

Intro to Julia

By Raj Rao Nadakuditi

An introduction to the Julia programming language. Introduces students to variables, arrays, functions, and everything else that they need to succeed!

© This content is the intellectual property of the respective contributors and may not be reproduced or distributed without permission. Please do not redistribute or host this material publicly.

1. Why use Julia for machine learning and data science?

2. Variables in Julia

2.1. Variable types in Julia

2.2. Converting a variable to a particular type

2.3. Parsing and concatenating strings

2.4. Complex numbers

3. Arrays: Vectors and Matrices

3.1. Vectors

3.2. Matrices

3.3. Slices of an array

3.4. DataFrames

3.5. Reshaping an array

4. Loops and Conditionals

4.1. Comparison

4.2. The Boolean `&&` and `||` conditional operators

4.3. `if-then` conditionals

5. Functions and Methods

5.1. Functions

5.2. Methods

5.3. Populating matrices with array comprehensions

6. Operations on arrays

6.1. `mapslices`: function on slices of an array

6.2. Broadcasting with `.`

Complete the code below to compute the submatrix `df_white` which contains data for all the white wines.

Hint: there are multiple ways to do this. You can explicitly take all the white wines, or compute the set of all wines that are not red...

```
1 df_white = df[df.Color .== "White" ,:]
```

✓ 2s

493 rows × 6 columns

| | fixedacidity | volatileacidity | freesulfurdioxide | density | alcohol | Color |
|----|--------------|-----------------|-------------------|---------|---------|--------|
| | Float64 | Float64 | Float64 | Float64 | Float64 | String |
| 1 | 6.6 | 0.425 | 23.0 | 0.99082 | 11.4 | White |
| 2 | 7.2 | 0.25 | 51.0 | 0.9964 | 9.2 | White |
| 3 | 6.9 | 0.25 | 28.0 | 0.99088 | 11.7 | White |
| 4 | 7.7 | 0.275 | 19.0 | 0.992 | 10.7 | White |
| 5 | 7.1 | 0.26 | 31.0 | 0.99644 | 11.2 | White |
| 6 | 6.0 | 0.24 | 34.0 | 0.9946 | 10.4 | White |
| 7 | 6.9 | 0.25 | 36.0 | 0.9948 | 10.7 | White |
| 8 | 7.5 | 0.17 | 65.0 | 0.997 | 10.0 | White |
| 9 | 5.1 | 0.26 | 26.0 | 0.99449 | 9.2 | White |
| 10 | 7.0 | 0.14 | 10.0 | 0.99352 | 9.9 | White |
| 11 | 6.4 | 0.31 | 12.0 | 0.9919 | 10.4 | White |
| 12 | 7.1 | 0.32 | 52.0 | 0.998 | 8.8 | White |
| 13 | 6.5 | 0.25 | 29.0 | 0.99776 | 10.1 | White |
| 14 | 5.8 | 0.12 | 35.0 | 0.9908 | 11.4 | White |
| 15 | 7.4 | 0.19 | 33.0 | 0.993 | 9.6 | White |
| 16 | 6.2 | 0.25 | 58.0 | 0.99454 | 10.4 | White |
| 17 | 5.9 | 0.27 | 43.0 | 0.9941 | 10.7 | White |
| 18 | 6.1 | 0.27 | 65.0 | 0.9957 | 9.0 | White |
| 19 | 8.0 | 0.4 | 27.0 | 0.9935 | 12.2 | White |
| 20 | 6.6 | 0.2 | 35.0 | 0.99396 | 9.4 | White |
| 21 | 5.1 | 0.14 | 15.0 | 0.9919 | 9.2 | White |
| 22 | 6.7 | 0.26 | 40.0 | 0.99479 | 10.4 | White |



| | fixedacidity | volatileacidity | freesulfurdioxide | density | alcohol | Color |
|----|--------------|-----------------|-------------------|---------|---------|--------|
| | Float64 | Float64 | Float64 | Float64 | Float64 | String |
| 23 | 6.4 | 0.17 | 33.0 | 0.99152 | 10.4 | White |
| 24 | 6.8 | 0.2 | 38.0 | 0.993 | 9.1 | White |
| 25 | 5.8 | 0.17 | 11.0 | 0.99202 | 10.4 | White |
| 26 | 7.1 | 0.27 | 26.0 | 0.99335 | 11.5 | White |
| 27 | 7.5 | 0.14 | 50.0 | 0.9945 | 9.6 | White |
| 28 | 6.5 | 0.19 | 23.0 | 0.9937 | 10.0 | White |
| 29 | 7.5 | 0.26 | 33.0 | 1.0011 | 8.8 | White |
| 30 | 6.2 | 0.35 | 33.0 | 0.99908 | 8.8 | White |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |



```
1 @assert all(df_white.Color .== "White")
2 println("White wines: $(size(df_white, 1))")
```

✓ 8ms

White wines: 493

7. Plotting

7.1. Multiple plots on the same axes

7.2. Multiple plots alongside each other using the `layout` option

8. The `@manipulate` macro

9. Additional Exercises

9.1. Exploratory visualization of a temperature dataset

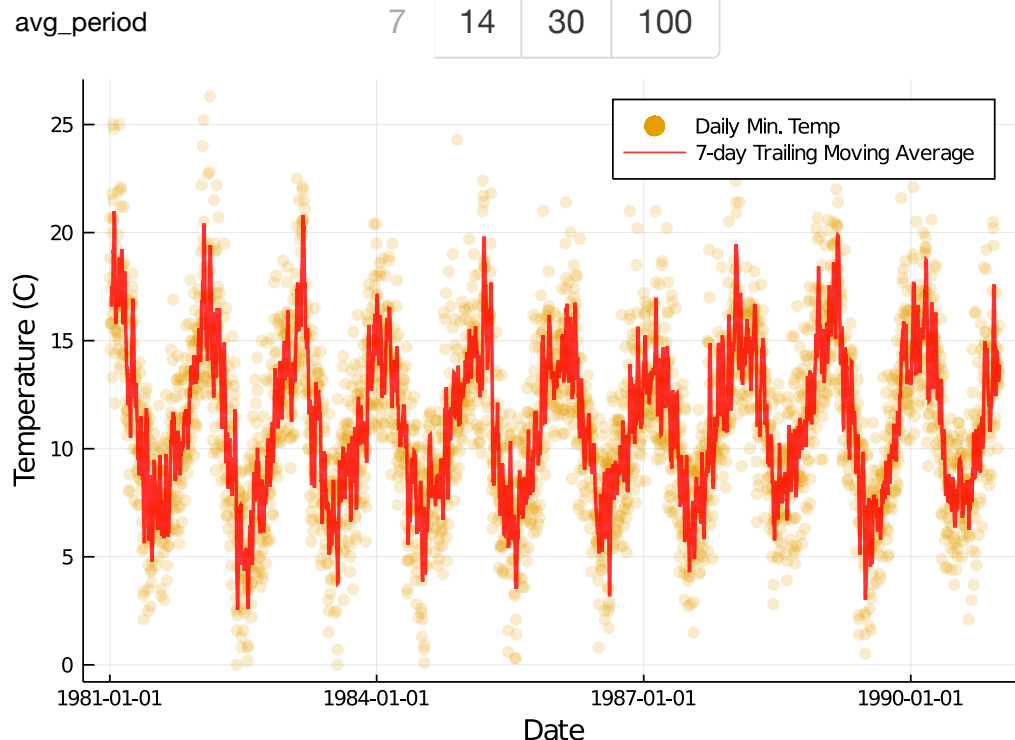
Complete the ?? below to create an interactive plot of temperature and moving average data.

Hint: recall that `df.col` may be used to index a column named `col` in a `DataFrame` called `df`.

```

1 using Statistics: mean
2
3 "Compute the `n`-step moving average for a vector `v`."
4 moving_average(v::AbstractVector, n::Integer) = [mean(v[(i - n +
5 1):i]) for i in 1:length(v)]
6
7 @manipulate for avg_period in (7, 14, 30, 100)
8     # scatterplot of raw temperature data vs date
9     scatter(
10         min_temp.Date, min_temp.Temp;
11         alpha=0.2, # partial transparency helps when there are
12         many overlapping points
13         xlabel="Date",
14         ylabel="Temperature (C)",
15         label="Daily Min. Temp"
16     )
17
18     # compute and plot moving average
19     moving_avg = moving_average(min_temp.Temp, avg_period)
20     plot!(
21         min_temp.Date[avg_period:end], moving_avg;
22         linewidth=2,
23         color=:red,
24         label="$(avg_period)-day Trailing Moving Average"
25     )
26 end

```



✓ 2s

Comment on the characteristics of the plot.

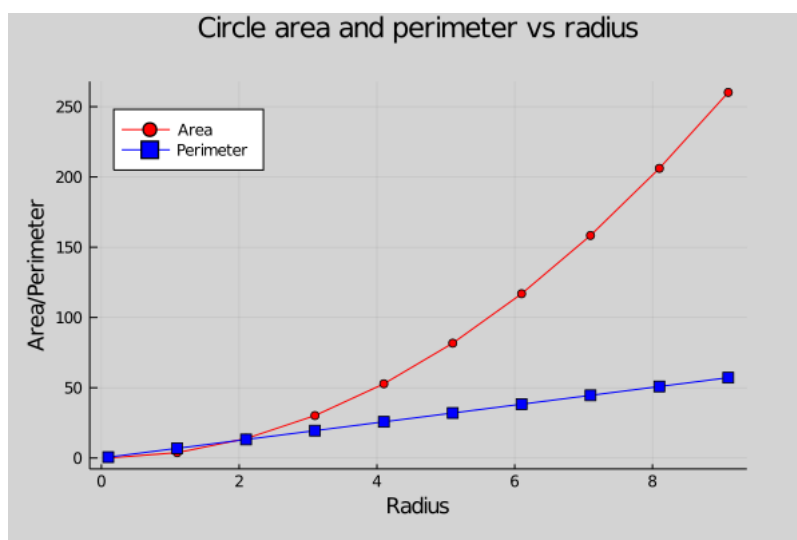
- Do the data appear to be periodic? If so, what is the approximate length of the period?
- What are the pros and cons of looking at the moving average rather than at the raw data?
- Why does the moving average "lag behind" the raw data, especially when `avg_period` is large?



1. Yes, the data appears to be periodic, the approximate length of period is 12 months (one year).
2. The one of advantages of using the moving average is that it can smooth out short-term fluctuations and highlight longer-term trends. On the other hand, it is slower to respond to rapid data, because it gives too much weight to old data.
3. Raw data leads and moving average follows because moverage is based on the past data, especially when the `avg_period` is large, even though the temperature is changed rapidly, the average temperature would not reflect the most recent trends.

9.2. Plotting multiple series

Complete the code below to compute the area and perimeter for a range of circle radius values, then plot the area and perimeter data on a single plot. Consider the following plot as a reference:

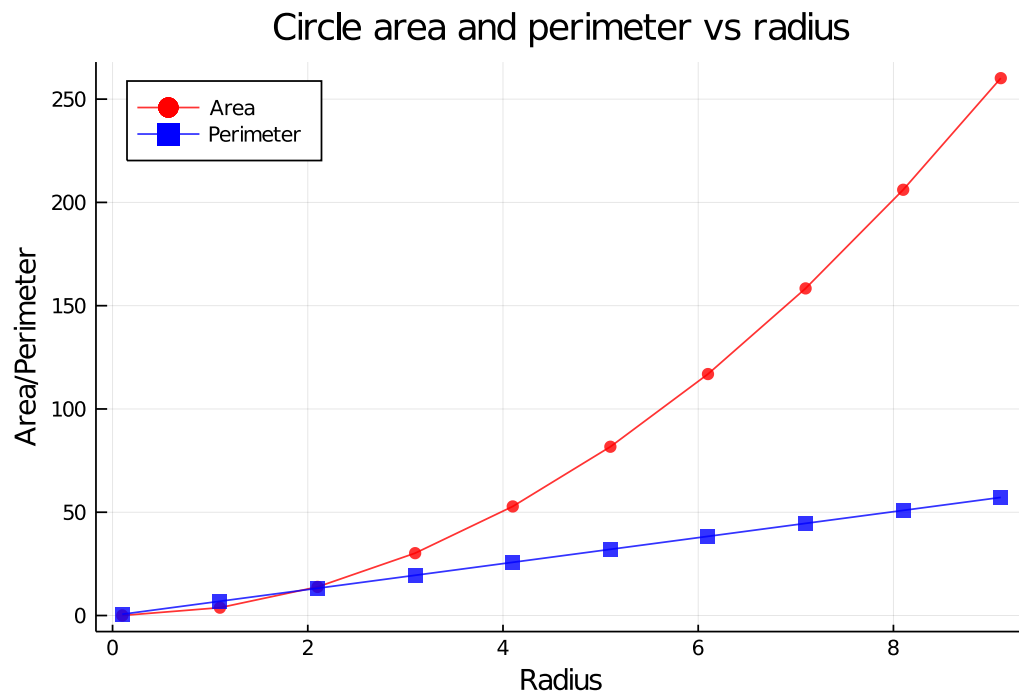


```
1 rr = 0.1:10
2
3 aa = [] # area vector
4 pp = [] # perimeter vector
5
6 for r in rr
7     a, p = circle_area_perimeter(r)
8     push!(aa, a)
9     push!(pp, p)
10 end
11
12 area_perimeter_vs_radius = plot(
13     rr, aa;
```



```
14 color=:red,  
15 marker=:circle,  
16 label="Area",  
17 xlabel="Radius",  
18 ylabel="Area/Perimeter",  
19 legend=:topleft,  
20 title="Circle area and perimeter vs radius"  
21 )  
22 plot!(  
23     rr, pp;  
24     color=:blue,  
25     marker=:square,  
26     label="Perimeter"  
27 )
```

✓ 1s




```
1 using PlotCheck
2 @check_plot area_perimeter_vs_radius
```

✓ 2s

PlotCheck Report

Subplot 1

Series 'Area': Checked.

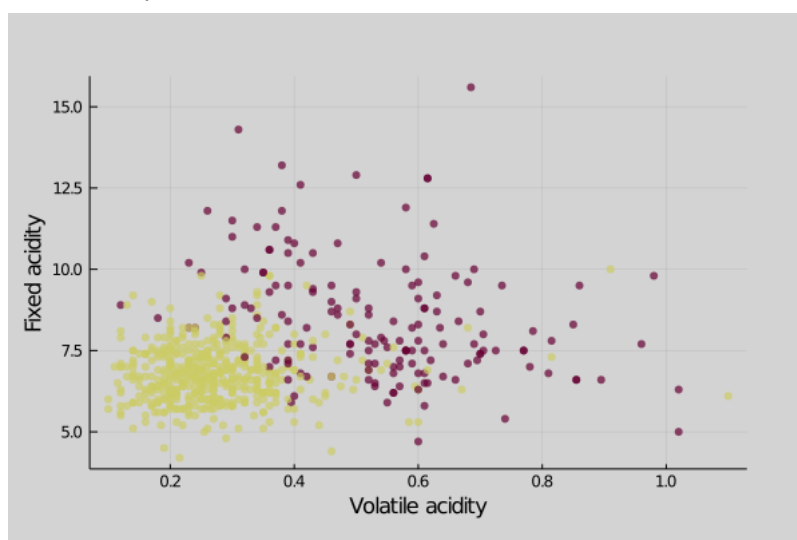
Series 'Perimeter': Checked.

Checked against reference plot.

9.3. A scatter plot with differently-colored points

Create a scatter plot of fixed acidity versus volatile acidity. Volatile acidity values should be on the x-axis, and fixed acidity values should be on the y-axis. Point color should indicate the type of wine: white or red.

When you're finished, the plot should look like this:

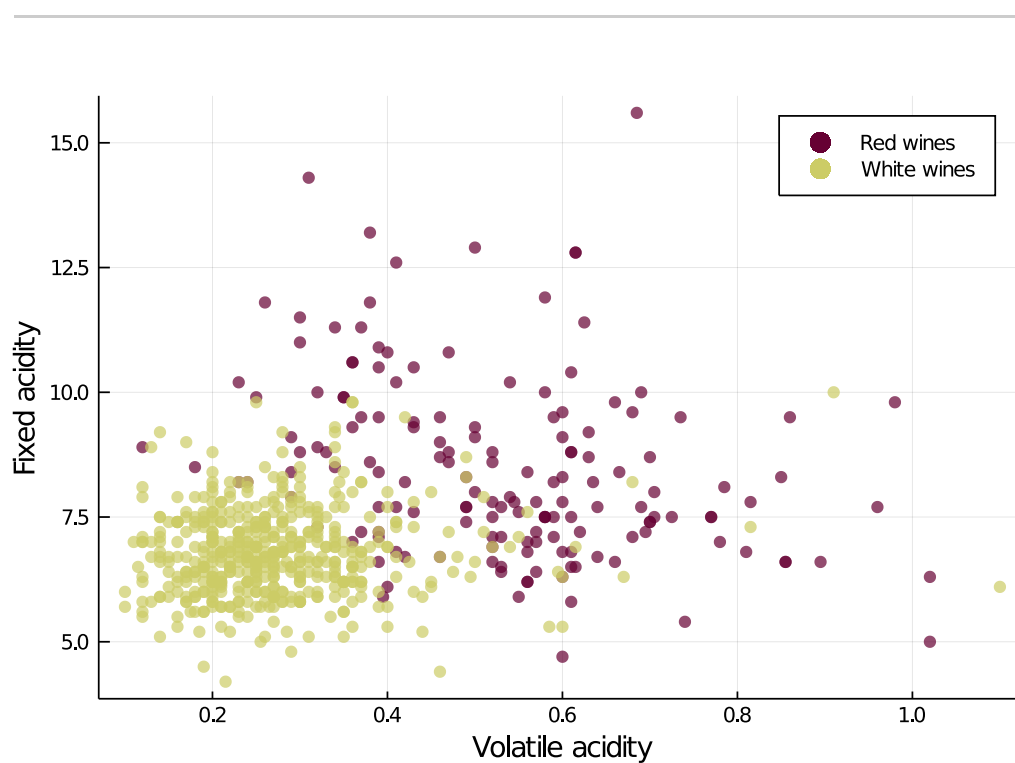


```
1 using Plots, CSV, DataFrames
2
3 default(
4     markerstrokewidth=0.3,
5     markerstrokecolor=:auto,
6     alpha=0.8,
7     label=""
8 )
9
10 wine = DataFrame!(CSV.File("wine.csv"))
11
12 # Select only the red and white wines respectively
13 wine_red = wine[wine.Color .== "Red" ,:]
14 wine_white = wine[wine.Color .== "White" ,:]
15
16 # Define appropriate colors to indicate white and red
```



```
17 c_white = RGB(0.8, 0.8, 0.4)
18 c_red = RGB(0.4, 0.0, 0.2)
19
20 # Plot the red wines
21 fixedacidity_vs_volatileacidity = scatter(
22     wine_red.volatileacidity, wine_red.fixedacidity;
23     color=c_red,
24     alpha=0.7,
25     xlabel="Volatile acidity",
26     ylabel="Fixed acidity",
27     label="Red wines"
28 )
29
30 # Add the white wines
31 scatter!(
32     wine_white.volatileacidity, wine_white.fixedacidity;
33     color=c_white,
34     alpha=0.7,
35     label="White wines"
36 )
```

✓ 28s



```
1 using PlotCheck
2 @check_plot fixedacidity_vs_volatileacidity
```

✓ 8s

PlotCheck Report

Subplot 1

Series 'White wines': Checked.

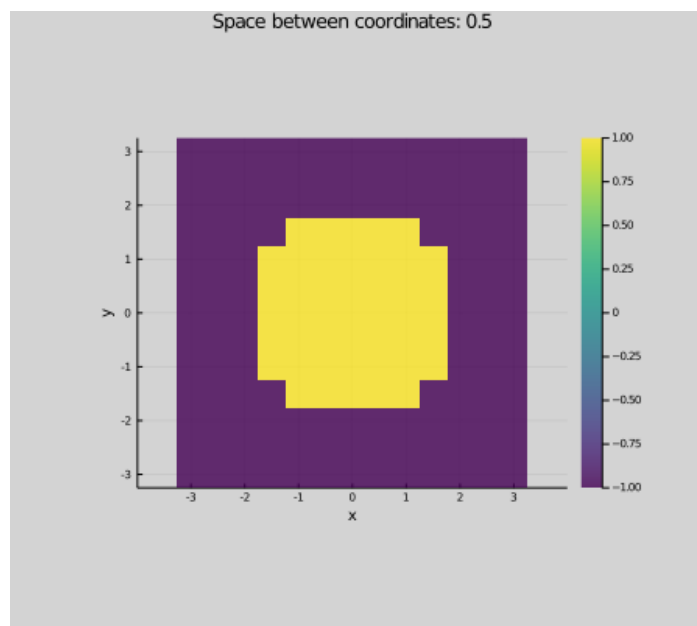
Series 'Red wines': Checked.

Checked against reference plot.



9.4. Visualizing decision boundaries

Complete the code below to visualize a circular decision boundary with radius 2 over a grid of points spaced 0.5 units apart. The finished plot should look like this:



```

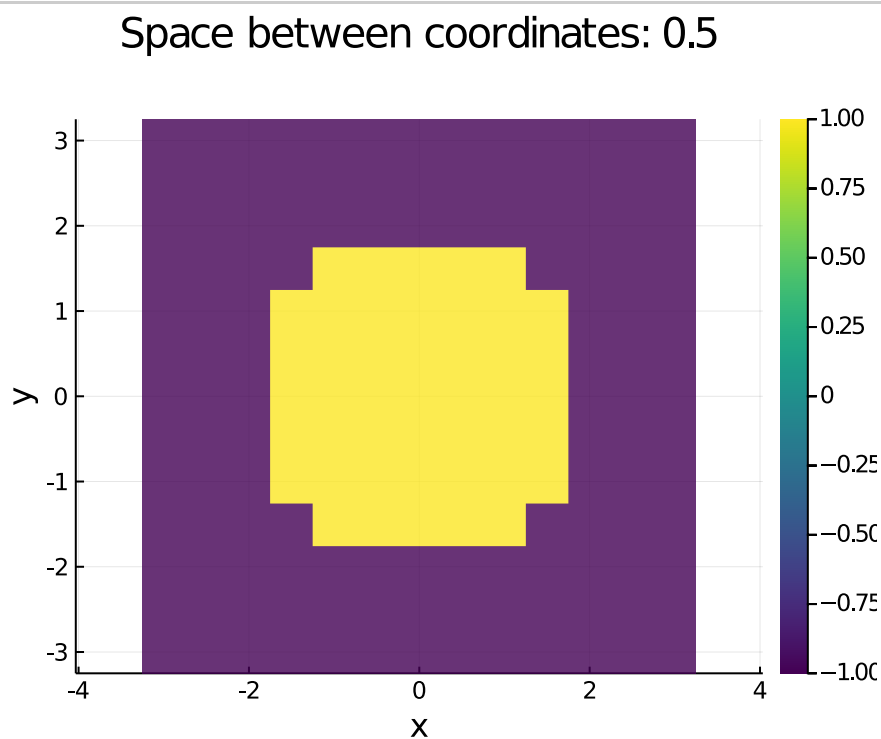
1  r = 2.0
2  s = 0.5
3  x_range = -3:s:3
4  y_range = copy(x_range)
5
6  D = circular_decision_boundary(r, x_range, y_range)
7
8  "Visualize a decision boundary encoded in `D` over a grid of
   points encoded in `x_range` and `y_range`"
9  function heatmap_decision_boundary(x_range::AbstractRange,
   y_range::AbstractRange, D::AbstractMatrix; kwargs...)
10     return heatmap(
11         x_range, y_range, D;
12         aspect_ratio=1.0,
13         size=(450, 400),

```



```
14     xlabel="x",
15     ylabel="y",
16     kwargs...
17 )
18 end
19
20 heatmap_decision_boundary_lowres = heatmap_decision_boundary(
21     x_range,
22     y_range,
23     D;
24     title="Space between coordinates: $(s)"
25 )
```

✓ 1s



```
1 using PlotCheck
2 @check_plot heatmap_decision_boundary_lowres
```

✓ 81ms



PlotCheck Report

Subplot 1

Series "": Checked.

Checked against reference plot.

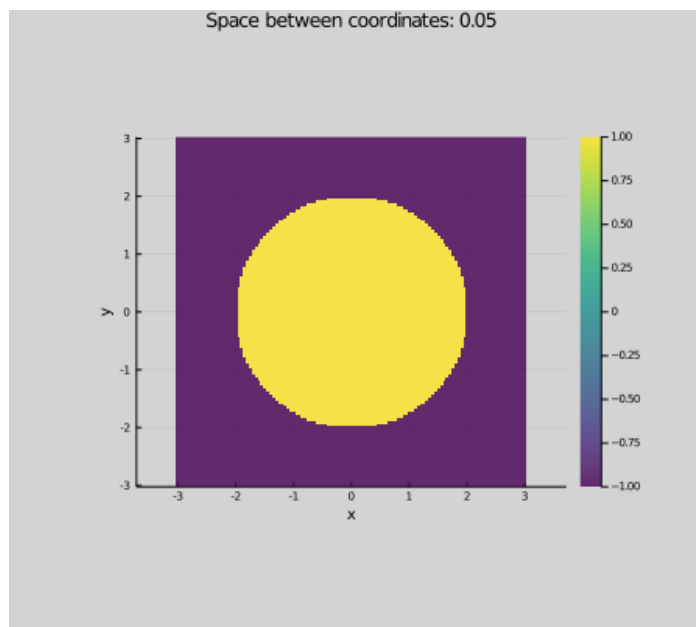


Why didn't we set `yflip=true` in the heatmap command?

Hint: See what happens if you set `yflip=true` and notice what axes is off.

We set `yflip=true` to let the heatmap upside down and it is exactly looking the same.

Now increase the point grid resolution by changing the step from 0.5 to 0.05, corresponding to ten times as many points in both dimensions. Visualize this refined decision boundary with another heatmap. This new heatmap should look like the following:



Finally, show both heatmaps next to each other to make comparison easier.

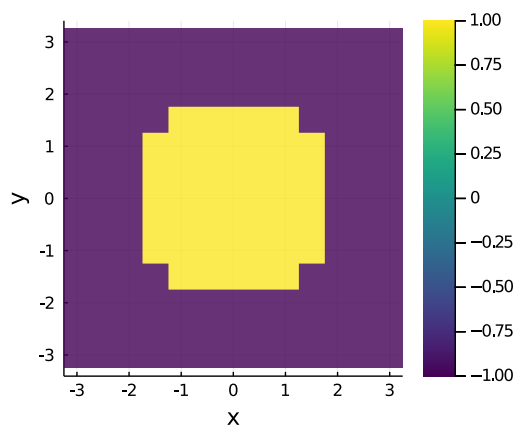
```

1  r = 2.0
2  s = 0.05
3  x_range = -3:0.05:3
4  y_range = copy(x_range)
5
6  D = circular_decision_boundary(r, x_range, y_range)
7  heatmap_decision_boundary_highres = heatmap_decision_boundary(
8      x_range,
9      y_range,
10     D;
11     title="Space between coordinates: $(s)"
12 )
13
14 plot(
15     heatmap_decision_boundary_lowres,
16     heatmap_decision_boundary_highres;
17     layout=(1, 2),
18     size=(750, 400)
19 )

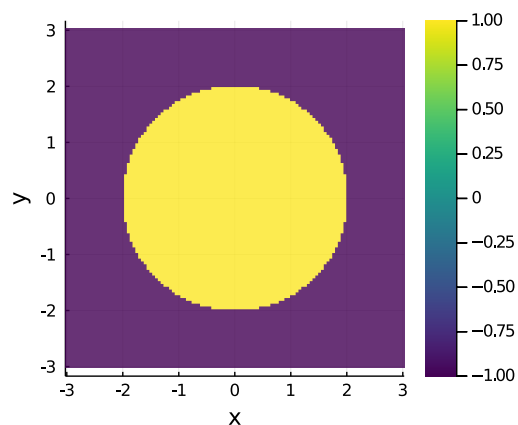
```

✓ 365ms

Space between coordinates: 0.5



Space between coordinates: 0.05



```
1 @check_plot heatmap_decision_boundary_highres
```

✓ 86ms

PlotCheck Report

Subplot 1

Series "": Checked.

Checked against reference plot.



9.5. Programmatically generating a file name

Run the following cell to ensure your function works.



```
1 prefix = "windspeed"
2 suffix = "2015-05-10"
3 index = 5
4 extension = ".tsv"
5
6 generated_filename = make_filename(prefix, suffix, index,
7 extension)
8 target_filename = "windspeed_2015-05-10_5.tsv"
9 isequal(generated_filename, target_filename) && println("Generated
10 file name matches target filename.")
```

✓ 69ms

Generated file name matches target filename.

9.6. Downloading files programmatically with try/ catch



- How many errors were there?
- Which image files were not found in the cloud folder?

There were 2 errors. "Data_file_0.png" and "Data_file_9.png" are not found in the could folder.



```
1 using Images
2
3 img_plots = []
4 for filename in readdir(folder_local)
5     if filename[end - 2:end] == ".png"
6         img = joinpath(folder_local, filename) |> load |> Array
7         push!(img_plots, heatmap(img; axis=false, grid=false,
8 aspect_ratio=1.0, size=(30, 30)))
9     end
10 end
11 plot(img_plots...; layout=(1, length(img_plots)), size=(850, 100))
```

✓ 1s



Completed on 2021-09-03 at 3:09PM
by Yuzhan Jiang

