

РОССИЙСКИЙ УНИВЕРСИТЕТ ДРУЖБЫ НАРОДОВ

Факультет физико-математических и естественных наук

Кафедра прикладной информатики и теории вероятностей

ОТЧЕТ

ПО ЛАБОРАТОРНОЙ РАБОТЕ № 7

дисциплина: Компьютерный практикум

по математическому моделированию

Студент: Ли Тимофей Александрович

Группа: НФИбд-01-18

МОСКВА

2021 г.

Постановка задачи

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

Выполнение работы

Сначала выполнил все примеры к лабораторной работе №7:

1. Ознакомился со считыванием данных.

```
# Обновление окружения:
using Pkg
Pkg.update
# Установка пакетов:
using Pkg
for p in ["CSV", "DataFrames", "RDatasets", "FileIO"]
    Pkg.add(p)
end
using CSV, DataFrames, DelimitedFiles
```

```
P = CSV.File("programminglanguages.csv") |> DataFrame
```

8 rows × 2 columns

	year	language
	Int64	String7
1	1950	basic
2	1951	fortran
3	1952	pascal
4	1953	c
5	1954	java
6	1955	prolog
7	1956	r
8	1957	python

```
function language_created_year(P, language::String)
    loc = findfirst(P[:,2].==language)
    return P[loc,1]
end
```

language_created_year (generic function with 1 method)

```
language_created_year(P, "python")
```

1957

```
language_created_year(P, "prolog")
```

1955

```
language_created_year(P, "Prolog")
```

MethodError: no method matching getindex(::DataFrame, ::Nothing, ::Int64)

```
function language_created_year_v2(P, language::String)
    loc = findfirst(lowercase.(P[:,2]).==lowercase.(language))
    return P[loc,1]
end
```

language_created_year_v2 (generic function with 1 method)

```
language_created_year_v2(P, "Prolog")
```

1955

```
Tx = readlm("programminglanguages.csv", ',')
```

9×1 Matrix{Any}:
"year;language"
"1950;basic"
"1951;fortran"
"1952;pascal"
"1953;c"
"1954;java"
"1955;prolog"
"1956;r"
"1957;python"

2. Затем ознакомился с записью данных в файл.

```
In [11]: CSV.write("programming_languages_data2.csv", P)
```

```
Out[11]: "programming_languages_data2.csv"
```

```
In [12]: writelml("programming_languages_data.txt", Tx, ',')
```

```
In [13]: # Пример записи данных в текстовый файл с разделителем '-':
writelml("programming_languages_data2.txt", Tx, '-')
```

```
In [14]: P_new_delim = readlm("programming_languages_data2.txt", '-')
```

```
Out[14]: 9×1 Matrix{Any}:
"year;language"
"1950;basic"
"1951;fortran"
"1952;pascal"
"1953;c"
"1954;java"
"1955;prolog"
"1956;r"
"1957;python"
```

3. Далее ознакомился со словарями.

```
In [15]: dict = Dict{Integer,Vector{String}}{()}
          dict2 = Dict{Integer,Vector{String}}{()}

Out[15]: Dict{Any, Any}{}

In [16]: 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

In [18]: dict[1954]

Out[18]: 1-element Vector{String}:
         "java"
```

4. Ознакомился с DataFrames

```
using DataFrames

df = DataFrame(year = P[:,1], language = P[:,2])

8 rows × 2 columns
```

	year	language
	Int64	String7
1	1950	basic
2	1951	fortran
3	1952	pascal
4	1953	c
5	1954	java
6	1955	prolog
7	1956	r
8	1957	python

```
df[!,:year]

8-element Vector{Int64}:
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957

describe(df)

2 rows × 7 columns
```

	variable	mean	min	median	max	nmissing	eltype
	Symbol	Union...	Any	Union...	Any	Int64	DataType
1	year	1953.5	1950	1953.5	1957	0	Int64
2	language		basic		r	0	String7

5. Ознакомился с RDatasets

```
using RDatasets

iris = dataset("datasets", "iris")

150 rows × 5 columns
```

	SepalLength	SepalWidth	PetalLength	PetalWidth	Species
	Float64	Float64	Float64	Float64	Cat...
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa
9	4.4	2.9	1.4	0.2	setosa
10	4.9	3.1	1.5	0.1	setosa

```
typeof(iris)

DataFrame

describe(iris)

5 rows × 7 columns
```

	variable	mean	min	median	max	nmissing	eltype
	Symbol	Union...	Any	Union...	Any	Int64	DataType
1	SepalLength	5.84333	4.3	5.8	7.9	0	Float64
2	SepalWidth	3.05733	2.0	3.0	4.4	0	Float64
3	PetalLength	3.758	1.0	4.35	6.9	0	Float64
4	PetalWidth	1.19933	0.1	1.3	2.5	0	Float64
5	Species	setosa		virginica		0	CategoricalValue{String, UInt8}

6. Ознакомился с работой с missing переменными

```
a = missing
typeof(a)
```

Missing

a+1

missing

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

```
4-element Vector{Union{Missing, Int64}}:
 missing
 47
 22
105
```

```
typeof(calories)
```

```
Vector<Union<Missing, Int64>> (alias for Array<Union<Missing, Int64>, 1))
```

```
using Statistics
```

```
mean(calories)
```

missing

```
mean(skipmissing(calories))
```

58.0

```
# 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')
```

4 rows \times 3 columns

	item	calories	price
	String	Int64?	Float64
1	apple	missing	0.85
2	cucumber	47	1.6
3	tomato	22	0.8
4	banana	105	0.6

7. Ознакомился с FileIO

```
using FileIO
```

```
import Pkg
Pkg.add("ImageIO")
```

```
Resolving package versions...
Installed Imath_jll ————— v3.1.2+0
Installed AbstractFFTs ——— v1.0.1
```

```
X1 = load("julialogo.png")
```

[illegible]

```
@show typeof(X1);
@show size(X1);

typeof(X1) = Matrix{ColorTypes.RGBA{FixedPointNumbers.N0f8}}
size(X1) = (776, 1200)
```

8. Ознакомился с кластеризацией методом k-means

```
import Pkg
Pkg.add("DataFrames")
Pkg.add("Statistics")
using DataFrames
using CSV
import Pkg
Pkg.add("Plots")
```

```

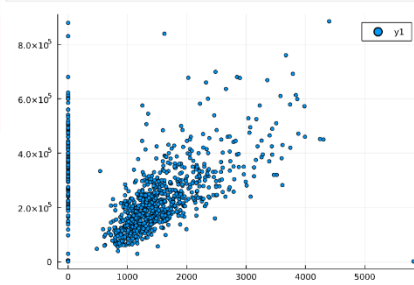
Resolving package versions...
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Project.toml'
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Manifest.toml'
Resolving package versions...
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Project.toml'
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Manifest.toml'
Resolving package versions...
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Project.toml'
No Changes to 'C:\Users\Xiaoaml\julia\environments\v1.6\Manifest.toml'

```

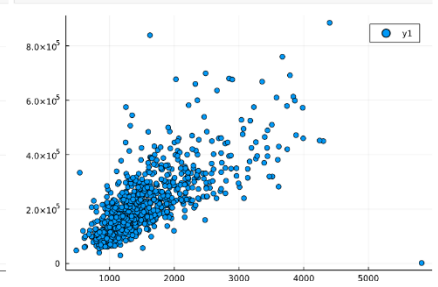
```
houses = CSV.File("houses.csv") |> DataFrame
985 rows × 12 columns (omitted printing of 6 columns)
```

	street	city	zip	state	beds	baths
	String63	String15	Int64	String3	Int64	Int64
1	3526 HIGH ST	SACRAMENTO	95838	CA	2	1
2	51 OMAHA CT	SACRAMENTO	95823	CA	3	1
3	2766 BRANCH ST	SACRAMENTO	95815	CA	2	1
4	2095 JANETTE WAY	SACRAMENTO	95845	CA	2	1

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



```
filter_houses = houses[houses[1,:sq_ft].>0,:]
# Построение графика:
x = filter_houses[1,:sq_ft]
y = filter_houses[1,:price]
scatter(x,y)
```



```
using Statistics
# Определить среднюю цену для определенного типа домов:
combine(filter_houses, :type, filter_houses -> mean(filter_houses[!, :price]))
```

814 rows × 2 columns

	type	x1
	String15	Float64
1	Residential	2.204480e5
2	Residential	2.204480e5
3	Residential	2.204480e5
4	Residential	2.204480e5
5	Residential	2.204480e5

```
import Pkg
Pkg.add("Clustering")
using Clustering
```

Resolving package versions...

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

814x2 Matrix{Float64}:

38.6319	-121.435
38.4789	-121.431
38.6183	-121.444
38.6168	-121.439
38.5195	-121.436
38.6626	-121.378
38.6817	-121.352
38.5351	-121.481
38.6212	-121.271
38.7009	-121.443
38.6177	-121.452
38.4707	-121.459
38.6187	-121.436

X=X'

```
2x814 adjoint(::Matrix{Float64}) with eltype Float64:
 38.6319 38.4789 38.6183 ... 38.7088 38.417 38.6552
-121.435 -121.431 -121.444 ... -121.257 -121.397 -121.076
```

```
k = length(unique(filter_houses[!, :zip]))
C = kmeans(X, k)
```

```
KmeansResult{Matrix{Float64}, Float64, Int64}([38.56742925 38.490447 ... 38.53068435714287 38.69875490476191; -121.49139358333333
-121.129337 ... -121.46543535714284 -121.37119928571428], [60, 42, 53, 60, 11, 4, 21, 29, 26, 33 ... 6, 39, 42, 14, 43, 54, 59, 2
8, 30, 13], [6.734555427101411e-5, 0.00037363141382229514, 0.00012962343315758052, 0.00017246662901015952, 0.00015398280185181
65, 0.0001611708285054192, 0.00025366112095071003, 0.00014704814748256467, 6.344488792819902e-5, 0.00011686658399412408 ... 5.7
84678793218538e-5, 5.2413251978578046e-5, 9.498398867435753e-6, 0.0001942098642757628, 0.00013415837747743353, 0.00014433951582
75962, 7.579613156849518e-5, 8.884414273779839e-5, 0.0003106035837845411, 0.0002109576998918783], [12, 1, 11, 18, 12, 28, 8, 1
9, 14, 28 ... 6, 3, 17, 18, 6, 4, 14, 1, 14, 21],
0.20318177624358214, 16, true)
```

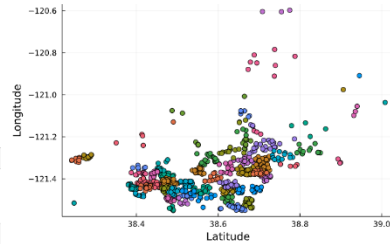
```
df = DataFrame(cluster = C.assignments, city = filter_houses[!, :city],
latitude = filter_houses[!, :latitude], longitude = filter_houses[!, :longitude],
zip = filter_houses[!, :zip])
```

814 rows × 5 columns

	cluster	city	latitude	longitude	zip
	Int64	String15	Float64	Float64	Int64
1	60	SACRAMENTO	38.6319	-121.435	95838
2	42	SACRAMENTO	38.4789	-121.431	95823
3	53	SACRAMENTO	38.6183	-121.444	95815
4	60	SACRAMENTO	38.6168	-121.439	95815
5	11	SACRAMENTO	38.5195	-121.436	95824

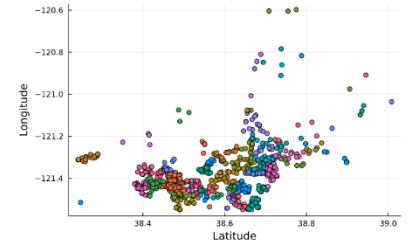
```
clusters_figure = plot(legend = false)
for i = 1:k
    clustered_houses = df[df[!, :cluster].== i, :]
    xvals = clustered_houses[!, :latitude]
    yvals = clustered_houses[!, :longitude]
    scatter!(clusters_figure, xvals, yvals, markersize=4)
end
xlabel!("Latitude")
ylabel!("Longitude")
title!("Houses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

Houses color-coded by cluster



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

Houses color-coded by zip code



9. Ознакомился с кластеризацией методом к ближайших соседей

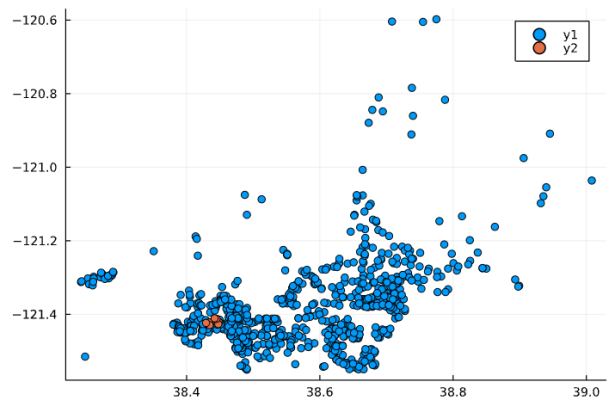
```
import Pkg
Pkg.add("NearestNeighbors")
using NearestNeighbors
```

Resolving package versions...
Updating `C:\Users\Xiaomi\.julia\environments\v1.6\Project.toml`
[b8a86587] + NearestNeighbors v0.4.9
No Changes to `C:\Users\Xiaomi\.julia\environments\v1.6\Manifest.toml`

```
knearest = 10
id = 70
point = X[:, id]
```

```
2-element Vector{Float64}:
 38.44004
-121.421012
```

```
# Поиск ближайших соседей:
kdtree = KDTree(X)
idxs, dists = knn(kdtree, point, knearest, true)
# Все объекты недвижимости:
x = filter_houses[!, :latitude];
y = filter_houses[!, :longitude];
scatter(x, y)
# Соседи:
x = filter_houses[idxs, :latitude];
y = filter_houses[idxs, :longitude];
scatter!(x, y)
```



```
cities = filter_houses[idxs, :city]
```

```
10-element PooledArrays.PooledVector{String15, UInt32, Vector{UInt32}}:
 "SACRAMENTO"
 "ELK GROVE"
 "SACRAMENTO"
 "SACRAMENTO"
 "SACRAMENTO"
 "SACRAMENTO"
 "ELK GROVE"
 "ELK GROVE"
 "ELK GROVE"
 "ELK GROVE"
```

10. Ознакомился с кластеризацией методом главных компонент

```
F = filter_houses[:,[:sq_ft,:price]]
F = Array{F}'

2×814 adjoint(::Matrix{Int64}) with eltype Int64:
 836  1167  796  852  797  1122  ...  1477  1216  1685  1362
59222 68212 68880 69307 81900 89921 234000 235000 235301 235738

convert(Array{Float64,2}, F)

2×814 Matrix{Float64}:
 836.0  1167.0  796.0  852.0  797.0  ...  1216.0  1685.0  1362.0
59222.0 68212.0 68880.0 69307.0 81900.0 235000.0 235301.0 235738.0

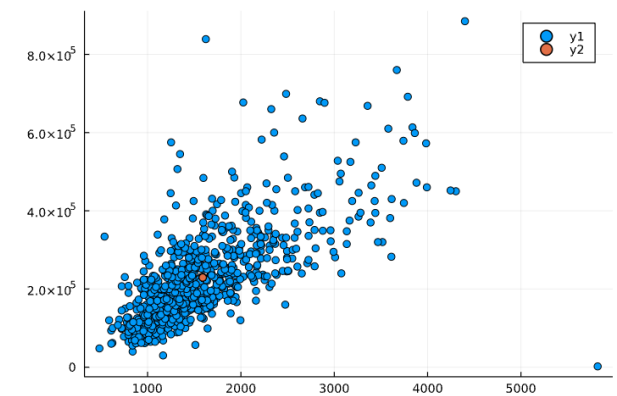
import Pkg
Pkg.add("MultivariateStats")
using MultivariateStats

Resolving package versions...
No Changes to `C:\Users\Xiaomi\.julia\environments\v1.6\Project.toml`
No Changes to `C:\Users\Xiaomi\.julia\environments\v1.6\Manifest.toml`

M = fit(PCA, F)

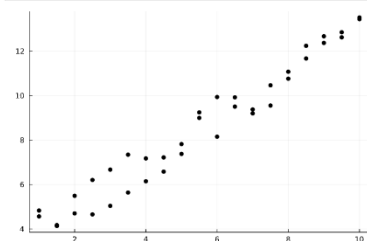
PCA(indim = 2, outdim = 1, principalratio = 0.9999840784692097)
```

```
# Выделение значений главных компонент в отдельную переменную:
Xr = reconstruct(M, y')
# Построение графика с выделением главных компонент:
scatter(F[1,:],F[2,:])
scatter!(Xr[1,:],Xr[2,:])
```



11. Ознакомился с линейной регрессией

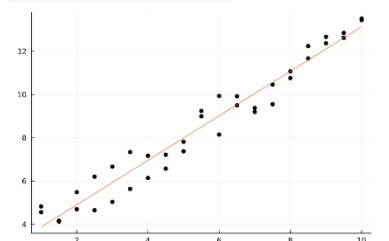
```
xvals = repeat(1:0.5:10,inner=2)
yvals = 3 + xvals + 2*rand(length(xvals)) .* 1
scatter(xvals,yvals,color=:black,leg=false)
```



```
function find_best_fit(xvals,yvals)
    meanx = mean(xvals)
    meany = mean(yvals)
    stdx = std(xvals)
    stdy = std(yvals)
    r = cor(xvals,yvals)
    a = r*stdy/stdx
    b = meany - a*meanx
    return a,b
end
```

```
find_best_fit (generic function with 1 method)
```

```
a,b = find_best_fit(xvals,yvals)
ynew = a * xvals + b
plot!(xvals,ynew)
```



```
xvals = 1:100000;
xvals = repeat(xvals,inner=3);
yvals = 3 + xvals + 2*rand(length(xvals)) .* 1;

@show size(xvals)
@show size(yvals)

size(xvals) = (300000,)
size(yvals) = (300000,)
(300000,)

@time a,b = find_best_fit(xvals,yvals)

0.136022 seconds (103.56 k allocations: 6.250 MiB, 98.48% compilation time)
(0.9999999583646021, 3.004228580146446)
```

```
import Pkg
Pkg.add("PyCall")
Pkg.add("Conda")
using PyCall
using Conda

Resolving package versions...
```

```
py"""
import numpy
def find_best_fit_python(xvals,yvals):
    meanx = numpy.mean(xvals)
    meany = numpy.mean(yvals)
    stdx = numpy.std(xvals)
    stdy = numpy.std(yvals)
    r = numpy.corrcoef(xvals,yvals)[0][1]
    a = r*stdy/stdx
    b = meany - a*meanx
    return a,b
"""

find_best_fit_python"""find_best_fit_python"""
xpy = PyObject(xvals)
ypy = PyObject(yvals)
@time a,b = find_best_fit_python(xpy,ypy)

0.667533 seconds (237.03 k allocations: 15.021 MiB, 23.34% gc time, 84.55% compilation time)
(0.9999999583646082, 3.004228580284689)
```

```
import Pkg
Pkg.add("BenchmarkTools")
using BenchmarkTools

Resolving package versions...
```

```
@time a,b = find_best_fit_python(xvals,yvals)
@time a,b = find_best_fit(xvals,yvals)

10.267 ms (27 allocations: 1.03 KiB)
1.568 ms (1 allocation: 32 bytes)
(0.9999999583646021, 3.004228580146446)
```

Далее выполнил поставленные задачи.

1. Кластеризовал набор данных iris методом k-means.

Для этого загрузил датасет, проиндексировал его, преобразовал в массив и транспонировал. Определил количество кластеров и нашел k-среднее.

Затем, создал новый датафрейм из изначального с добавлением столбца, в котором хранится метка кластера. Далее просто построил точечную диаграмму кластеров.

Далее распараллелил цикл отрисовки траекторий с помощью `Threads.@threads`

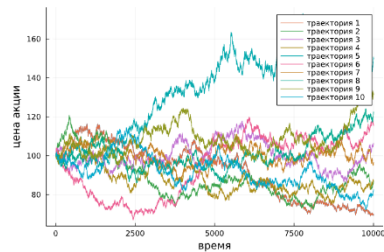
```
S=100
T=1
n=10000
h=T/n
sigma=0.3
r=0.08
u=exp(r*h + sigma*sqrt(h))
d=exp(r*h - sigma*sqrt(h))
p=(exp(r*h)-d)/(u-d)
0.4992500005625153
```

```
line=[]
c=0
append!(line,S)
for i in 1:n
    k=rand()
    if k < p
        append!(line,S*u*(1-c)*d*c)
    else
        append!(line,S*u*(1-c-1)*d*(c+1))
        c+=1
    end
end
plot(line, title="траектория курса акций", xlabel="время", ylabel="цена акции")
```



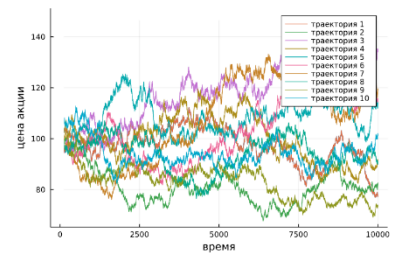
```
function CreatePath(S,r,sigma,T,n)
    h=T/n
    u=exp(r*h + sigma*sqrt(h))
    d=exp(r*h - sigma*sqrt(h))
    p=(exp(r*h)-d)/(u-d)
    line=[]
    c=0
    append!(line,S)
    for i in 1:n
        k=rand()
        if k < p
            append!(line,S*u*(1-c)*d*c)
        else
            append!(line,S*u*(1-c-1)*d*(c+1))
            c+=1
        end
    end
    return line
end
CreatePath (generic function with 1 method)
```

```
for i in 1:10
    IJulia.clear_output(true)
    path=CreatePath(100,0.08,0.3,1,10000)
    if i==1
        p=plot(path, label=false)
    end
    p=plot!(path, label="траектория $i", xlabel="время", ylabel="цена акции")
    display(p)
end
```



```
using Base.Threads

Threads.@threads for i in 1:10
    IJulia.clear_output(true)
    path=CreatePath(100,0.08,0.3,1,10000)
    if i==1
        p=plot(path, label=false)
    end
    p=plot!(path, label="траектория $i", xlabel="время", ylabel="цена акции")
    display(p)
end
```



Выводы

Изучил специализированные пакеты Julia для обработки данных.