

# Лабораторная работа №7

Введение в работу с данными

---

Легиньких Г.А.

Российский университет дружбы народов, Москва, Россия

## Информация

---

- Легиньких Галина Андреевна
- НФИбд-02-21
- Российский университет дружбы народов
- 1032216447@pfur.ru
- <https://github.com/galeginkikh>

## Основная информация

---

Основной целью работы является специализированных пакетов Julia для обработки данных.

1. Используя Jupyter Lab, повторите примеры из раздела 7.2.
2. Выполните задания для самостоятельной работы (раздел 7.4).

## Выполнение лабораторной работы

---

1. Для начала я повторила примеры по данной теме. Начала с раздела считывание данных. В Julia для работы с такого рода структурами данных используют пакеты CSV, DataFrames, RDatasets, FileIO. Скачала эти пакеты и попробовала прочитать данные из скаченного датафрейма.

```
[6]:  
# Считывание данных и их запись в структуру:  
P = CSV.File("programminglanguages.csv") |> DataFrame  
  
[6]:  
73×2 DataFrame                                     48 rows omitted  
Row  year  language  
   Int64  String31  
-----  
 1  1951  Regional Assembly Language  
 2  1952  Autocode  
 3  1954  IPL
```

Рис. 1: Считывание данных



Далее приведём пример функции, в которой на входе указывается название языка программирования, а на выходе — год его создания.

```
[10]:  
# 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")  
  
[10]:  
1991
```

Рис. 2: Функция язык-год

При этом, чтобы убрать в функции зависимость данных от регистра, необходимо изменить исходную функцию следующим образом.

```
[14]:  
# 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  
# Пример вызова функции и определение даты создания языка julia:  
language_created_year_v2(P, "julia")  
  
[14]:  
2012
```

Рис. 3: Зависимость от регистра

2. Далее попробовала записать данные в файл, а так же задала при сохранении тип файла и разделитель.

3. При работе с данными бывает удобно записать их в формате словаря. При инициализации словаря можно задать конкретные типы данных для ключей и значений, а можно инициализировать пустой словарь, не задавая строго структуру.

```
[24]:  
# Инициализация словаря:  
dict = Dict{Integer, Vector{String}}()  
  
[24]:  
Dict{Integer, Vector{String}}()  
  
[26]:  
# Инициализация словаря:  
dict2 = Dict()  
  
[26]:  
Dict{Any, Any}()
```

Рис. 4: Инициализация словаря

В результате при вызове словаря можно, выбрав любой год, узнать, какие языки программирования были созданы в этом году.

```
[28]:  
  
# Заполнение словаря данными:  
for i = 1:size(P,1)  
    year,lang = P[i,i]  
    if year in keys(dict)  
        dict[year] = push!(dict[year],lang)  
    else  
        dict[year] = [lang]  
    end  
end  
  
dict[2003]  
  
[28]:  
  
2-element Vector{String}:  
 "Groovy"  
 "Scala"
```

Рис. 5: Работа со словарем

4. На примере с данными о языках программирования и годах их создания зададим структуру DataFrame. Попробовала вывести один столбец.

```
[32]:  
# задаём переменную со структурой DataFrame:  
df = DataFrame(year = P[i,1], language = P[i,2])  
  
[32]:  
73×2 DataFrame                                     48 rows omitted  


| Row | year  | language                   |
|-----|-------|----------------------------|
|     | Int64 | String31                   |
| 1   | 1951  | Regional Assembly Language |
| 2   | 1952  | Autocode                   |


```

Рис. 6: Переменная со структурой датафрейм

А так же получила статистические сведения о фрейме.

```
[36]:  
# Получение статистических сведений о фрейме:  
describe(df)
```

[36]:  
2x7 DataFrame

Row	variable	mean	min	median	max	nmissing	eltype
	Symbol	Union...	Any	Union...	Any	Int64	DataType
1	year	1982.99	1951	1986.0	2014	0	Int64
2	language		ALGOL 58		dBase III	0	String31

Рис. 7: Статистические сведения о фрейме

5. С данными можно работать также как с наборами данных через пакет RDatasets языка R.



6. Попробовала поработать с переменными отсутствующего типа.

```
[46]:  
# Отсутствующий тип:  
a = missing  
typeof(a)
```

```
[46]:  
Missing
```

**Рис. 8:** Переменная отсутствующего типа

7. Перешла к обработке данных. И первый метод это метод k-средних. Задача кластеризации данных заключается в формировании однородной группы упорядоченных по какому-то признаку данных. Метод k-средних позволяет минимизировать суммарное квадратичное отклонение точек кластеров от центров этих кластеров.

Рассмотрела задачу кластеризации данных на примере данных о недвижимости.

```
[72]:
```

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

```
[72]:
```

985×12 DataFrame *960 rows omitted*

Row	street	city	zip	state	beds	baths	sq_ft	type	s
	String	String15	Int64	String3	Int64	Int64	Int64	String15	String15
1	3526 HIGH ST	SACRAMENTO	95838	CA	2	1	836	Residential	1
2	51 OMAHA CT	SACRAMENTO	95823	CA	3	1	1167	Residential	1
3	2796 BRANCH ST	SACRAMENTO	95815	CA	2	1	796	Residential	1

Рис. 9: Кластеризация данных на примере данных о недвижимости

## Выполнение лабораторной работы

[74]:

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

[74]:

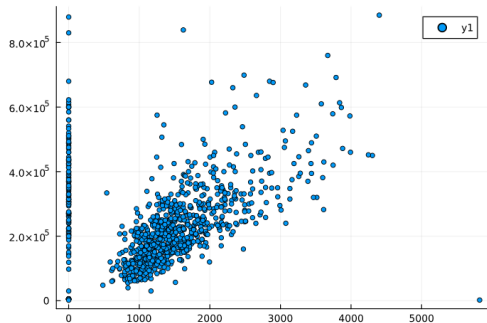


Рис. 10: Цены на недвижимость в зависимости от площади

## Выполнение лабораторной работы

[76]:

```
# 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)
```

[76]:

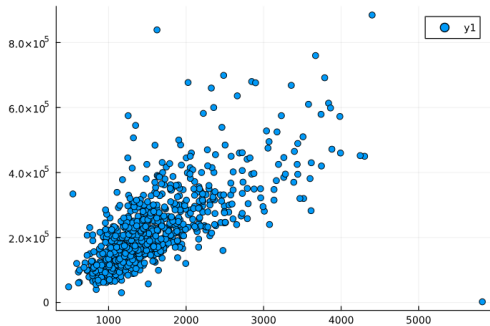


Рис. 11: Цены на недвижимость в зависимости от площади (исключены артефакты данных)

Сначала подключаем необходимые пакеты и формируем данные в нужном виде:

```
[86]:  
# Добавление данных :latitude и :longitude в новый фрейм:  
X = filter_houses[!,:latitude,:longitude]
```

[86]:  
814×2 DataFrame 789 rows omitted

Row	latitude	longitude
	Float64	Float64
1	38.6319	-121.435
2	38.4789	-121.431

Рис. 12: Добавление данных

```
[92]:  
# Конвертация данных в матричный вид:  
X = Matrix(X)
```

```
[92]:  
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  -121.328
```

Рис. 13: Конвертация данных в матричный вид

В качестве критерия для формирования кластеров данных и определения количества кластеров попробуем использовать количество почтовых индексов:

```
[96]:  
# Задание количества кластеров:  
k = length(unique(filter_houses[,zip]))  
[96]:  
66
```

Рис. 14: Задание количества кластеров

Далее сформируем новый фрейм, включающий исходные данные о недвижимости и столбец с данными о назначенном каждому дому кластере

```
[100]:  
  
# Формирование фрейма данных:  
df = DataFrame(cluster = C.assignments,city = filter_houses[:,city],  
               latitude = filter_houses[:,latitude],longitude = filter_houses[:,longitude],zip  
               )
```

[100]:  
814x5 DataFrame 789 rows omitted

Row	cluster	city	latitude	longitude	zip
	Int64	String15	Float64	Float64	Int64
1	28	SACRAMENTO	38.6319	-121.435	95838
2	34	SACRAMENTO	38.4789	-121.431	95823

Рис. 15: Формирование фрейма данных



Плстроила график, обозначив каждый кластер отдельным цветом

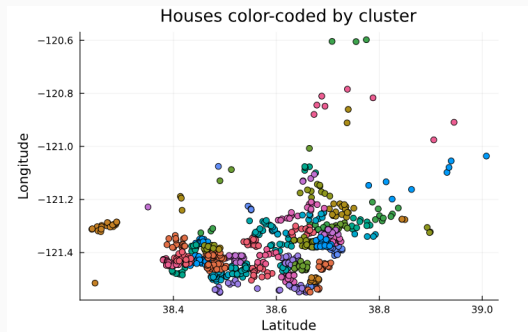


Рис. 16: График кластеров

Построила график, раскрасив кластеры по почтовому индексу

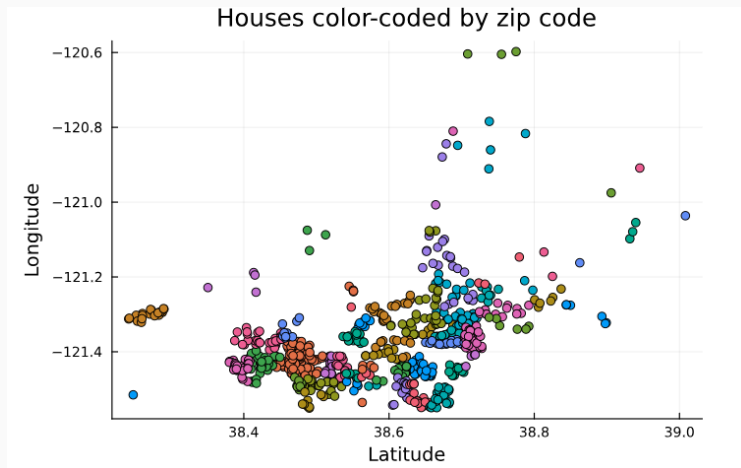


Рис. 17: График кластеров 2

8. Кластеризация данных. Метод k ближайших соседей. Данный метод заключается в отнесении объекта к тому из известных классов, который является наиболее распространённым среди  $k$  соседей данного элемента.

В случае использования метода для регрессии, объекту присваивается среднее значение по ? ближайшим к нему объектам.

```
[112]:
knearest = 10
id = 70
point = X[:,id]

[112]:
2-element Vector{Float64}:
 38.44004
-121.421012

[114]:
# Поиск ближайших соседей:
kdtree = KDTree(X)
idxs, dists = knn(kdtree, point, knearest, true)

[114]:
([70, 764, 196, 125, 557, 368, 415, 92, 112, 683], [0.0, 0.006264891539364138, 0.00825320259050462, 0.008473585132630057, 0.009164073548370188, 0.009405065124697706, 0.009921759722950759, 0.009941028618812013, 0.010332637707777167, 0.011168993911721985])
```

Рис. 18: Метод главных компонент

[116]:

```
# 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)
```

[116]:

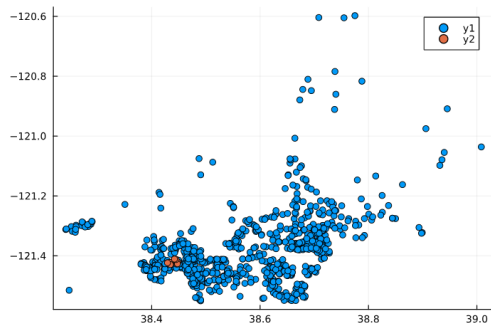


Рис. 19: Определение соседей объекта недвижимости

На примере с данными о недвижимости попробуем уменьшить размеры данных о цене и площади из набора данных домов

```
[126]:  
  
# Фрейм с указанием площади и цены недвижимости:  
F = filter_houses[I, [:sq__ft,:price]]  
# Конвертация данных в массив:  
F = Array(F)  
  
[126]:  
  
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
```

Рис. 20: Создание фрейма

[130]:

```
# Приведение типов данных к распределению для PCA:  
M = fit(PCA, F)
```

[130]:

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

Pattern matrix (unstandardized loadings):

	PC1
1	460.52
2	1.19826e5

Importance of components:

	PC1
SS Loadings (Eigenvalues)	1.43584e10
Variance explained	0.999984
Cumulative variance	0.999984
Proportion explained	1.0
Cumulative proportion	1.0

Рис. 21: Приведение типов данных к распределению для PCA

## Выполнение лабораторной работы

[140]:

```
# 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,:])
```

[140]:

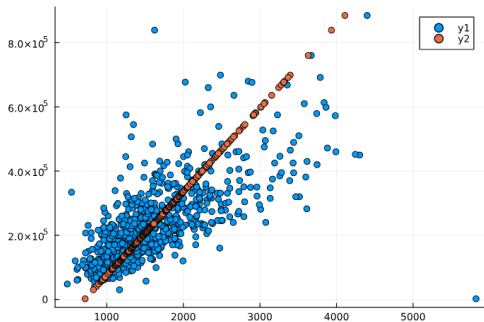


Рис. 22: Определение главных компонент для данных по объектам недвижимости



## Выполнение лабораторной работы

Зададим случайный набор данных (можно использовать и полученные экспериментальным путём какие-то данные). Попробуем найти для данных лучшее соответствие

[142]:

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

[142]:

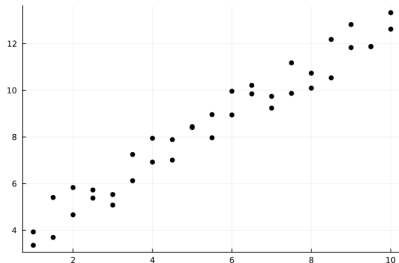


Рис. 23: Исходные данные

Применим функцию линейной регрессии для построения соответствующего графика значений

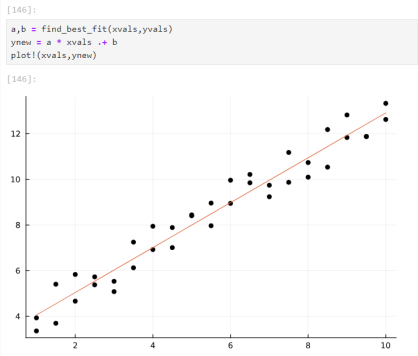


Рис. 24: Линейная регрессия

11. Перешла к заданиям для самостоятельного выполнения. Нумерация соответствует.

## Задания для самостоятельной работы

- Задание на Кластеризацию.

```
[162]:  
  
# Загрузка данных iris  
iris = dataset("datasets", "iris")  
  
# Извлечение числовых данных (первые четыре столбца) и преобразование в массив  
data = Matrix(iris[, 1:4])  
# Транспонируем массив для использования в k-средних (объекты – столбцы)  
  
# Кластеризация методом k-средних на 3 кластера (так как в данных есть 3 вида ирисов)  
k = 3  
result = kmeans(data, k)  
  
# Получение меток кластеров  
labels = result$assignments  
  
# Построение точечной диаграммы (используем первые два признака для визуализации)  
scatter(  
  data[1, ], data[2, ],  
  group = labels,  
  title = "K-means Clustering of Iris Dataset",  
  xlabel = "Feature 1: Sepal Length",  
  ylabel = "Feature 2: Sepal Width",  
  legend = :topright,  
  markersize = 5  
)
```

Рис. 25: Задание 1. Код

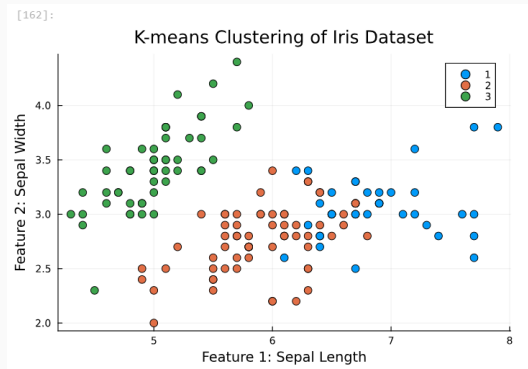


Рис. 26: Задание 1. График

# Задания для самостоятельной работы

- Задание на регрессию (часть 1)

```
[172]: using LinearAlgebra
using DataFrames
using GLM

# Генерация данных
N = 1000 # Количество наблюдений
X = randn(N, 3) # Матрица факторов (N x 3)
a0 = rand(3) # Истинные коэффициенты
y = X * a0 + 0.1 * randn(N) # Линейная зависимость с добавлением шума

# Добавляем столбец единиц для свободного члена
X2 = hcat(ones(N), X)

# Проверка матрицы  $X^T X$ 
XtX = X2' * X2
if !isposdef(XtX)
    println("Матрица  $X^T X$  не положительно определена. Добавляем регуляризацию.")
end

# Регуляризация вручную: Ridge Regression
λ = 1e-6 # Маленькое значение регуляризации
XtX_reg = XtX + λ * I # Добавляем λI к  $X^T X$ 
beta_ridge = XtX_reg \ (X2' * y) # Решение системы
println("Оценка коэффициентов (с регуляризацией): ", beta_ridge)

# Сравнение с GLM.jl
df = DataFrame(hcat(X, y), [:X1, :X2, :X3, :Y]) # Создаем DataFrame
model = lm(@formula(Y ~ X1 + X2 + X3), df) # Модель линейной регрессии
println("Оценка с использованием GLM.jl: ", coef(model))

# Проверка схождения коэффициентов
println("Разница между оценками (GLM - Ridge): ", coef(model) - beta_ridge)
```

Оценка коэффициентов (с регуляризацией): [-0.004361583263215564, 0.6433713344418354, 0.6787211781696519, 0.6094844818961423]  
Оценка с использованием GLM.jl: [-0.0043615833086847935, 0.643371335071556, 0.6787211788728496, 0.6094844825440442]  
Разница между оценками (GLM - Ridge): [-4.546922907633366e-11, 6.297206089911356e-10, 7.031976112514826e-10, 6.479018432870021e-10]

Рис. 28: Задание 2.1. Решение

- Задание на регрессию (часть 2)

```
[188]: # Часть 2: Линейная регрессия для одномерного случая
# Создание данных
X = rand(100)
y = 2X + 0.1 * randn(100)

# Добавляем столбец единиц
X2 = hcat(ones(length(X)), X)

# МНК-оценка
beta_mnk = (X2' * X2) \ (X2' * y)
println("МНК-оценка коэффициентов: ", beta_mnk)

# Построение графика
scatter(X, y, label="Данные", title="График регрессии", xlabel="X", ylabel="Y")

МНК-оценка коэффициентов: [0.006898153545410991, 1.9878127109267747]
```

Рис. 29: Задание 2.2. Код



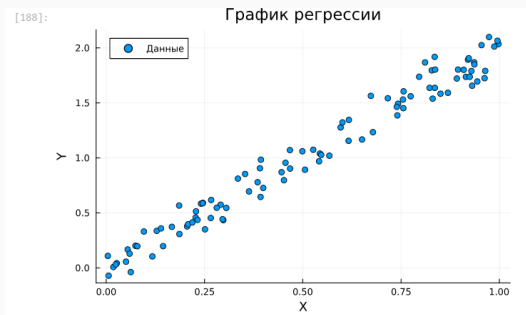


Рис. 30: Задание 2.2. График

- Модель ценообразования биномиальных опционов (а)

```
[192]: using Plots
       using Random

# Параметры
S = 100.0      # Начальная цена акции
T = 1.0        # Длина биномиального дерева в годах
n = 10000      # Количество периодов
σ = 0.3        # Волатильность
r = 0.08       # Годовая процентная ставка

# Вычисления
h = T / n      # Длина одного периода
u = exp(r * h + σ * sqrt(h)) # Рост
d = exp(r * h - σ * sqrt(h)) # Падение
p_star = (exp(r * h) - d) / (u - d) # Вероятность роста

# Построение одной траектории
function generate_path(S, u, d, p_star, n)
    path = [S] # Начальная цена
    for _ in 1:n
        rand() < p_star ? push!(path, path[end] * u) : push!(path, path[end] * d)
    end
    return path
end

# Генерация и построение траектории
path = generate_path(S, u, d, p_star, n)
plot(0:n, path, xlabel="Периоды", ylabel="Цена акции",
     title="Траектория курса акций")
```

Рис. 31: Задание 3.1. Код

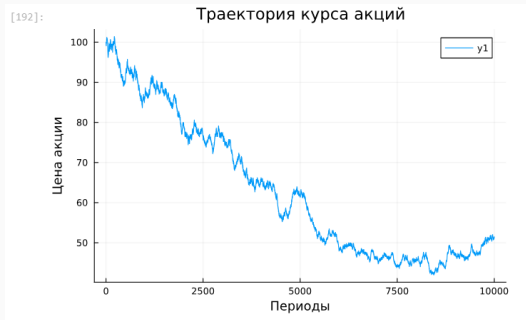


Рис. 32: Задание 3.1. График

- Модель ценообразования биномиальных опционов (b)

```
•[214]: # Функция для создания траектории
function createPath(S::Float64, r::Float64, σ::Float64, T::Float64, n::Int64)
    h = T / n
    u = exp(r * h + σ * sqrt(h))
    d = exp(r * h - σ * sqrt(h))
    p_star = (exp(r * h) - d) / (u - d)
    return generate_path(S, u, d, p_star, n)
end

# Генерация 10 траекторий
plot(title="10 траекторий курса акций", xlabel="Периоды", ylabel="Цена акции")
for _ in 1:10
    path = createPath(S, r, σ, T, n)
    plot!(0:n, path, legend=false)
end
```

Рис. 33: Задание 3.2. Код

# Задания для самостоятельной работы

- Модель ценообразования биномиальных опционов (с)

```
using Base.Threads
# Функция для генерации траектории
function createPath(S, r, σ, T, n)
    dt = T / n
    path = zeros(Float64, n+1)
    path[1] = S
    for i in 2:n+1
        path[i] = path[i-1] * exp((r - 0.5 * σ^2) * dt + σ * sqrt(dt) * randn())
    end
    return path
end

# Параметры для генерации траекторий
S = 100.0 # Начальная цена акции
r = 0.05 # Безрисковая ставка
σ = 0.2 # Волатильность
T = 1.0 # Время до экспирации (в годах)
n = 252 # Количество шагов (дней в году)

# Многопоточная генерация траекторий
n_trajectories = 10
trajectories = Vector{Vector{Float64}}(undef, n_trajectories)

# Параллельный расчет траекторий
Threads.@threads for i in 1:n_trajectories
    trajectories[i] = createPath(S, r, σ, T, n)
end

# Построение графика всех траекторий
p = plot(title="Распараллеленные траектории курса акций", xlabel="Периоды", ylabel="Цена акции")

# Добавление траекторий на график
for path in trajectories
    plot!(p, 0:n, path, legend=false)
end

# Явный вывод графика
display(p)
```

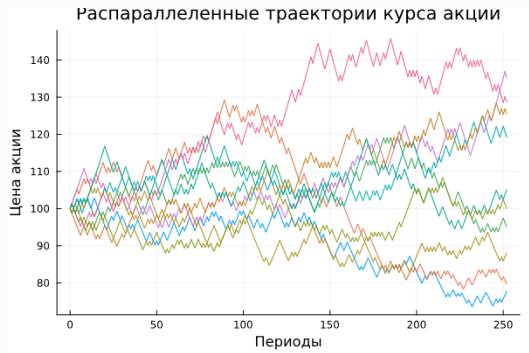


Рис. 35: Задание 3.3. График

## Задания для самостоятельной работы

- Модель ценообразования биномиальных опционов (d)

```
[224]: # Параметры
S = 100.0 # Начальная цена акции
T = 1.0 # Длина биномиального дерева в годах
n = 10000 # Количество периодов
sigma = 0.3 # Волатильность
r = 0.08 # Годовая процентная ставка

# Вычисления
h = T / n # Длина одного периода
u = exp(r * h + sigma * sqrt(h)) # Рост
d = exp(r * h - sigma * sqrt(h)) # Падение
p_star = (exp(r * h) - d) / (u - d) # Вероятность роста
path = createPath(S, r, sigma, T, n)
plot(0:n, path, xlabel="Периоды", ylabel="Цена акции", title="Одна траектория курса акций")
```

[224]:



## Вывод

---



Изучила специализированные пакеты Julia для обработки данных.