

Лабораторная работа № 7. Введение в работу с данными

Компьютерный практикум по статистическому анализу данных

Демидова Е. А.

24 ноября 2024

Российский университет дружбы народов, Москва, Россия

Информация

- Демидова Екатерина Алексеевна
- студентка группы НКНбд-01-21
- Российский университет дружбы народов
- <https://github.com/eademidova>



Введение

Цель работы

Основной целью работы является специализированных пакетов Julia для обработки данных.

Задачи

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

Выполнение лабораторной работы

```
D> using CSV, DataFrames, DelimitedFiles, Plots, Statistics, GLM
[296]

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 = readln("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 ',':
writeln("programming_languages_data.txt", Tx, ',')
# Пример записи данных в текстовый файл с разделителем '-':
writeln("programming_languages_data2.txt", Tx, '-')

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

Рис. 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. 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: Примеры

Выполнение примеров

```
# R script example:
houses = CSV.File("houses.csv") |> DataFrame

# Basic house plot:
using Plots
plot(houses[1:100, :], log=false)
x = houses[1:100, :].x
y = houses[1:100, :].y
scatter(x, y, markersize=3)

# Basic house plot with longitude:
filter_houses = filter(houses[1:100, :])
# Basic house plot:
x = filter_houses[1:100, :].x
y = filter_houses[1:100, :].y
scatter(x, y)

# Basic house plot Statistics:
using Statistics
# Basic house plot:
combine(groupby(filter_houses, [:type]), filter_houses -> mean(filter_houses[1:100, :].price))

using Clustering
# Basic house plot:
X = filter_houses[1:100, :].latitude, longitude
# Basic house plot:
X = Matrix(X)

# Basic house plot:
X = X'

# Basic house plot:
k = length(unique(filter_houses[1:100, :].zip))

# Basic house plot:
C = kmeans(X, k)

# Basic house plot:
df = DataFrame(cluster = C.assignments, city = filter_houses[1:100, :].city,
latitude = filter_houses[1:100, :].latitude, longitude =
filter_houses[1:100, :].longitude, zip = filter_houses[1:100, :].zip)

clusters_figure = plot(legend = false)
for i = 1:k
    clustered_houses = df[df[1:100, :].city == i, :]
    xvals = clustered_houses[1:100, :].latitude
    yvals = clustered_houses[1:100, :].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[1:100, :].zip)
zips_figure = plot(legend = false)
for zip in unique_zips
    subs = filter_houses[filter_houses[1:100, :].zip == zip, :]
    x = subs[1:100, :].latitude
    y = subs[1:100, :].longitude
    scatter!(zips_figure, x, y)
end
xlabel!("Latitude")
ylabel!("Longitude")
title!("houses color-coded by zip code")
display(zips_figure)
```

Рис. 4: Примеры

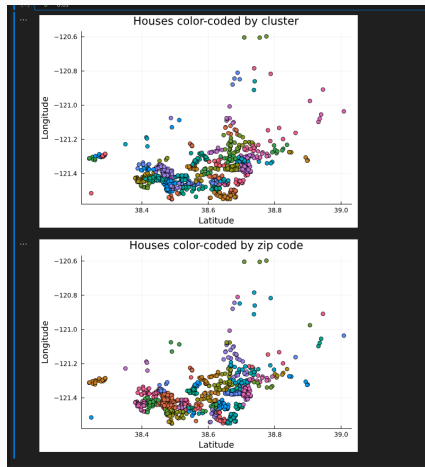


Рис. 5: Примеры



Рис. 6: Примеры

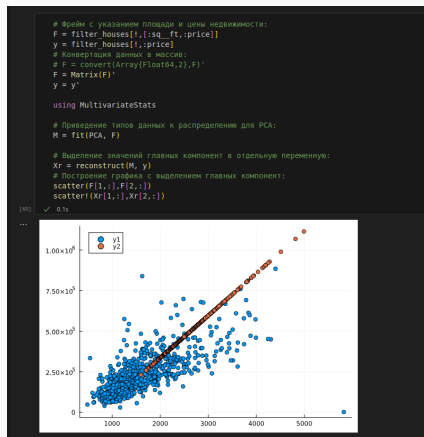


Рис. 7: Примеры

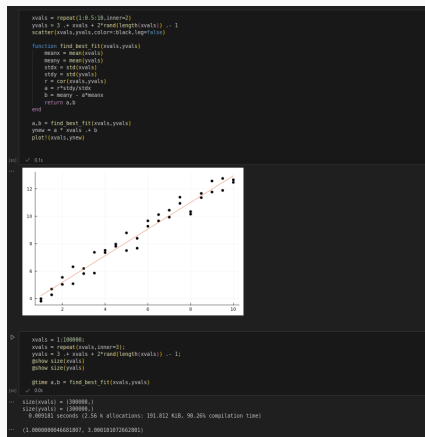


Рис. 8: Примеры

Выполнение заданий для самостоятельной работы



Рис. 9: Задание 1

Выполнение заданий для самостоятельной работы

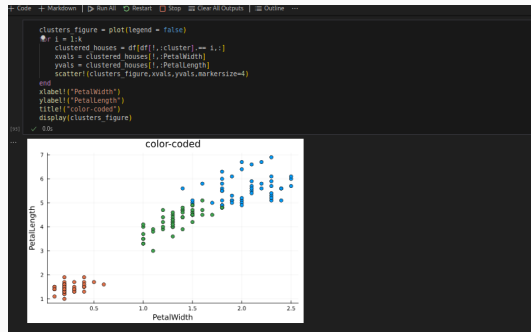


Рис. 10: Задание 1

Выполнение заданий для самостоятельной работы



Рис. 11: Задания 2

Выполнение заданий для самостоятельной работы



Рис. 12: Задание 2

Выполнение заданий для самостоятельной работы

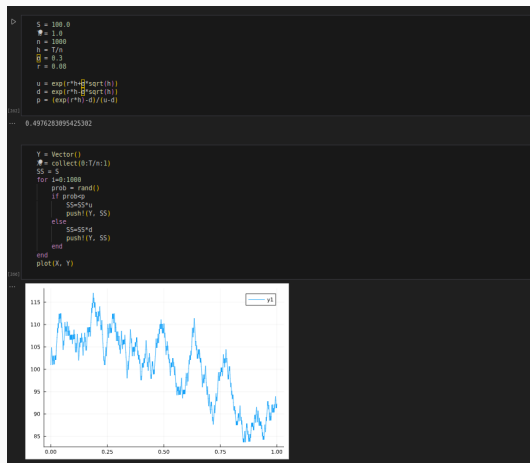


Рис. 13: Задание 3

Выполнение заданий для самостоятельной работы

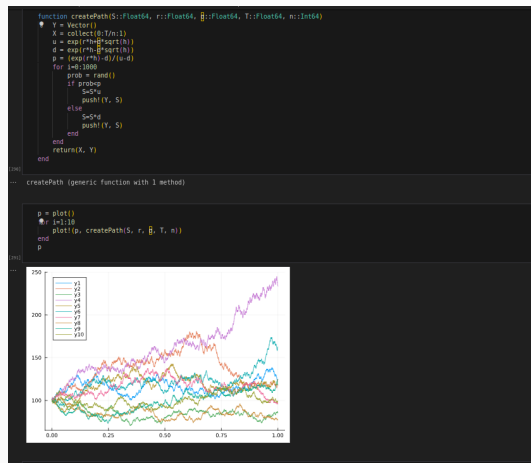


Рис. 14: Задание 3

Выполнение заданий для самостоятельной работы

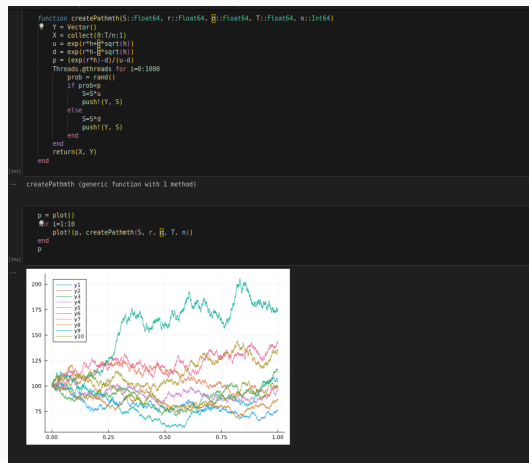


Рис. 15: Задание 3

Выводы

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

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).