
scatter(x,y,markersize=3)

# Фильтрация данных по заданному условию:
filter_houses = houses[houses[!,:sq_ft].>0,:]
# Построение графика:
x = filter_houses[!,:sq_ft]
y = filter_houses[!,:price]
scatter(x,y)

# Подключение пакета Statistics:
using Statistics
# Определение средней цены для определённого типа домов:
combine(groupby(filter_houses,[:type]),filter_houses->mean(filter_houses[!,:price]))

using Clustering
# Добавление данных :latitude и :longitude в новый фрейм:
X = filter_houses[!,:latitude,:longitude]
# Конвертация данных в матричный вид:
X = Matrix{X}

# Транспонирование матрицы с данными:
X = X'

# Задание количества кластеров:
k = length(unique(filter_houses[!,:zip]))

# Определение k-среднего:
C = kmeans(X,k)

# Формирование фрейма данных:
df = DataFrame(cluster = C.assignments,city = filter_houses[!,:city],
latitude = filter_houses[!,:latitude],longitude =
filter_houses[!,:longitude],zip = filter_houses[!,:zip])

clusters_figure = plot(legend = false)
for i = 1:k
    clustered_houses = df[df[!,:cluster].== i,:]
    xvals = clustered_houses[!,:latitude]
    yvals = clustered_houses[!,:longitude]
    scatter!(clusters_figure,xvals,yvals,markersize=4)
end
xlabel!("Latitude")
ylabel!("Longitude")
title!("Houses color-coded by cluster")
display(clusters_figure)

unique_zips = unique(filter_houses[!,:zip])
zips_figure = plot(legend = false)
for uzip in unique_zips
    subs = filter_houses[filter_houses[!,:zip].==uzip,:]
    x = subs[!,:latitude]
    y = subs[!,:longitude]
    scatter!(zips_figure,x,y)
end
xlabel!("Latitude")
ylabel!("Longitude")
title!("Houses color-coded by zip code")
display(zips_figure)

```

Рис. 3.4: Примеры

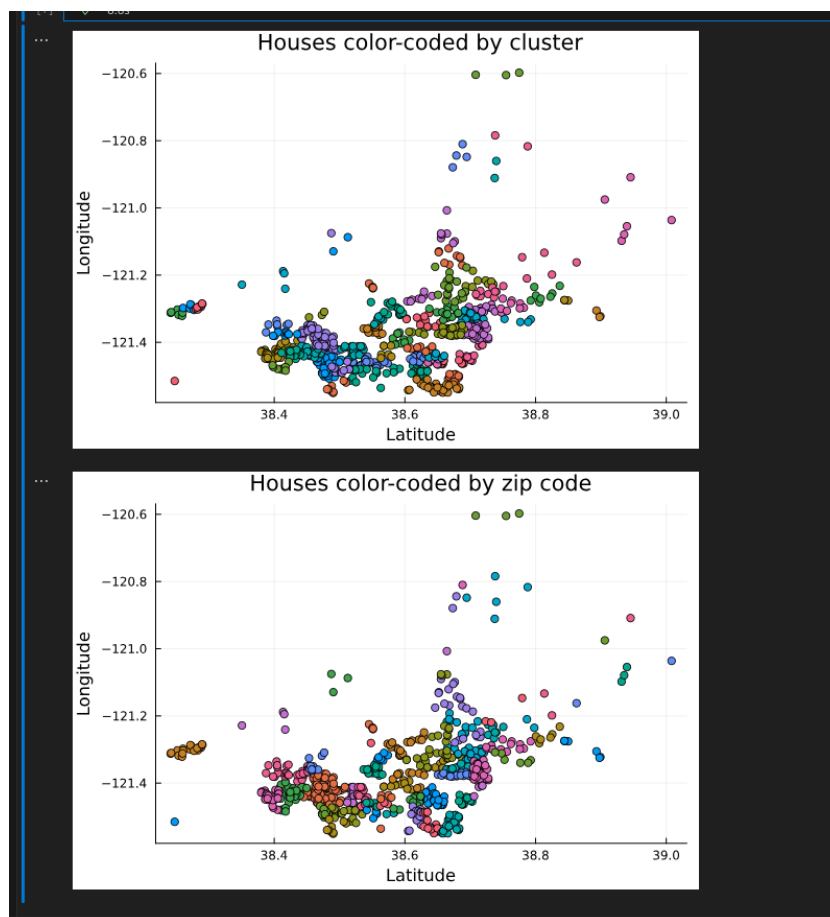


Рис. 3.5: Примеры



Рис. 3.6: Примеры

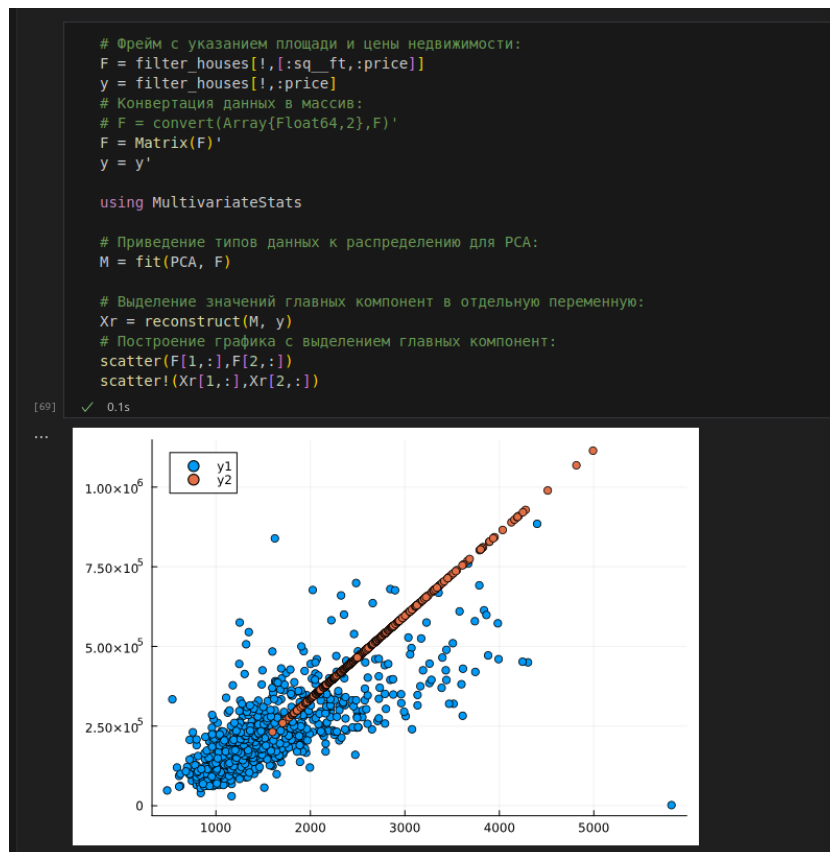


Рис. 3.7: Примеры

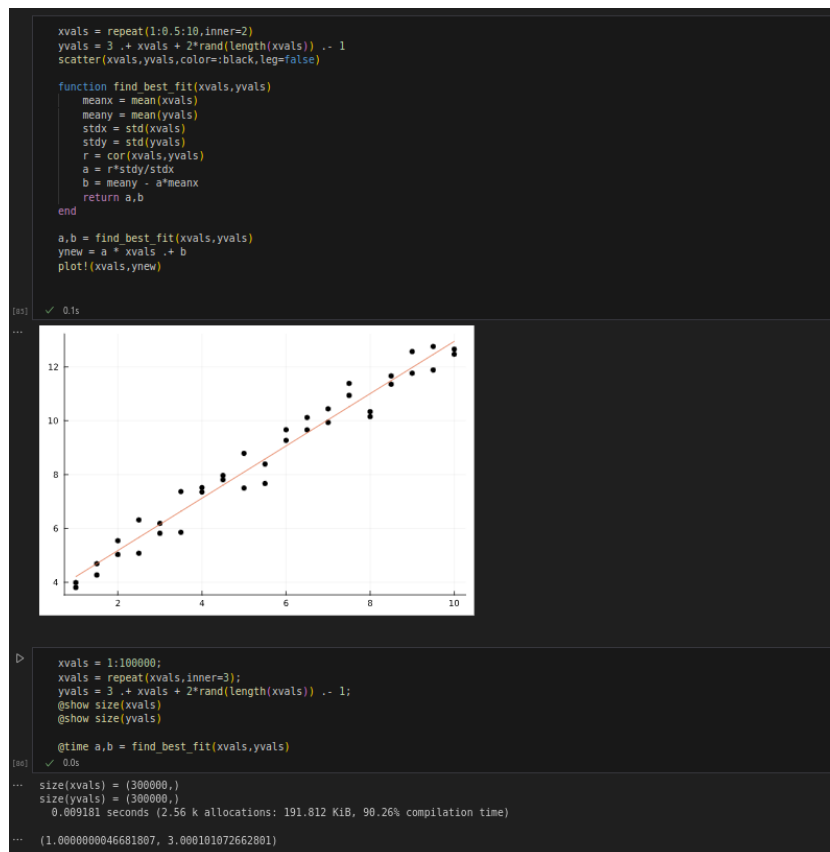


Рис. 3.8: Примеры

Затем выполним задания(рис. 3.9 - 3.15)



Рис. 3.9: Задание 1

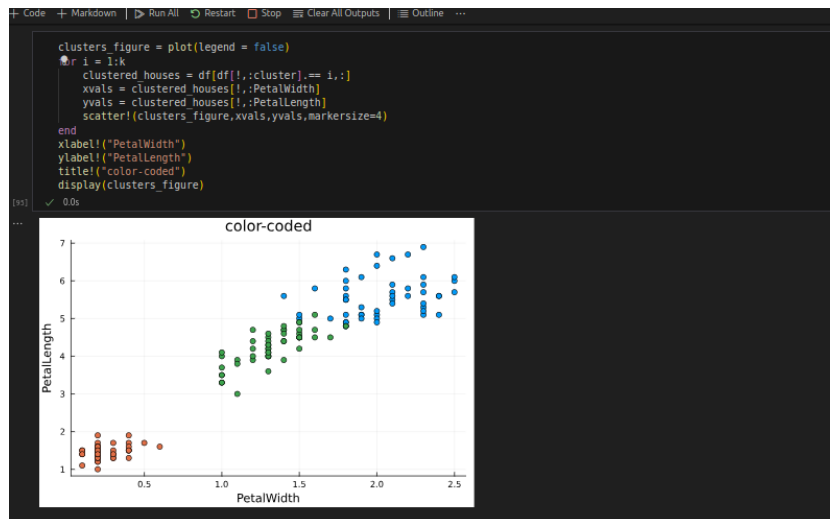


Рис. 3.10: Задание 1



Рис. 3.11: Задания 2



Рис. 3.12: Задание 2

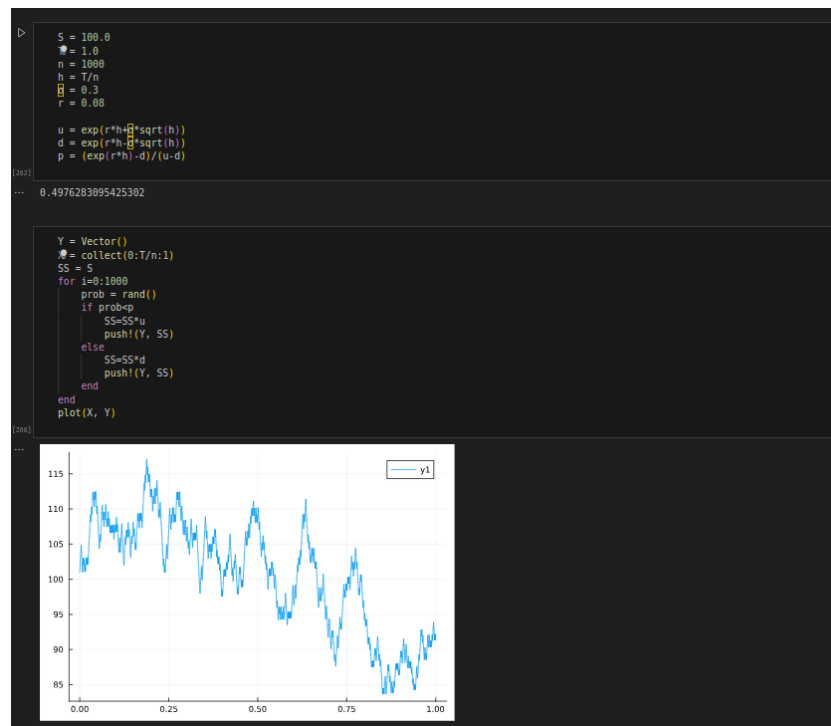


Рис. 3.13: Задание 3

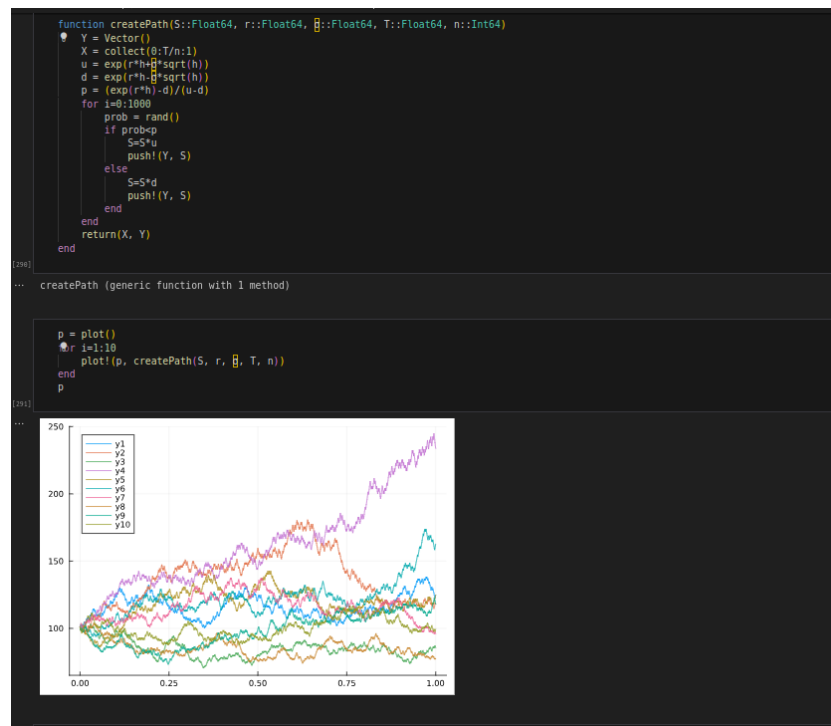


Рис. 3.14: Задание 3

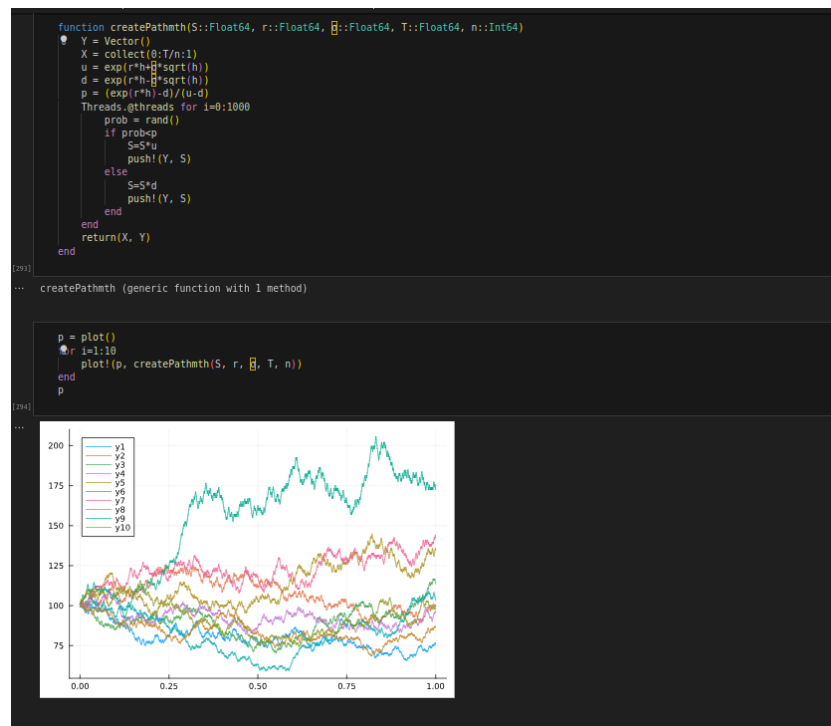


Рис. 3.15: Задание 3

4 Выводы

В результате выполнения работы освоили использование специализированных пакетов Julia для обработки данных.

Список литературы

1. JuliaLang [Электронный ресурс]. 2024 JuliaLang.org contributors. URL: <https://julialang.org/> (дата обращения: 11.10.2024).
2. Julia 1.11 Documentation [Электронный ресурс]. 2024 JuliaLang.org contributors. URL: <https://docs.julialang.org/en/v1/> (дата обращения: 11.10.2024).