

Лабораторная работа № 7

Введение в работу сданными

Беличева Д. М.

Российский университет дружбы народов, Москва, Россия

Информация

- Беличева Дарья Михайловна
- студентка
- Российский университет дружбы народов
- 1032216453@pfur.ru
- <https://dmbelicheva.github.io/ru/>



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

1. Используя JupyterLab, повторите примеры. При этом дополните графики обозначениями осей координат, легендой с названиями траекторий, названиями графиков и т.п.
2. Выполните задания для самостоятельной работы.

Выполнение лабораторной работы

```
[3]: using CSV, DataFrames, DelimitedFiles

Считывание данных

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

[4]: 73x2 DataFrame                                     48 rows omitted
      Row  year  language
      ----  ---  -
      Int64  String31
1      1   1951  Regional Assembly Language
2      2   1952  Autocode
3      3   1954  IPL
4      4   1955  FLOW-MATIC
5      5   1957  FORTRAN
6      6   1957  COMTRAN
7      7   1958  LISP
8      8   1958  ALGOL 58

[4]: # 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

[4]: language_created_year (generic function with 1 method)

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

[7]: 1991

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

[8]: 2012

[9]: language_created_year(P,"julia")

MethodError: no method matching getindex{::DataFrame, ::Nothing, ::Int64}
The function 'getindex' exists, but no method is defined for this combination of argument types.
Closest candidates are:
  getindex{::DataFrame, ::Union{Missing, ::Union{Signed, Unsigned}}}
```

Рис. 1: Считывание данных

Запись данных в файл

```
[14]: # Запись данных в CSV-файл:
      CSV.write("programming_languages_data2.csv", P)

[14]: "programming_languages_data2.csv"

[15]: # Пример записи данных в текстовый файл с разделителем ',':
      writedlm("programming_languages_data.txt", Tx, ',')

[16]: # Пример записи данных в текстовый файл с разделителем '-':
      writedlm("programming_languages_data2.txt", Tx, '-')

[17]: # Построчное считывание данных с указанием разделителя:
      P_new_delim = readdlm("programming_languages_data2.txt", '-')

[17]: 74x2 Matrix{Any}:
      "year"      "language"
1951      "Regional Assembly Language"
1952      "Autocode"
1954      "IPL"
1955      "FLOW-MATIC"
1957      "FORTRAN"
1957      "COMTRAN"
1958      "LISP"
1958      "ALGOL 58"
1959      "FACT"
1959      "COBOL"
1959      "RPG"
1962      "APL"
      :
2003      "Scala"
2005      "F#
2006      "PowerShell"
2007      "Clojure"
2009      "Go"
2010      "Rust"
2011      "Dart"
2011      "Kotlin"
2011      "Red"
2011      "Elixir"
2012      "Julia"
2014      "Swift"
```

Рис. 2: Запись данных в файл

Словари

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

[18]: Dict{Any, Any}{}

[19]: # 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

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

[20]: 2-element Vector{String}:
      "Groovy"
      "Scala"
```

Рис. 3: Словари

Выполнение лабораторной работы

```
DataFrames

[11]: # Подготовим наш DataFrame:
      using DataFrames

[5]: # Задаем переменную со структурой DataFrame:
      df = DataFrame{year = P{1}, language = P{1}}

[5]: 73x2 DataFrame
      Row  year  language
      Int64  String31
      ────
      1  1951  Regional Assembly Language
      2  1952  Autocode
      3  1954  IPL
      4  1955  FLOW-MATIC
      5  1957  FORTRAN
      6  1957  COMTRAN
      7  1958  LISP
      8  1958  ALGOL 58

[12]: # Вывод всех значений столбца year:
      df[:, :year]

[12]: 73-element Vector{Int64}:
      1951
      1952
      1954
      1955
      1957
      1957
      1957
      1958
      1958
      1959
      1959
      1959
      1962
      1
      2003
      2005
      2006

[32]: # Получим статистические сведения о файле:
      describe(df)

[32]: 2x7 DataFrame
      Row  variable  mean  min  median  max  missing  eltype
      ────
      1  Symbol  Union...  Any  Union...  Any  Int64  DataType
      1  year  1982.99  1951  1986.0  2014  0  Int64
      2  language  ALGOL 58  dBase-II  0  String31
```

Рис. 4: DataFrames

```
RDatasets

[33]: # Подгружаем пакет RDatasets:
      using RDatasets

[34]: # Задаём структуру данных в виде набора данных:
      iris = dataset("datasets", "iris")

[34]: 150x5 DataFrame                                     12
      Row  Sepal.Length  Sepal.Width  Petal.Length  Petal.Width  Species
      <dbl>         <dbl>         <dbl>         <dbl>         <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

[37]: # Определяем тип переменной:
      typeof(iris)

[37]: DataFrame

[38]: describe(iris)                                     12
[38]: 5x7 DataFrame
      Row  variable  mean  min  median  max  nmissing  eltype
      <dbl>      <str>  <dbl> <dbl> <dbl>  <dbl>  <dbl>  <str>
1  Sepal.Length  5.84333  4.3   5.8   7.9     0  Float64
2  Sepal.Width   3.05733  2.0   3.0   4.4     0  Float64
3  Petal.Length  3.758    1.0  4.35  6.9     0  Float64
4  Petal.Width   1.19933  0.1   1.3   2.5     0  Float64
5  Species       setosa   virginica 0  CategoricalValue(String, UInt8)
```

Рис. 5: RDatasets

Работа с переменными отсутствующего типа (Missing Values)

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

40]: Missing

41]: # Пример операции с переменной отсутствующего типа:
a + 1

41]: missing

43]: # Определение перечня продуктов:
foods = ["apple", "cucumber", "tomato", "banana"]
# Определение калорий:
calories = [missing, 47, 22, 105]
# Определение типа переменной:
typeof(calories)

43]: Vector{Union{Missing, Int64}} (alias for Array{Union{Missing, Int64}, 1})

44]: # Подключен пакет Statistics:
using Statistics
# Определение среднего значения:
mean(calories)

44]: missing

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

45]: 58.0

49]: # Задание средних 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 = outerjoin(dataframe_calories, dataframe_prices, on=:item)

49]: 4x3 DataFrame
      Row  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
```

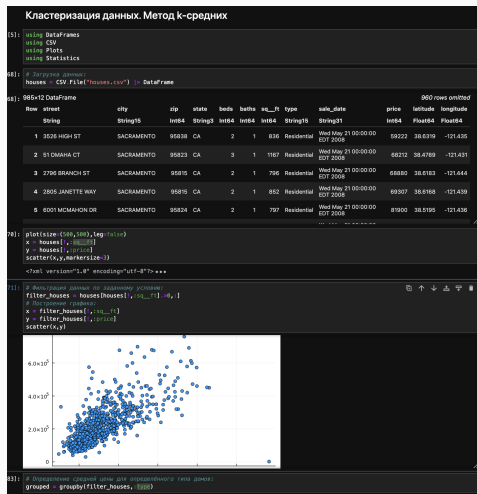
Рис. 6: Работа с переменными отсутствующего типа (Missing Values)

Выполнение лабораторной работы

[illegible]

Рис. 7: FileIO

Выполнение лабораторной работы



Выполнение лабораторной работы

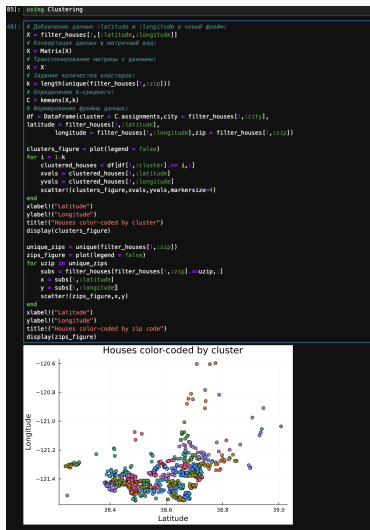


Рис. 9: Кластеризация данных. Метод k-средних

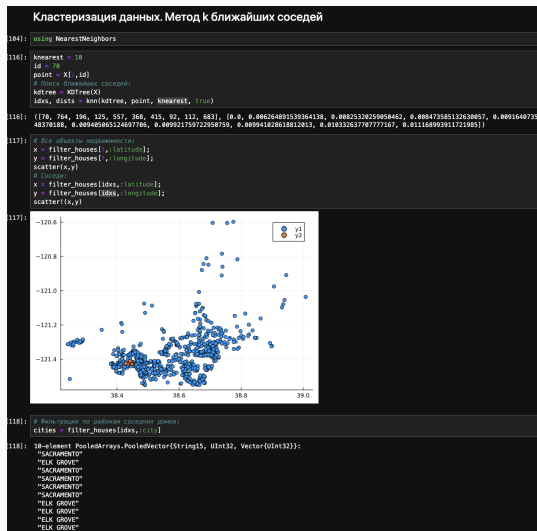


Рис. 10: Кластеризация данных. Метод к ближайших соседей

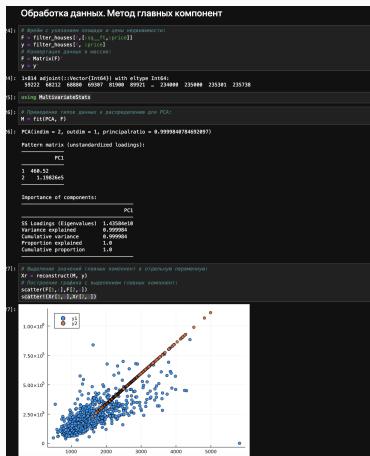


Рис. 11: Обработка данных. Метод главных компонент



Рис. 12: Обработка данных. Линейная регрессия

Выполнение лабораторной работы

```
333]: xvals = 1:100000;
      xvals = repeat(xvals,inner=3);
      yvals = 3 + xvals + 2*rand(length(xvals)) .* 1;
      @show size(xvals)
      @show size(yvals)
      @time a,b = find_best_fit(xvals,yvals)

      size(xvals) = (300000,)
      size(yvals) = (300000,)
      0.016692 seconds (29.73 k allocations: 1.531 MiB, 96.41% compilation time)
333]: (1.00000000292669111, 2.9970284742303193)

334]: using PyCall
      using Conda

337]: 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
      """
      xpy = PyObject(xvals)
      ypy = PyObject(yvals)
      @time a,b = py"find_best_fit_python"(xpy,ypy)

      0.043921 seconds (112.61 k allocations: 5.636 MiB, 92.86% compilation time)
337]: (1.0000000029266911, 2.997028474244871)

339]: using BenchmarkTools

340]: @time a,b = py"find_best_fit_python"(xvals,yvals)
      @time a,b = find_best_fit(xvals,yvals)

      0.011368 seconds (11.77 k allocations: 685.641 KiB, 70.18% compilation time)
      337.375 μs (1 allocation: 32 bytes)
340]: (1.00000000292669111, 2.9970284742303193)
```

Рис. 13: Обработка данных. Линейная регрессия

Выполнение лабораторной работы

[illegible]

Рис. 14: Кластеризация

Выполнение лабораторной работы

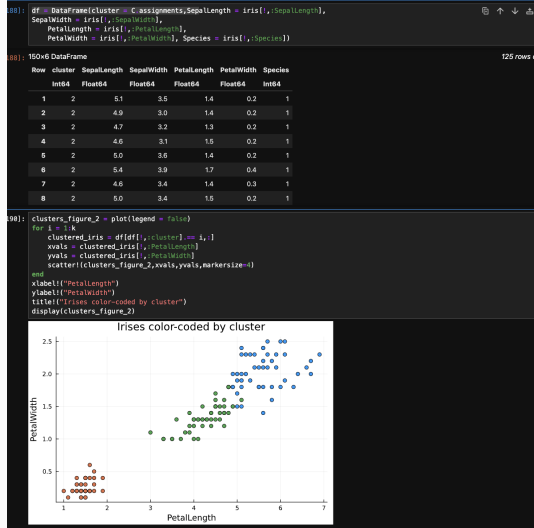


Рис. 15: Кластеризация

Выполнение лабораторной работы

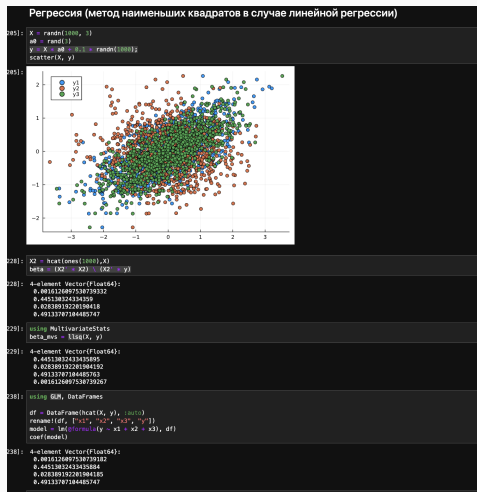


Рис. 16: Регрессия

Выполнение лабораторной работы



Рис. 17: Регрессия

Выполнение лабораторной работы

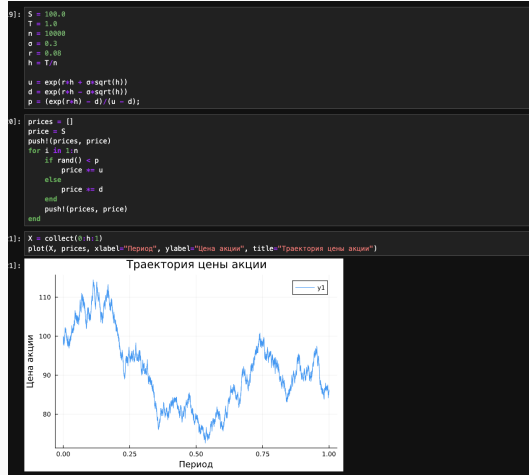


Рис. 18: Модель ценообразования биномиальных опционов

Выполнение лабораторной работы

```
[24]: function createPath(S::Float64, r::Float64, σ::Float64, T::Float64, n::Int64)
      h = T/n
      X = collect{0:h:1}
      u = exp(r+h + σ*sqrt(h))
      d = exp(r+h - σ*sqrt(h))
      p = (exp(r+h) - d)/(u - d);
      prices = []
      price = S
      push!(prices, price)
      for i in 1:n
          if rand() < p
              price *= u
          else
              price *= d
          end
          push!(prices, price)
      end
      return (X, prices)
end
```

```
[24]: createPath (generic function with 1 method)
```

```
[25]: p = plot()
      for i=1:10
          plot!(p, createPath(S, r, σ, T, n))
      end
      p
```

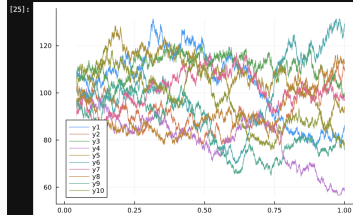


Рис. 19: Модель ценообразования биномиальных опционов

Выполнение лабораторной работы

```
26]: function createPath_parallel(S::Float64, r::Float64, d::Float64, T::Float64, n::Int64)
    h = T/n
    u = exp(r*h + d*sqrt(h))
    d = exp(r*h - d*sqrt(h))
    p = (exp(r*h) - d)/(u - d)
    prices = []
    price = S
    push!(prices, price)
    Threads.@threads for i in 1:n
        if rand() < p
            price *= u
        else
            price *= d
        end
        push!(prices, price)
    end
    return prices
end
```

26]: createPath_parallel (generic function with 1 method)

```
28]: p = plot()
    for i=1:10
        plot!(p, createPath_parallel(S, r, d, T, n))
    end
    p
```

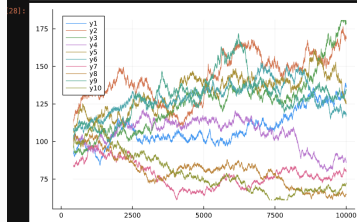


Рис. 20: Модель ценообразования биномиальных опционов

В результате выполнения данной лабораторной работы я освоила специализированные пакеты 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).