

Introdução à Ciência de Dados com R

Nelson Quesado
Caio Gustavo

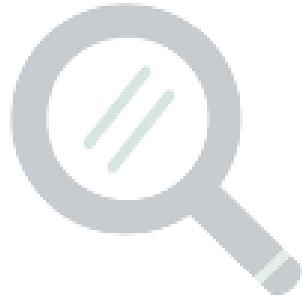


R

**A informação no mundo
DOBRA a cada 20 meses.**

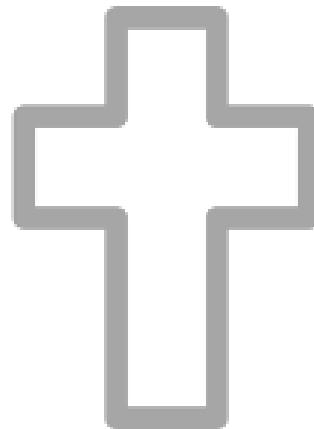
**Thomas Runkle, 2012.
Data Analytics.**

EMPÍRICO



**OBSERVAÇÃO E
INTERAÇÃO**

RELIGIOSO



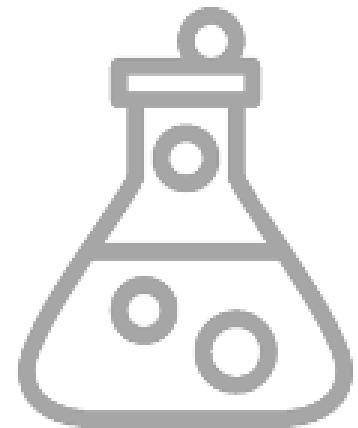
FÉ E DOGMA

FILOSÓFICO



**A NATUREZA DAS
COISAS**

CIENTÍFICO



**LÓGICA
VERIFICÁVEL**

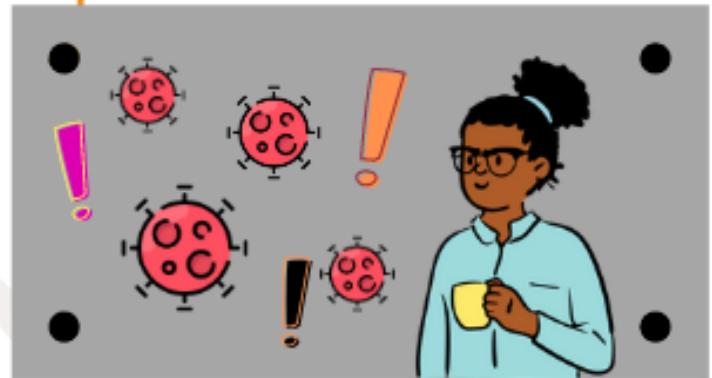
Observação de um fenômeno



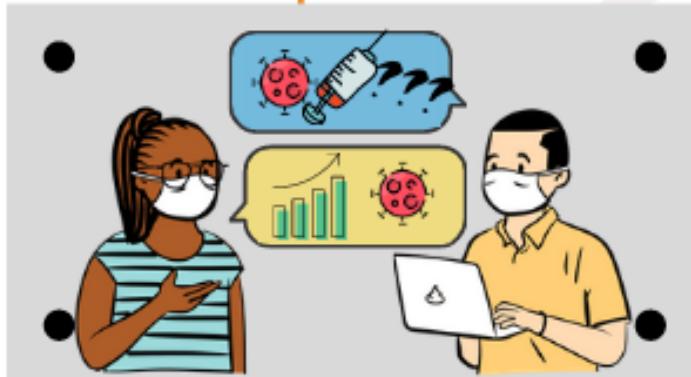
Elaboração de perguntas



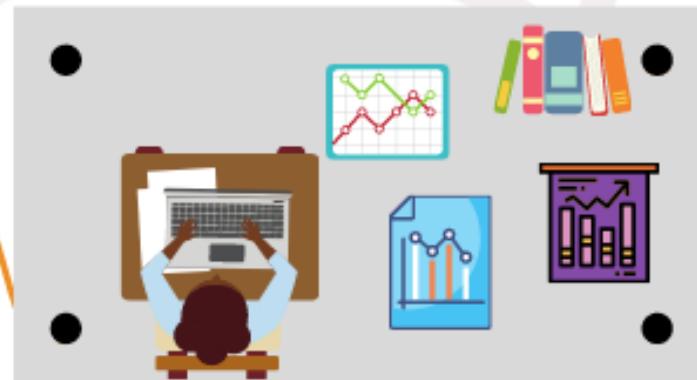
Elaboração de hipóteses



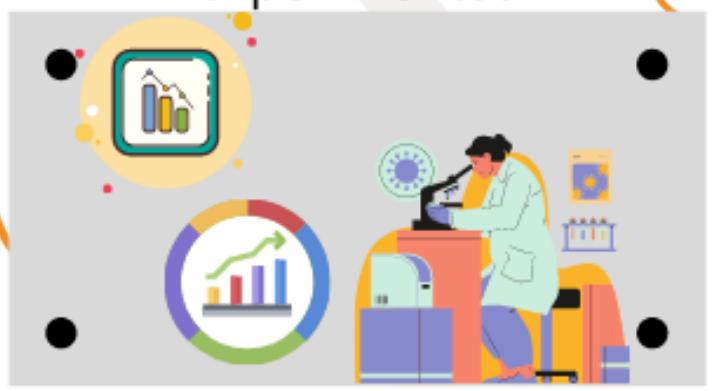
Compartilhar conhecimento



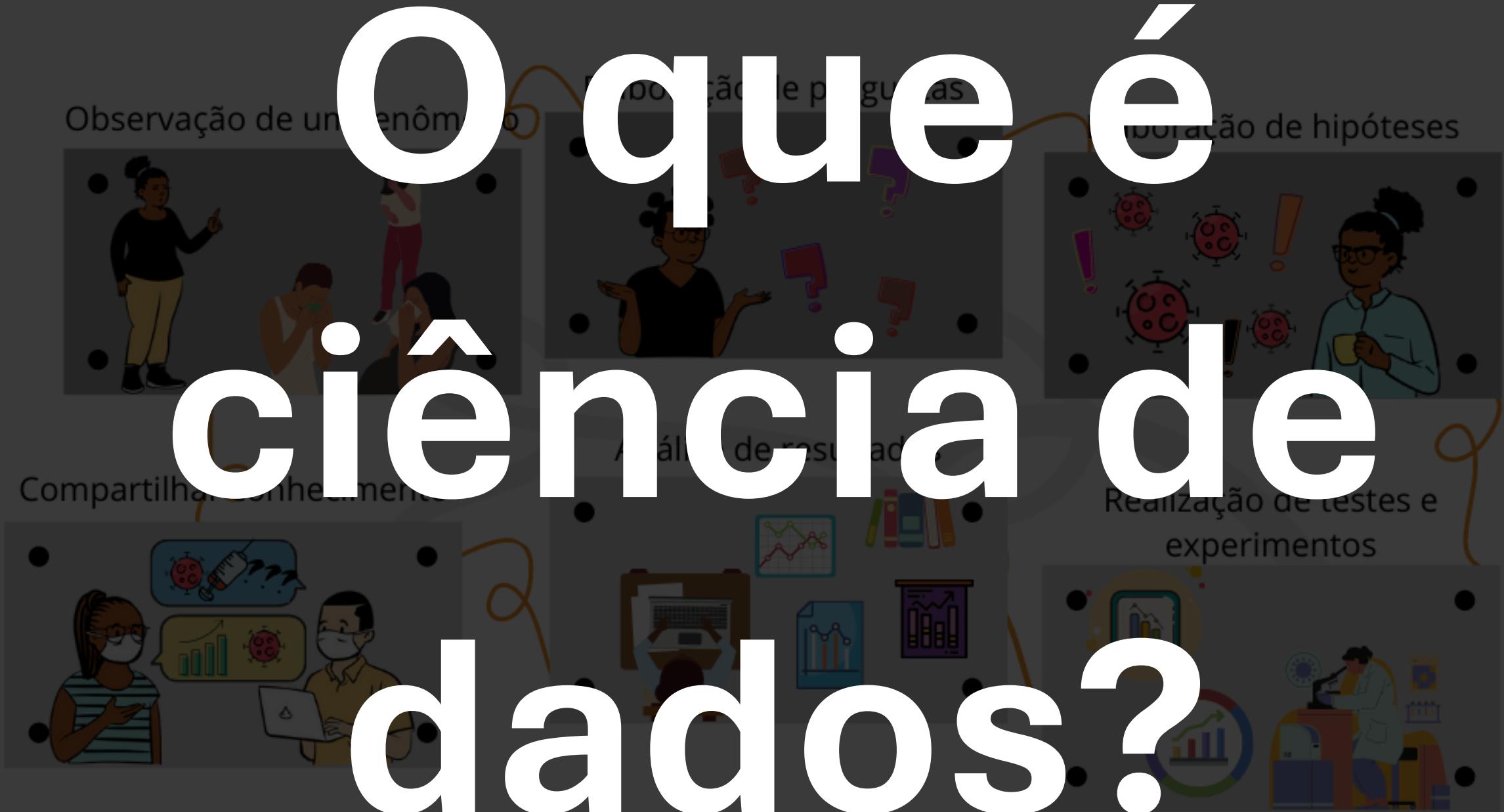
Análise de resultados

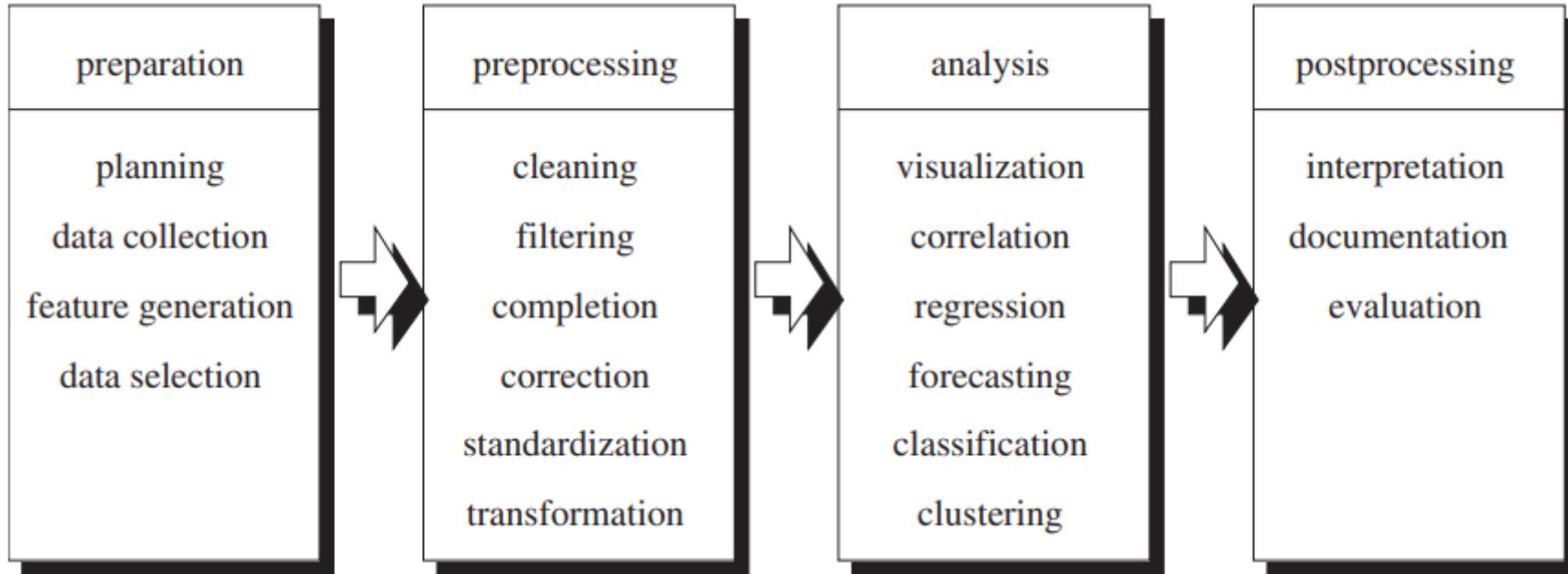


Realização de testes e experimentos



O que é ciência de dados?





Thomas Runkle, 2012.
Data Analytics.



ESTATÍSTICA



ESTATÍSTICA

PROGRAMAÇÃO



ESTATÍSTICA

FENÔMENO

PROGRAMAÇÃO

ESTATÍSTICA

FENÔMENO

DATA
SCIENCE

PROGRAMAÇÃO



Fraud detection



Speech recognition



Internet search



Finance



Predictive maintenance



Genomics



Search engines



Healthcare



Airline route planning



Logistics



Video game



Supply chain optimization



Manufacturing



Sentiment analysis



Recommendation systems



Pattern recognition



E-commerce



Exploratory data analysis



Targeted advertising



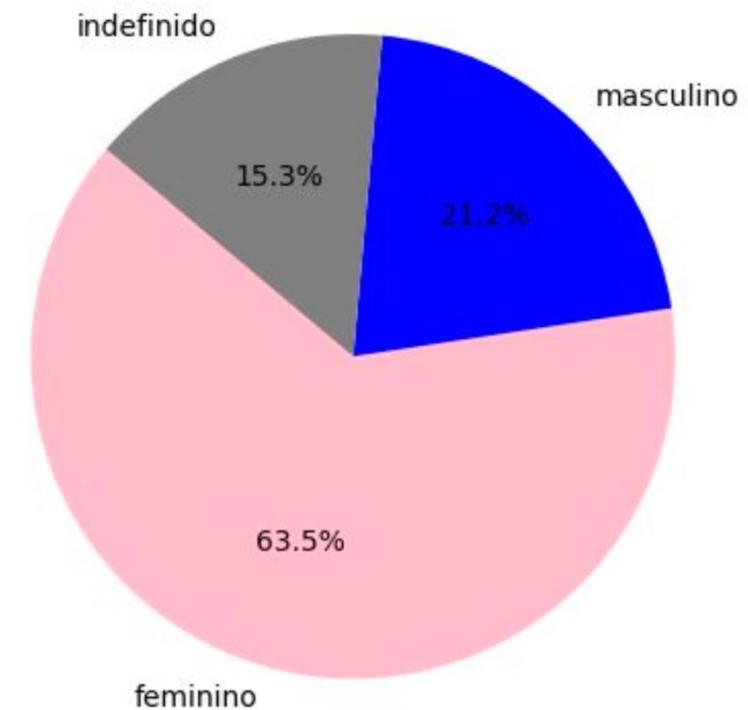
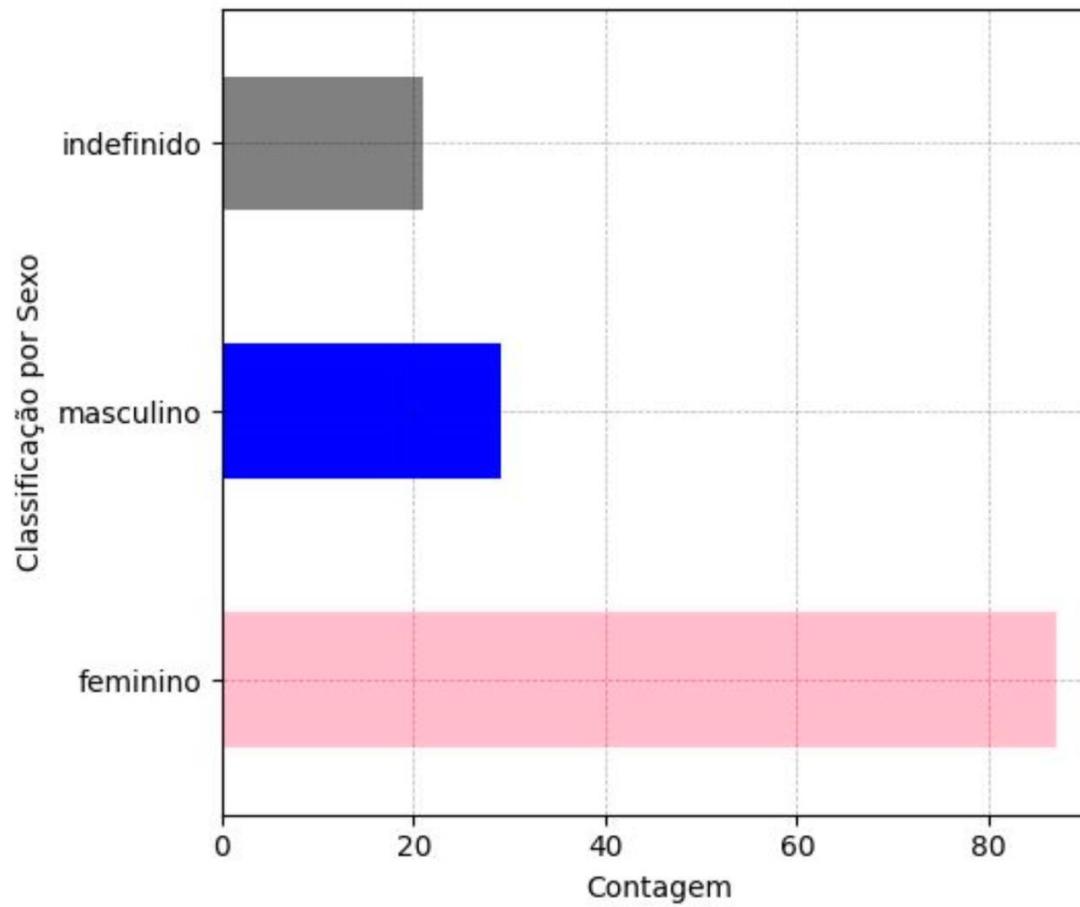
Retail analytics



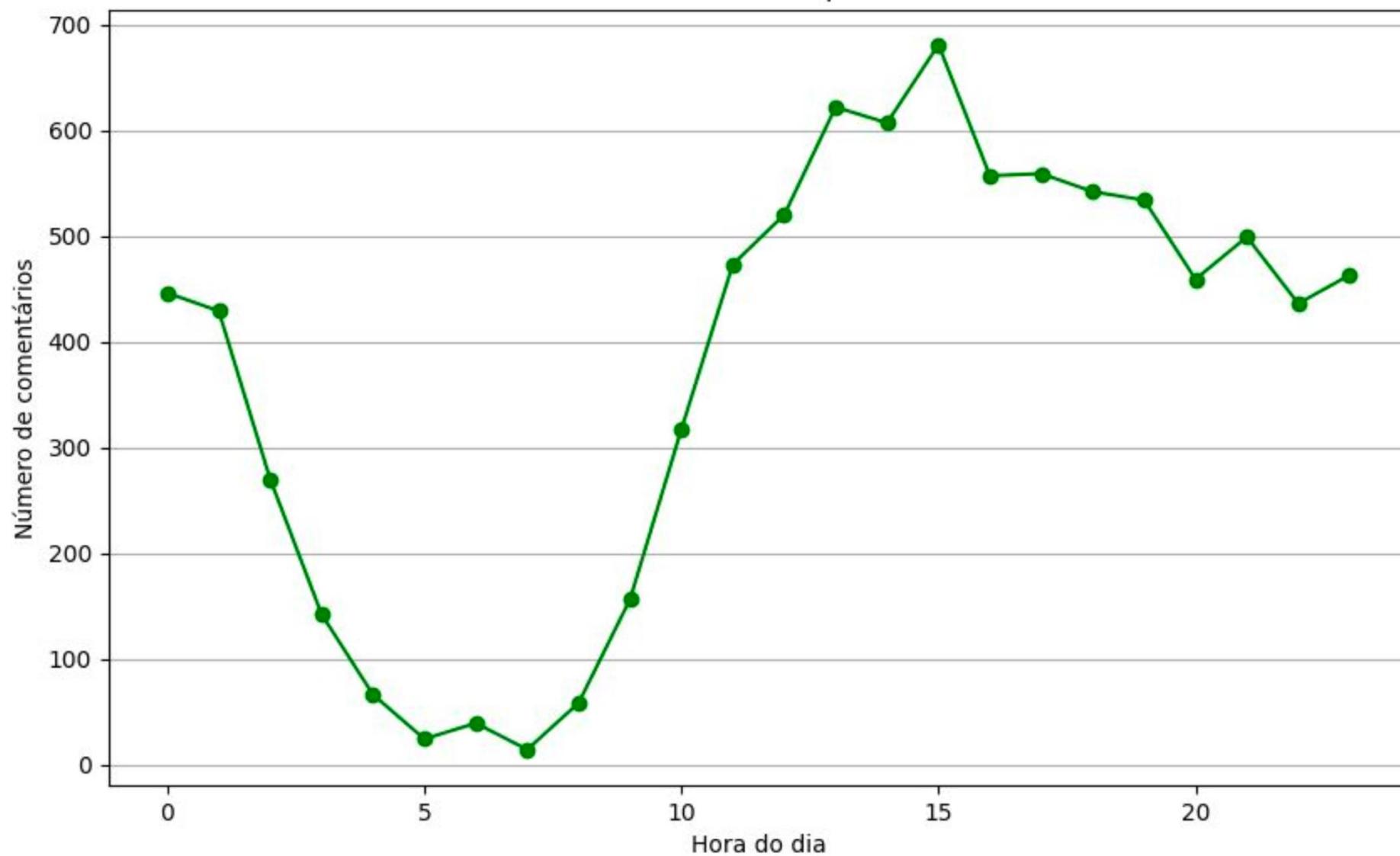
Augmented reality



Distribuição dos Comentários por Sexo - prefeitura_saude_sl.xlsx



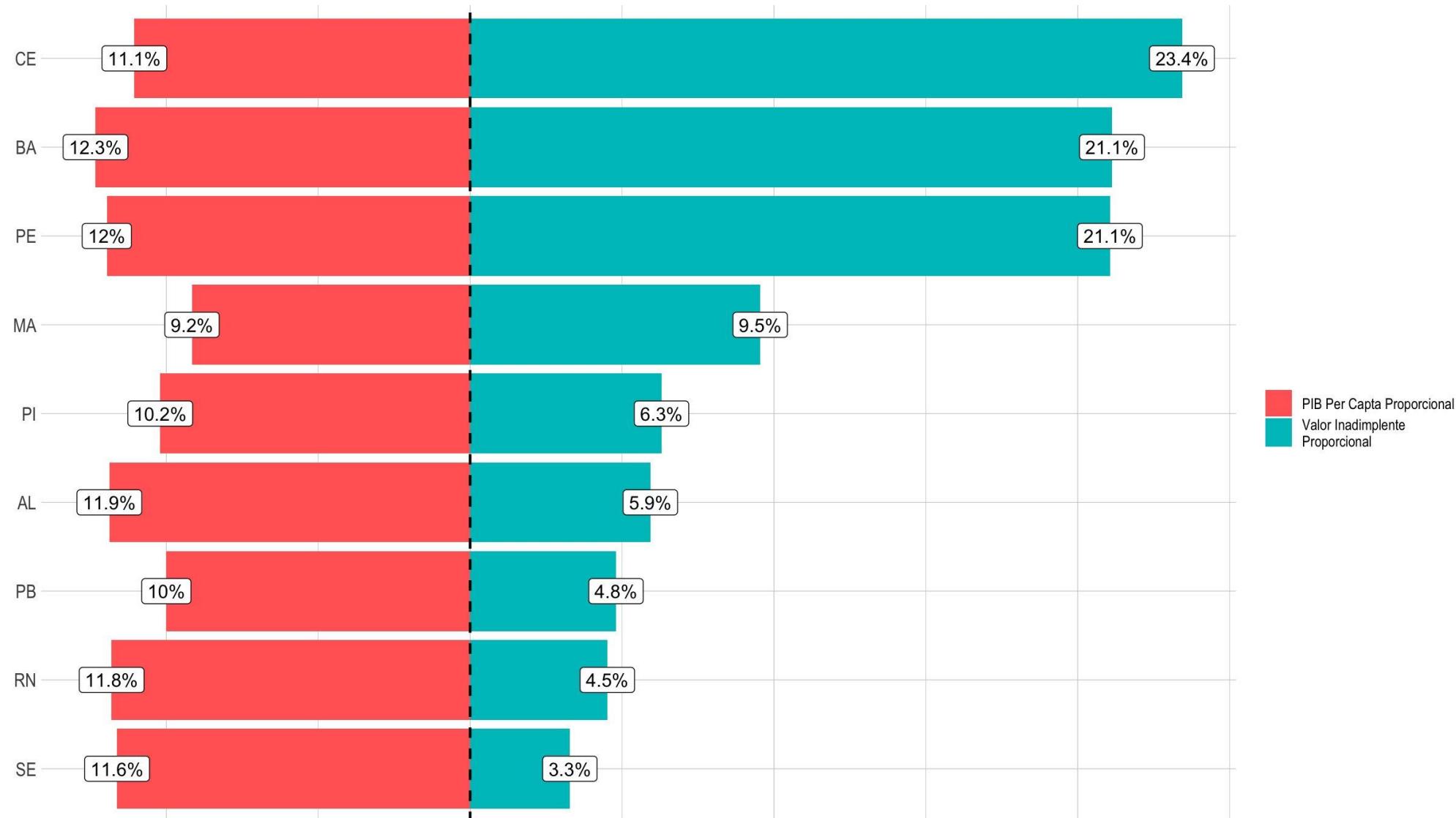
Número de comentários por hora do dia



amorDeus OrgulhoCorreios
CADEMINHAENCOMENDA tvglobo globonews Orgulhoperpender
Cadeencomenda celsorussomanno Orgulho
Quandovãoentregaremcomenda SP Quemaravilha
ParabénsCorreios Palhaçada Parabénsiniciativa correiosoficialdirect
Ondeencomenda Podemresponderdirect MERESPONDAM LIBEREMASENCOMENDAS Porgentileza favor
Cadê lamentável correiosoficialrespondedirect respondedm respondam Porgentileza favor
LIBEREMINHAENCOMENDA Euqueroencomenda respondam Liberapedido receita_federal Encomendaparaadadesdedia12
correiosoficialrespondam obrigada Porque Respondamdirect Atendimentoopéssimo 10agoranada
Show vergonha LIBEREMNOSSAENCOMENDA Agilidade
PeloamorDeus Ouseja Querocompra
jornaloglobo globonews Oacontecendo Querocomprad
Absurdo Mandeidirect Cadêpedido
Precisoencomenda Desrepeitoconsumidor Respondedirect Serviçopéssimo
Porfavor ConcursoPúblicoUrgente Mentira
LIBEREMNOSSASENCOMENDAS Liberemencomenda Entreguemencomenda querodinheirovolta Desdedia29
Carteiroatendido Entreguemencomenda querodinheirovolta Desdedia29
Olhemdirect Péssimoserviço correiosoficialolhadirect gentileza
procon_rj procons Issoabsurdo Faltarespeito
Descasototal jornaloglobo tvgloboBoatarde Queabsurdo
socorro Fabiano operacional Queremosencomendas
Desdedia26 Piorserviço Eai CADÊMINHAENCOMENDA pfv
paguei Palhaçad proconarioofficial procon_rj Descaso atendem precisoajuda misericórdia
correiosoficialcadêencomenda Minhaencomenda Querosaberencomenda
pagementoconfirmado encomendasparadas

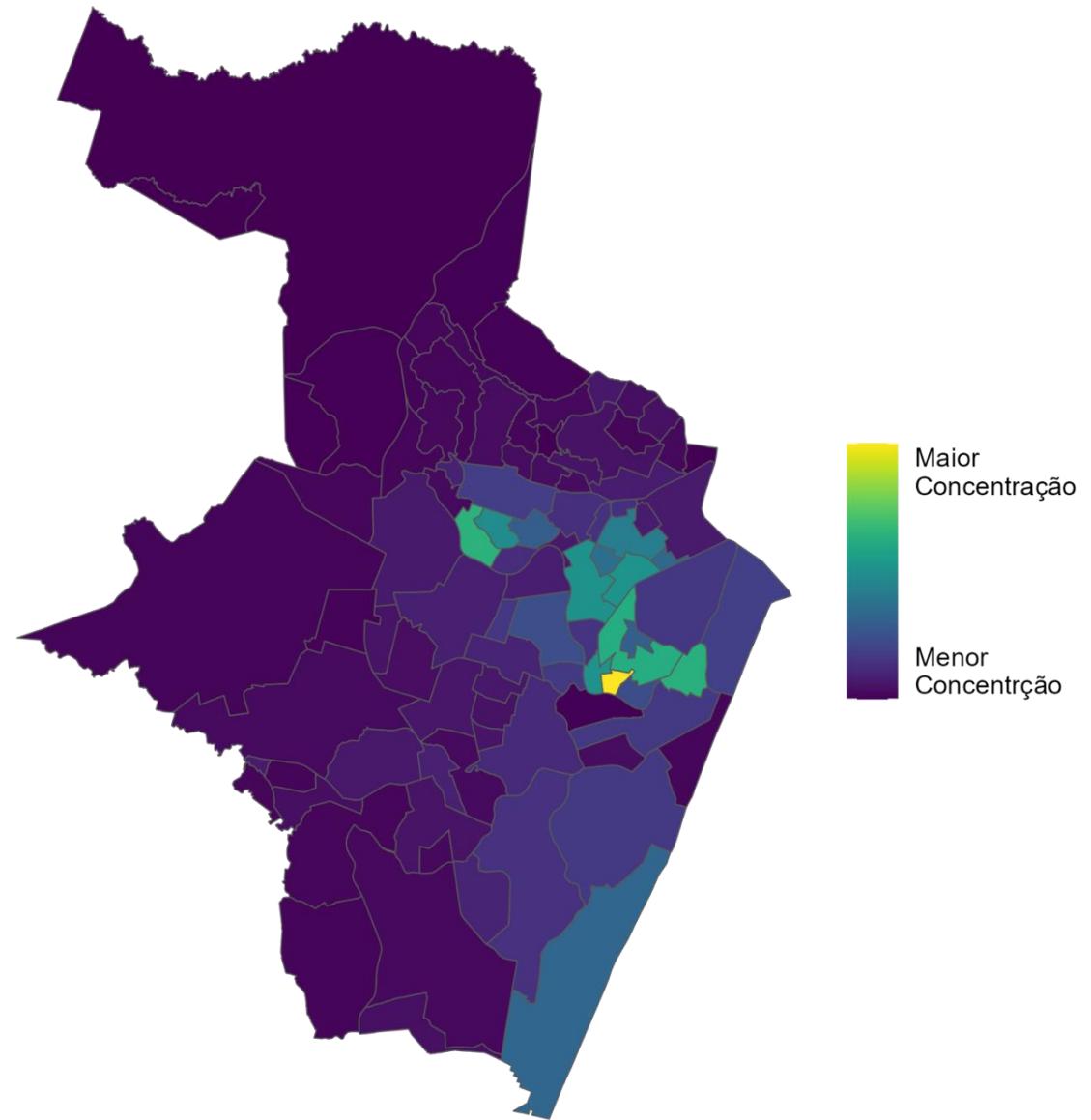
Capacidade de Pagamento de cada Estado Pib Per Capta Proporcional vs. Valor Inadimplente Proporcional

Fonte: Banco do Nordeste e IBGE



Concentração de Firmas por Bairro em Recife

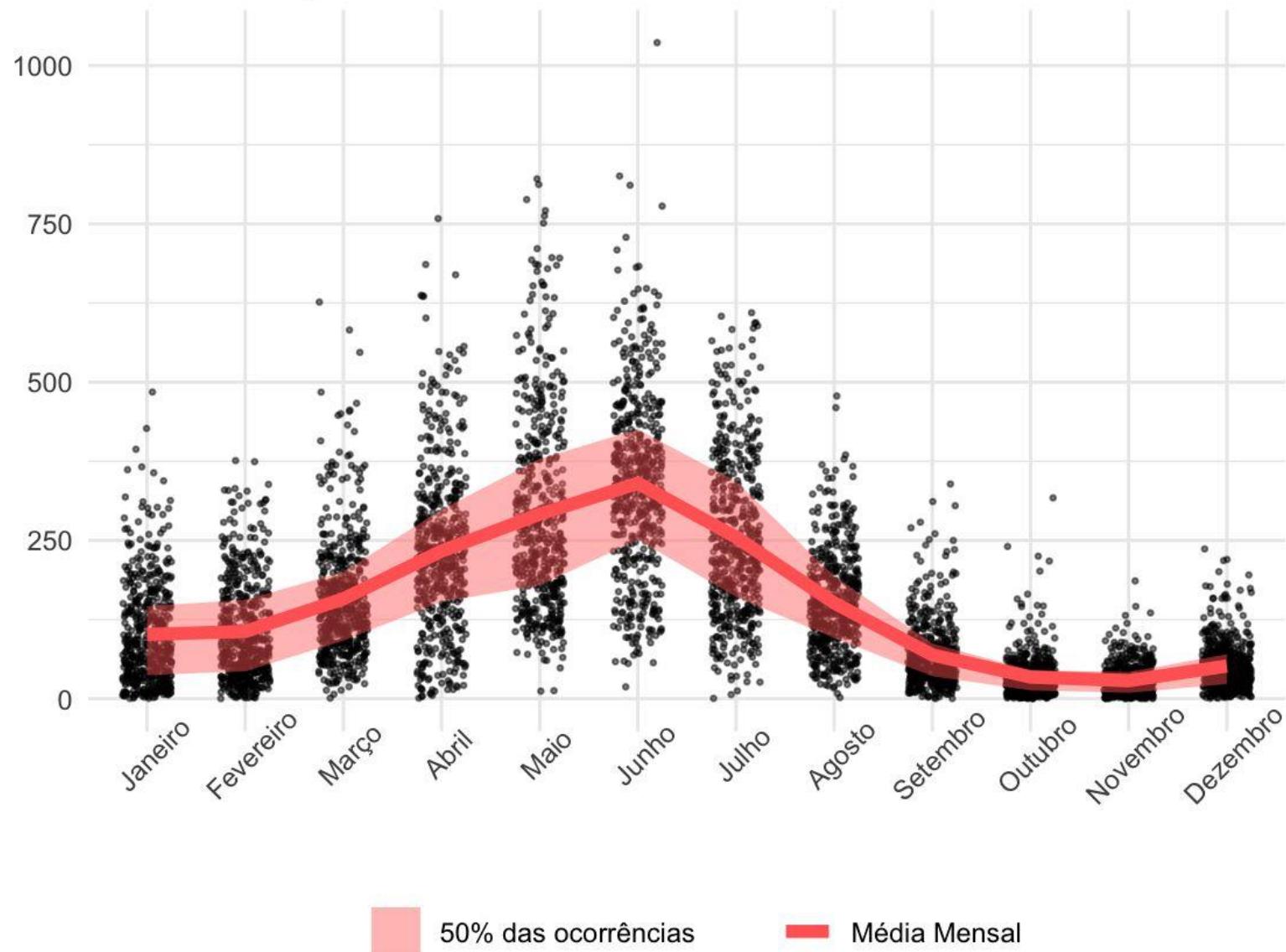
Fonte: Dados Abertos - Prefeitura de Recife



Evolução Pluviométrica em Recife - PE

Histórico de 2003 a 2023

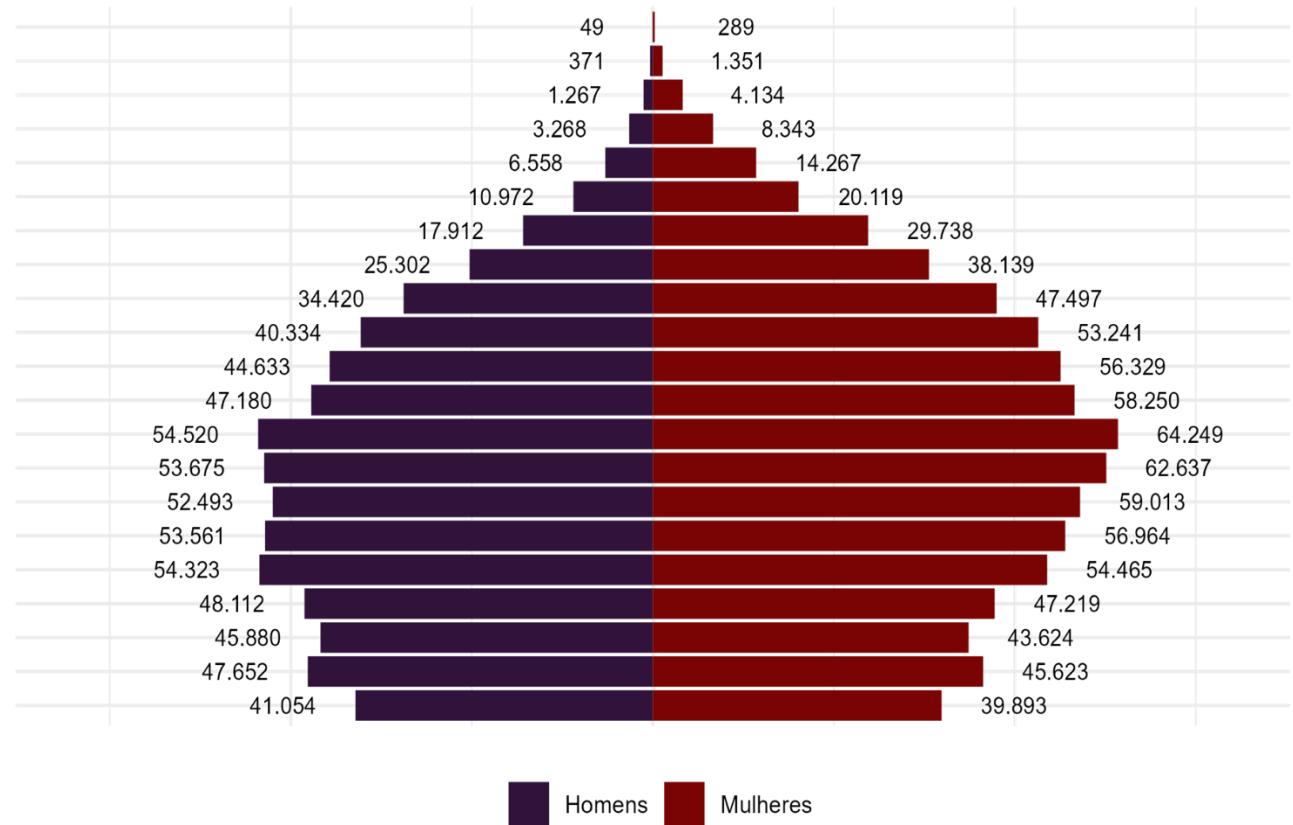
6.222 observações



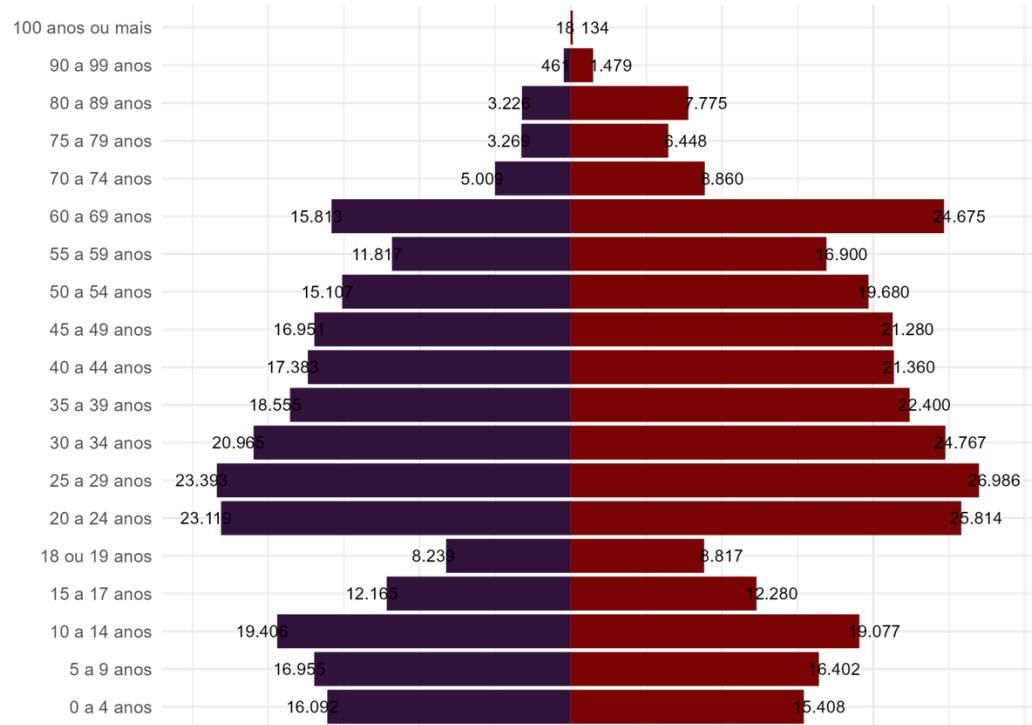
Prâmide Etária por Gênero, Recife

Valores Absolutos

Fonte: Censo 2022, IBGE



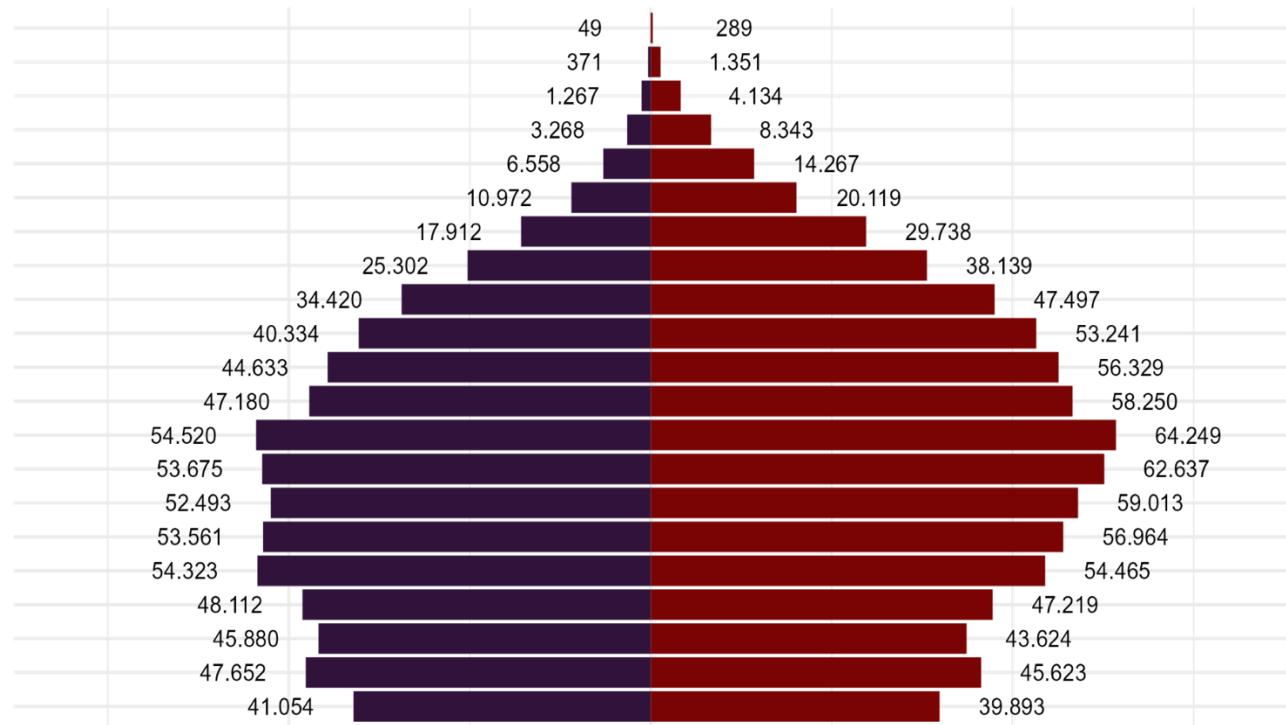
Prâmide Etária por Gênero - Bairros com Alagamento, Recife
Fonte: Censo 2010, IBGE



Homens Mulheres

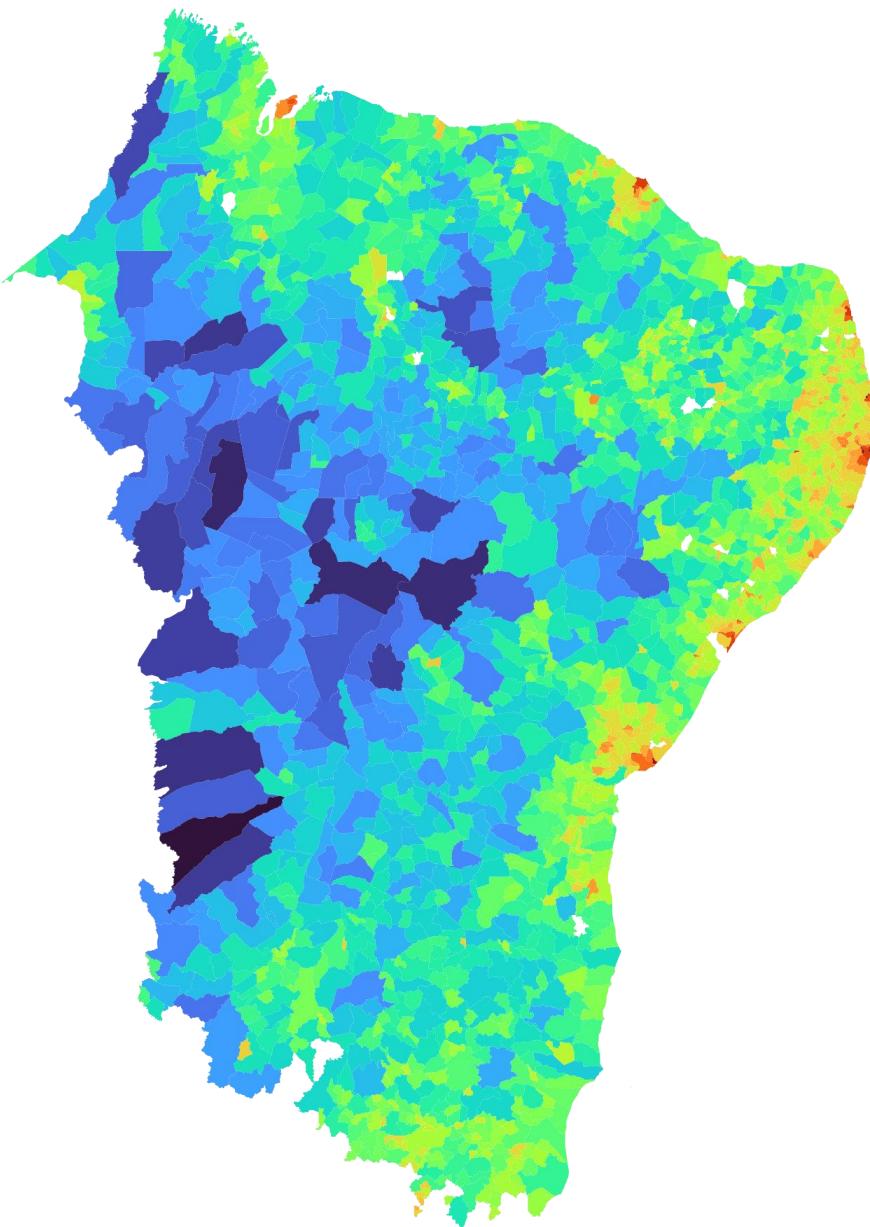
Prâmide Etária por Gênero, Recife

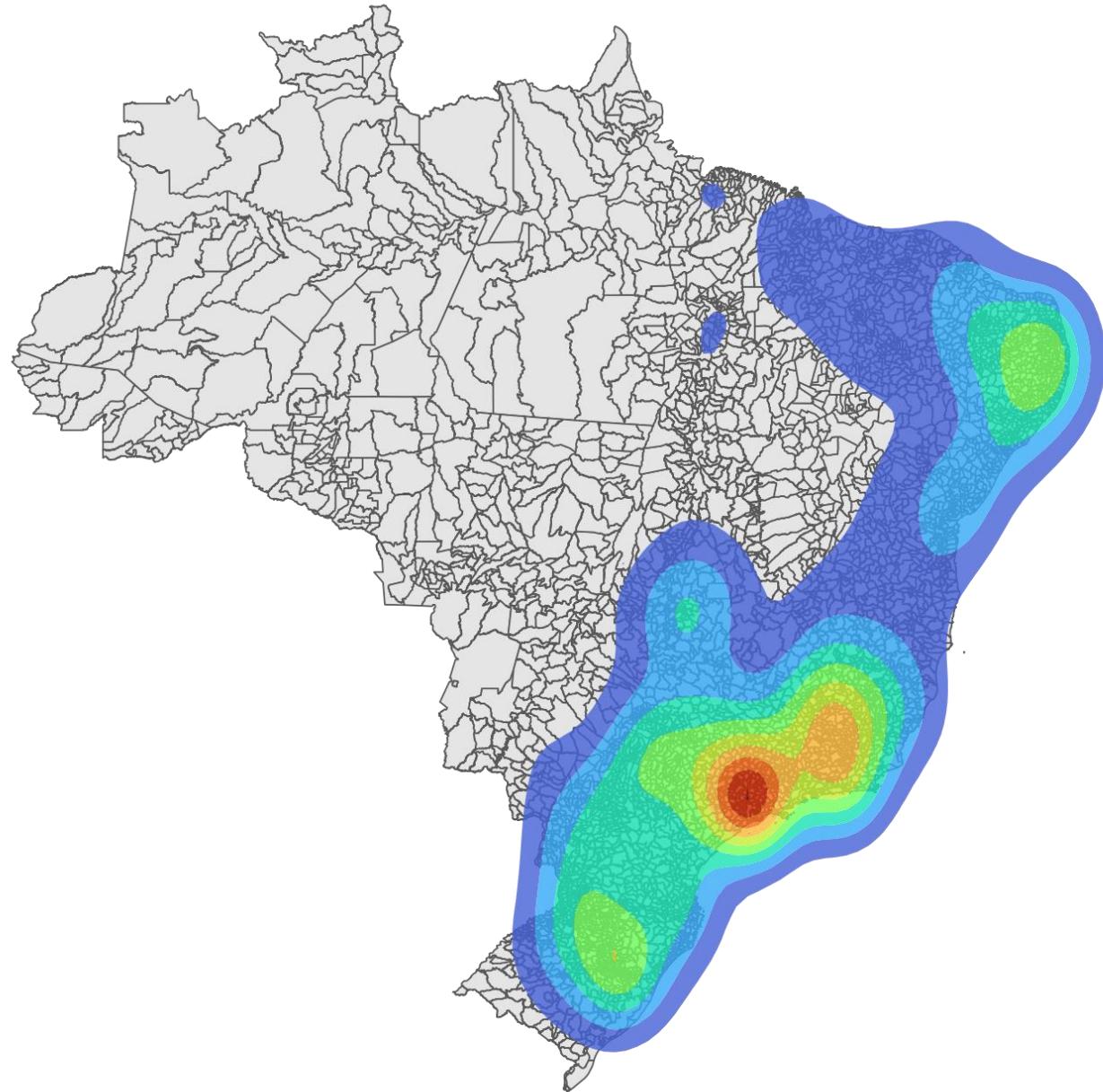
Valores Absolutos
Fonte: Censo 2022, IBGE



Homens Mulheres

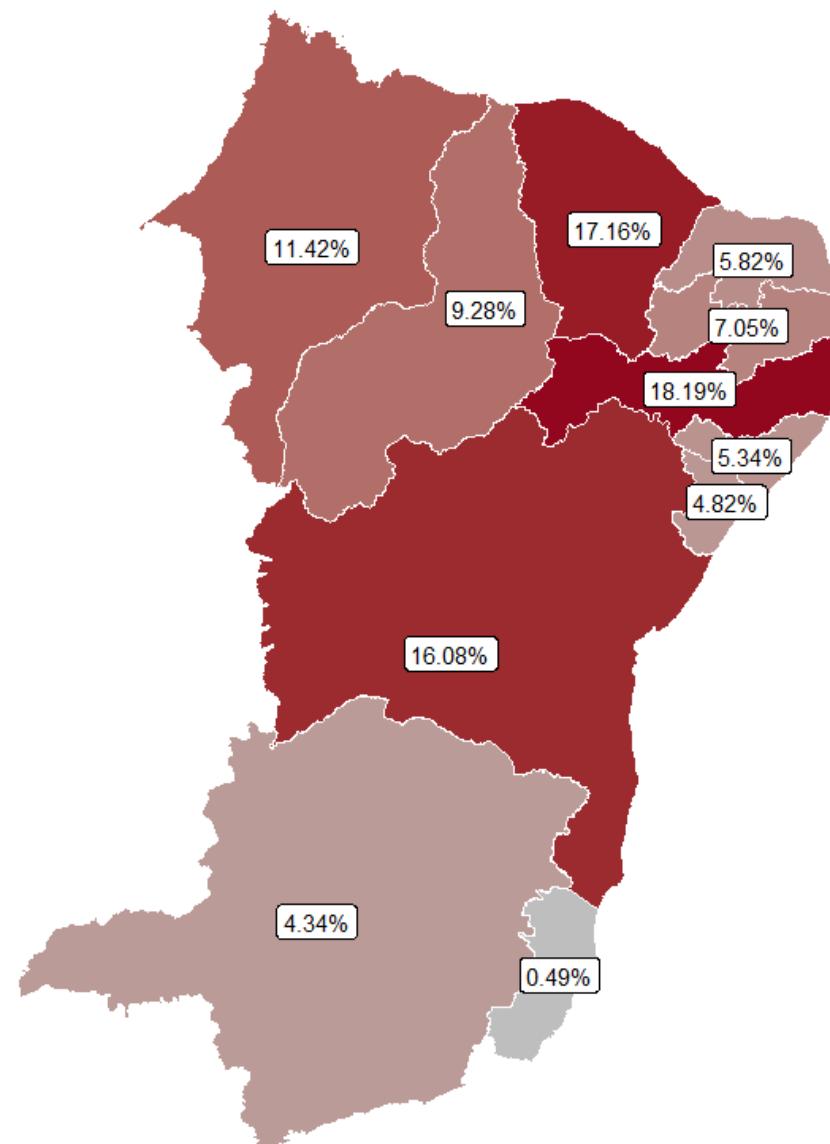
Densidade de Candidaturas a Vereador
Fonte: Estatísticas Eleitorais Candidatos 2020 - TSE 2024



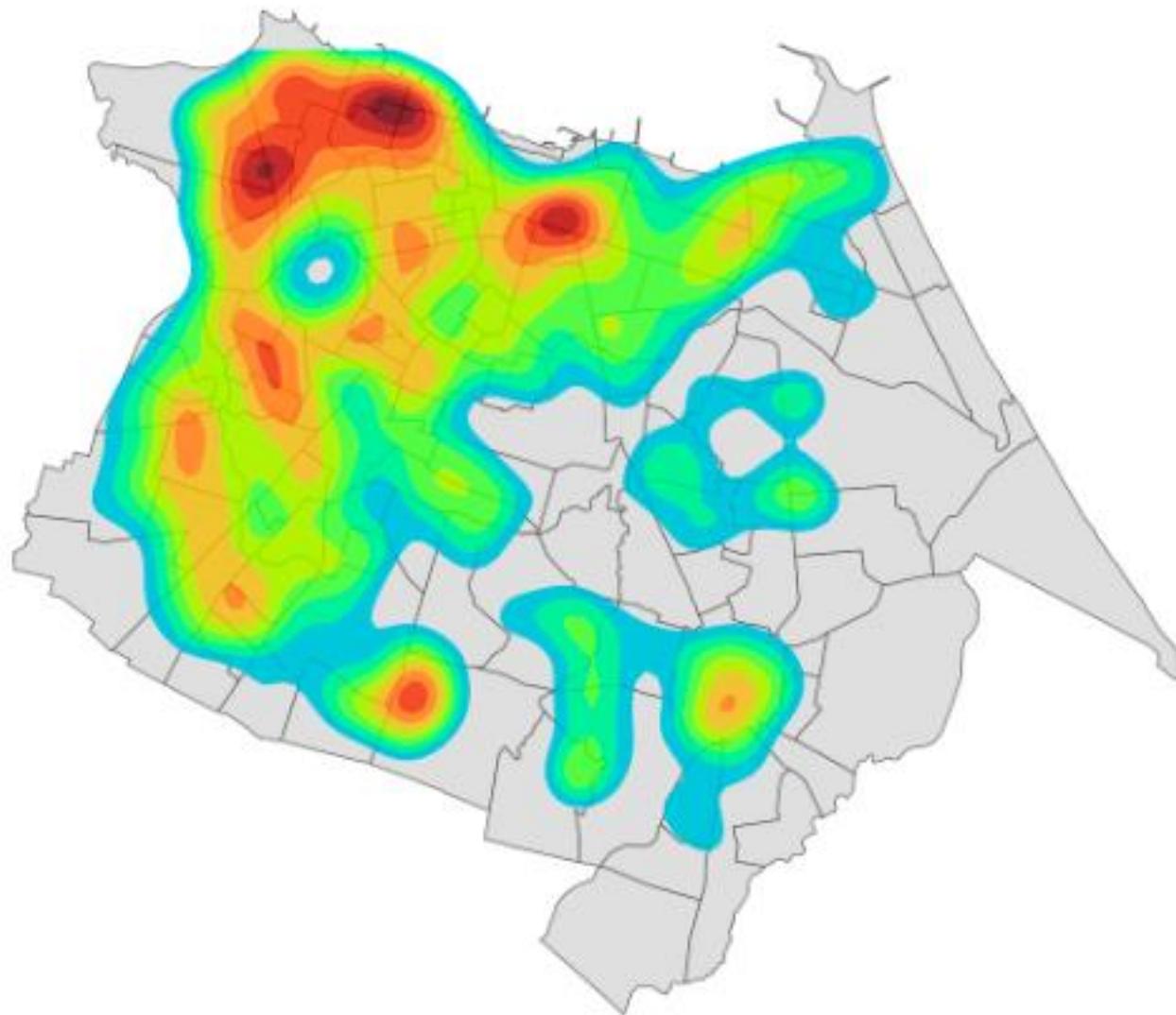


Recomendação Técnica de Investimento por Estado

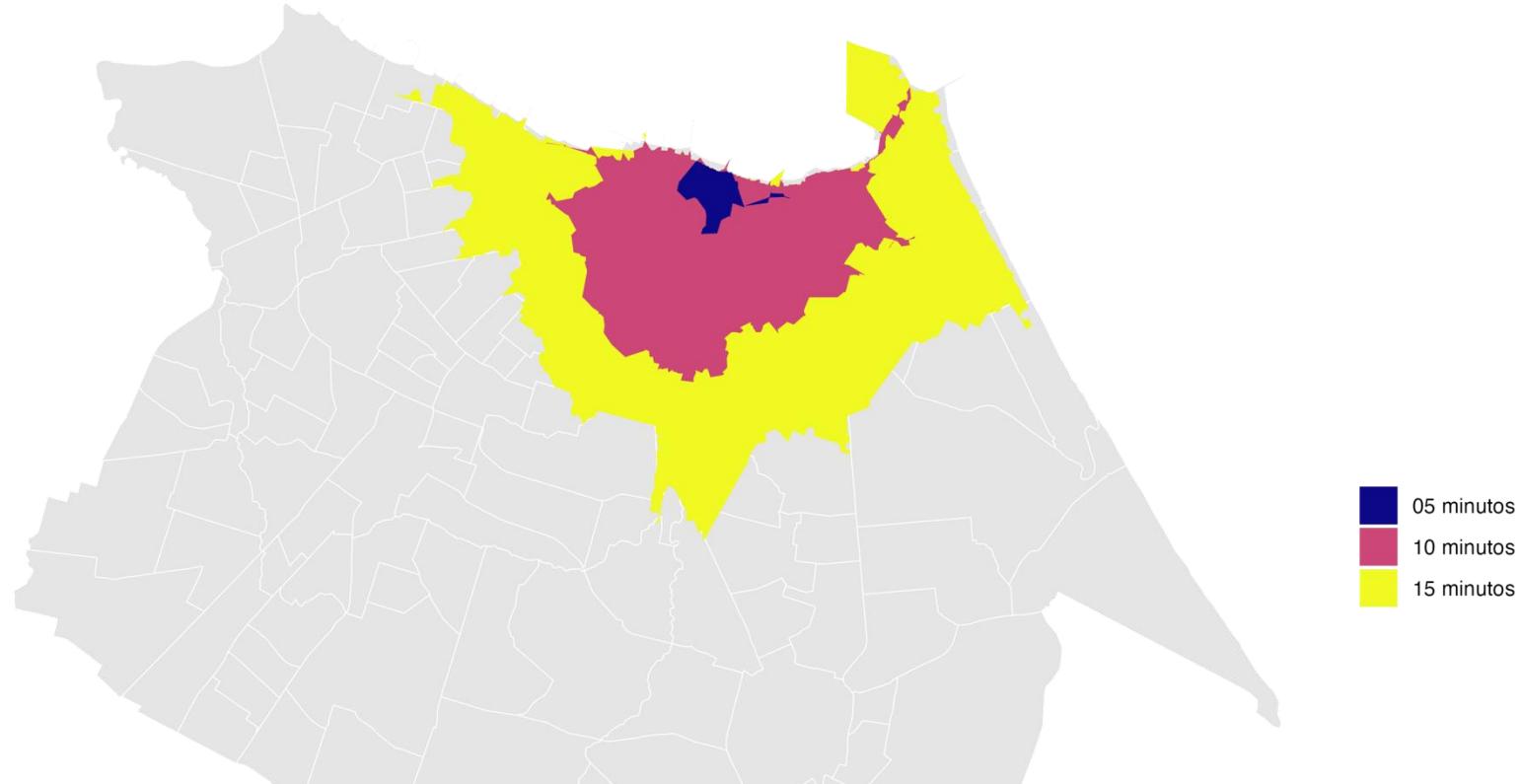
Fonte: Banco do Nordeste



Votos Elmano



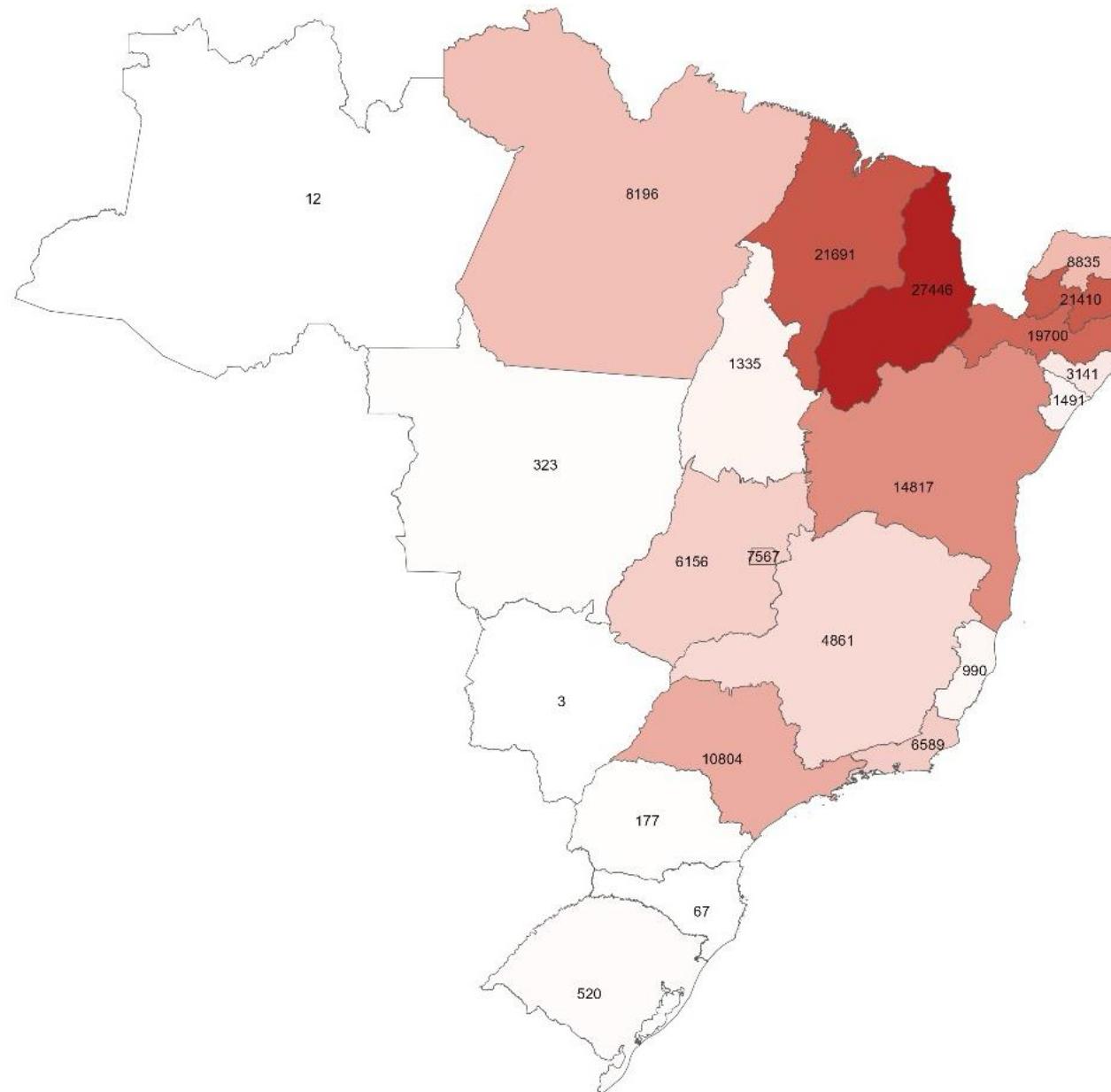
Zonas de Interesse

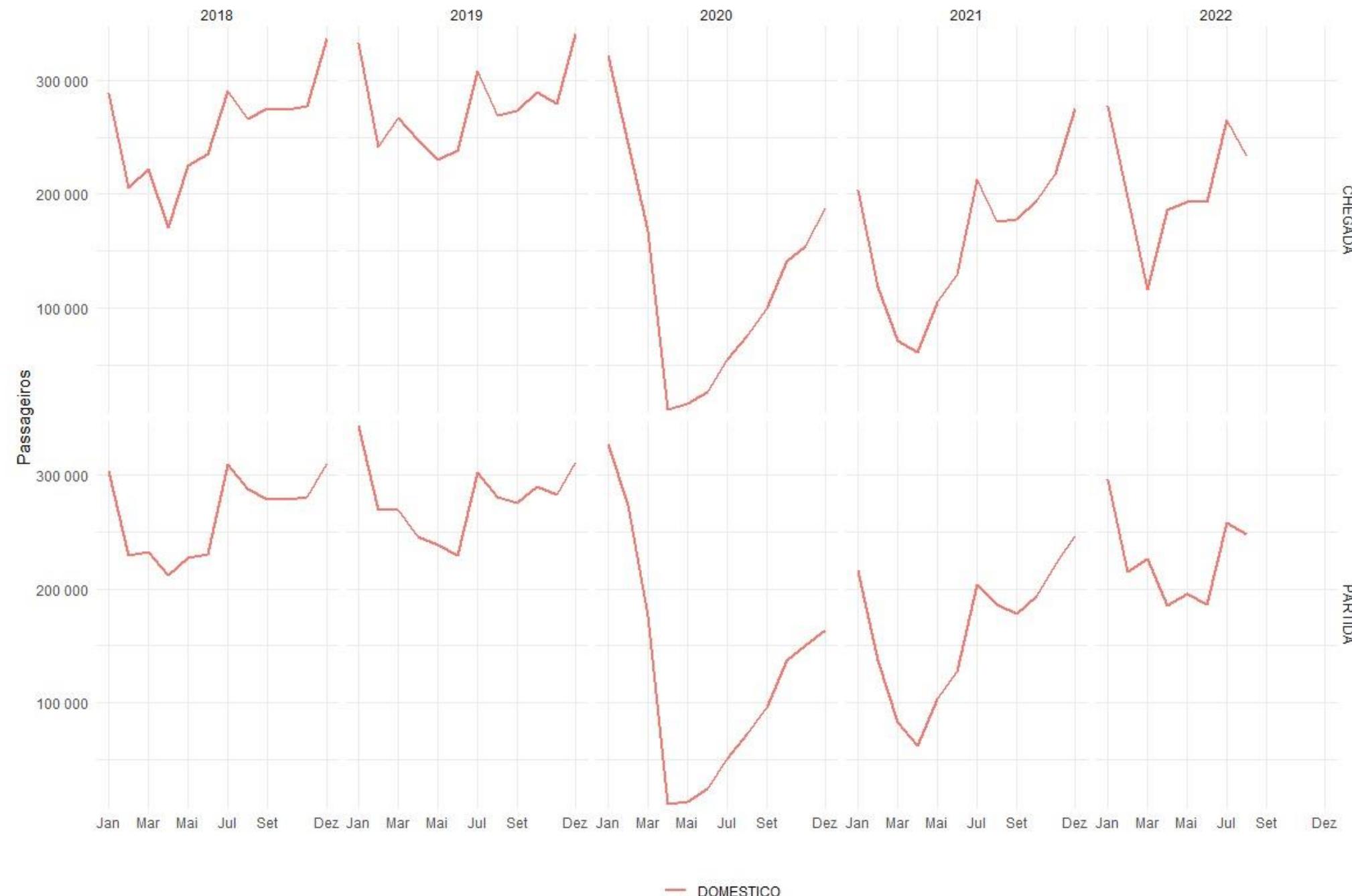


■ 05 minutos
■ 10 minutos
■ 15 minutos

| Zona de Interesse | Condomínio | Total de Unidades | Área Média (m ²) | Unidades/Lote |
|-------------------|------------|-------------------|------------------------------|---------------|
| 05 minutos | Sim | 301 | 88,8 | 21,5 |
| | Não | 70 | 83,8 | 1,3 |
| 10 minutos | Sim | 5.840 | 88,3 | 28,8 |
| | Não | 3.830 | 85,3 | 1,44 |
| 15 minutos | Sim | 11.527 | 86,4 | 31,3 |
| | Não | 12.166 | 85,0 | 1,27 |

Mapa de Estados de Destino de Viagens de Ônibus Saindo do Ceará





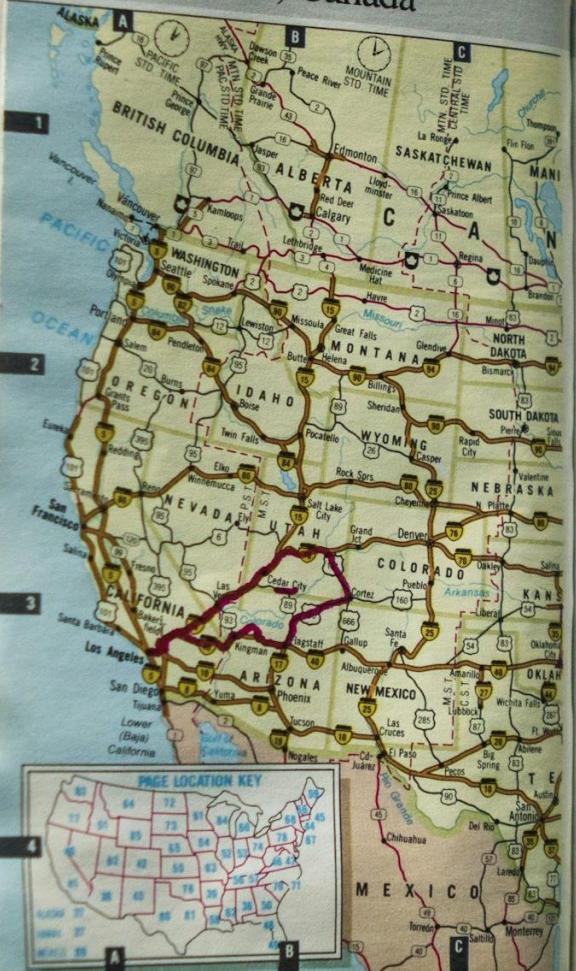
Estudo de caso motivador: Mudança de carreira!



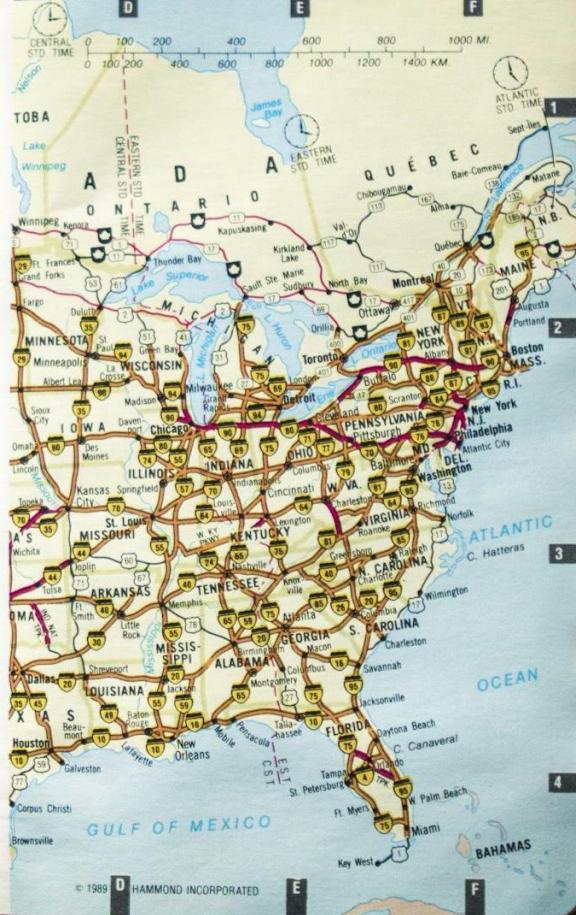




United States, Canada



35



© 1989 HAMMOND INCORPORATED

A on one?



1. Conceitos Básicos

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Project: (None)

Untitled1

Source on Save Run Source

1

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |
```

Environment History Connections Tutorial

Import Dataset 250 MiB Grid C

R Global Environment

Name Type Length Size Value

Environment is empty

Files Plots Packages Help Viewer Presentation

Zoom Export C



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

+ Go to file/function Addins Project: (None)

Untitled1 × Source on Save Run Source

1 |

1:1 (Top Level) R Script

Console Terminal × Background Jobs ×

R 4.3.1 ~/

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |
```

Environment History Connections Tutorial

Import Dataset 250 MiB Grid C

R Global Environment

Name Type Length Size Value

Environment is empty

Files Plots Packages Help Viewer Presentation

Zoom Export

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

+ Go to file/function Addins Project: (None)

Untitled1 × Source on Save Run Source Environment History Connections Tutorial Import Dataset 250 MiB Grid Global Environment

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/ ↻

```
> 2 + 2
[1] 4
> 5 * 3
[1] 15
> 10 / 2
[1] 5
> 2^3
[1] 8
>
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Untitled1 x

Source on Save Run Source

1

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]
```

> |

Project: (None)

Environment History Connections Tutorial

Import Dataset 250 MiB Grid C

R Global Environment

Name Type Length Size Value

Environment is empty

Files Plots Packages Help Viewer Presentation

Zoom Export C

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

+ Go to file/function Addins

Untitled1 × Source on Save Run Source

1 |

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/ →

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |
```

Project: (None)

Environment History Connections Tutorial

Import Dataset 250 MiB Grid C

R Global Environment

Name Type Length Size Value

Environment is empty

Files Plots Packages Help Viewer Presentation

Zoom Export

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

+ Go to file/function Addins Project: (None)

Untitled1 × Source on Save Run Source

1 |

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/ →

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |
```

Environment History Connections Tutorial

Import Dataset 250 MiB Grid C

R Global Environment

Name Type Length Size Value

Environment is empty

Files Plots Packages Help Viewer Presentation

Zoom Export

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

New Session Addins Project: (None)

Untitled1 x Sources

Interrupt R Terminate R...
Restart R Ctrl+Shift+F10
Set Working Directory...
Load Workspace...
Save Workspace As...
Clear Workspace...
Quit Session... Ctrl+Q

1:1 (Top Level) R Script

Console Terminal Background Jobs

R 4.3.1 ~/

```
R version 4.3.1 (2023-06-16 ucrt) -- "Beagle Scouts"
Copyright (C) 2023 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> |
```

Environment History Connections Tutorial

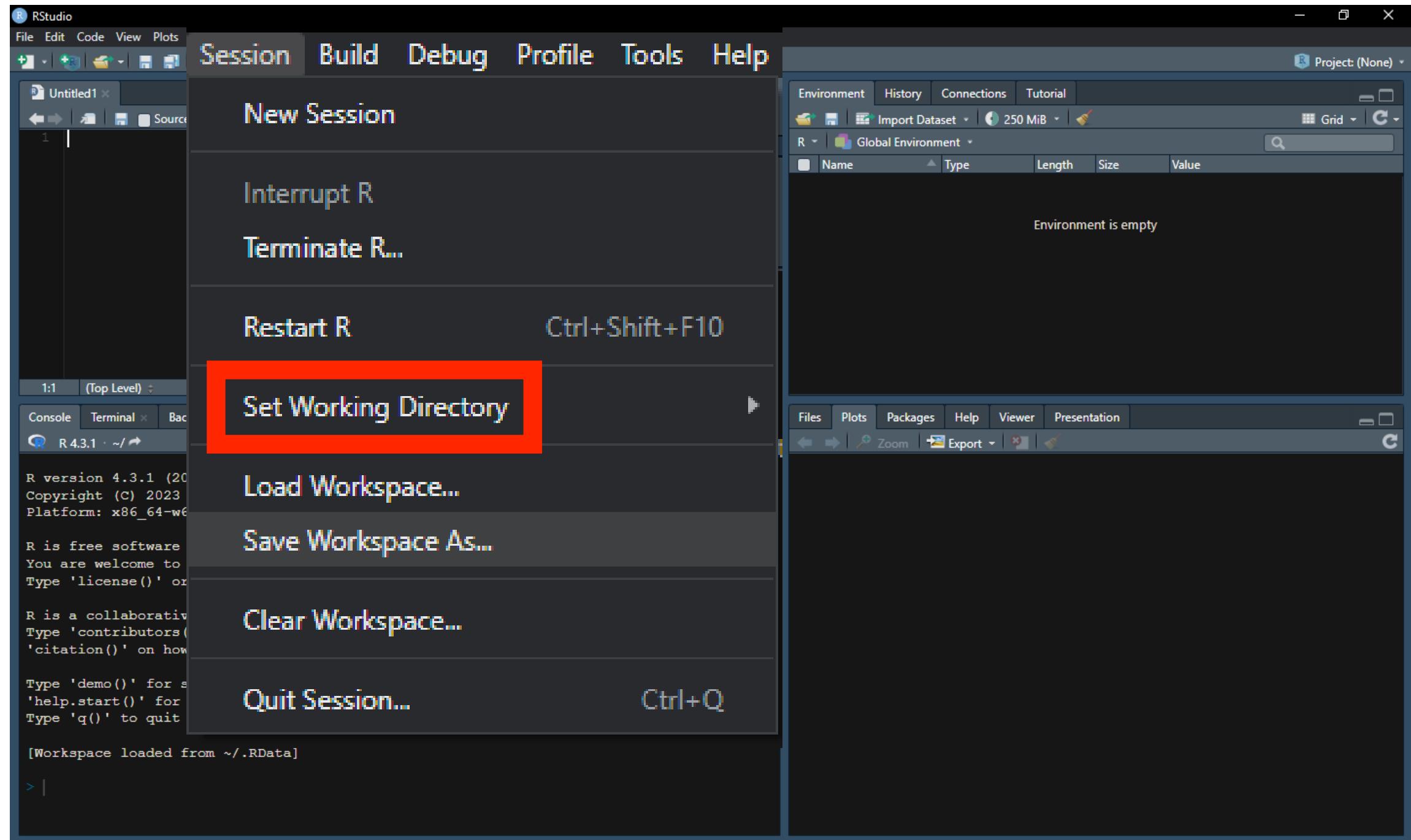
Run Import Dataset 250 MiB Grid C

R Global Environment

| Name | Type | Length | Size | Value |
|----------------------|------|--------|------|-------|
| Environment is empty | | | | |

Files Plots Packages Help Viewer Presentation

Zoom Export C



```
read.csv(file = 'data.csv')
```

```
> read.csv('data.csv')
```

| | state | abb | region | population | total |
|----|----------------------|-----|---------------|------------|-------|
| 1 | Alabama | AL | South | 4779736 | 135 |
| 2 | Alaska | AK | West | 710231 | 19 |
| 3 | Arizona | AZ | West | 6392017 | 232 |
| 4 | Arkansas | AR | South | 2915918 | 93 |
| 5 | California | CA | West | 37253956 | 1257 |
| 6 | Colorado | CO | West | 5029196 | 65 |
| 7 | Connecticut | CT | Northeast | 3574097 | 97 |
| 8 | Delaware | DE | South | 897934 | 38 |
| 9 | District of Columbia | DC | South | 601723 | 99 |
| 10 | Florida | FL | South | 19687653 | 669 |
| 11 | Georgia | GA | South | 9920000 | 376 |
| 12 | Hawaii | HI | West | 1360301 | 7 |
| 13 | Idaho | ID | West | 1567582 | 12 |
| 14 | Illinois | IL | North Central | 12830632 | 364 |
| 15 | Indiana | IN | North Central | 6483802 | 142 |

Tipos de banco de dados

```
read.csv('data.csv')
```

```
read.csv(path = 'data.csv')
```

```
?read.csv()
```

R RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins Project: (None)

Files Plots Packages Help Viewer Presentation

R: Data Input Find in Topic

read.table {utils}

R Documentation

Data Input

Description

Reads a file in table format and creates a data frame from it, with cases corresponding to lines and variables to fields in the file.

Usage

```
read.table(file, header = FALSE, sep = "", quote = "\t",
           dec = ".", numerals = c("allow.loss", "warn.loss", "no.loss"),
           row.names, col.names, as.is = !stringsAsFactors, tryLogical = TRUE,
           na.strings = "NA", colClasses = NA, nrow = -1,
           skip = 0, check.names = TRUE, fill = !blank.lines.skip,
           strip.white = FALSE, blank.lines.skip = TRUE,
           comment.char = "#",
           allowEscapes = FALSE, flush = FALSE,
           stringsAsFactors = FALSE,
           fileEncoding = "", encoding = "unknown", text, skipNul = FALSE)

read.csv(file, header = TRUE, sep = ",", quote = "\t",
         dec = ".", fill = TRUE, comment.char = "", ...)

read.csv2(file, header = TRUE, sep = ";", quote = "\t",
          dec = ",", fill = TRUE, comment.char = "", ...)

read.delim(file, header = TRUE, sep = "\t", quote = "\t",
           dec = ".", fill = TRUE, comment.char = "", ...)

read.delim2(file, header = TRUE, sep = "\t", quote = "\t",
            dec = ",", fill = TRUE, comment.char = "", ...)
```

Arguments

the name of the file which the data are to be read from. Each row of the table appears as one line of the file. If it does not contain an *absolute path*, the file name is *relative* to the current working directory.

Description

Reads a file in table format and creates a data frame from it, with cases corresponding to lines and variables to fields in the file.

Usage

```
read.table(file, header = FALSE, sep = "\t", quote = "\r\n",
          dec = ".", numerals = c("allow.dots", "warn.loss", "no.lo...),
          row.names, col.names, na.strings = "", colClasses = "nrclass", nrclass,
          skip = 0, check.names = TRUE, check.names = !is.na(readLines(f)),
          strip.white = FALSE, comment.char = "#", strip.white = TRUE,
          comment.char = "#", allowEscapes = TRUE, allowEscapes = TRUE,
          stringsAsFactors = TRUE, stringsAsFactors = TRUE,
          fileEncoding = "", encoding = "unknown", text, skipNul = FALSE)

read.csv(file, header = TRUE, sep = "\t", quote = "\r\n",
         dec = ".", fill = TRUE, comment.char = "#", ...)

read.csv2(file, header = TRUE, sep = ";", quote = "\r\n",
          dec = ",", fill = TRUE, comment.char = "#", ...)

read.delim(file, header = TRUE, sep = "\t", quote = "\r\n",
           dec = ".", fill = TRUE, comment.char = "#", ...)

readLines(f, n = -1, sep = "\r\n", quote = "\r\n",
          dec = ".", fill = TRUE, comment.char = "#", ...)
```

Arguments

the name of the file we want to read from. Each row of the table appears as one line of the file. If it does not contain an *absolute path*, the file name is *relative* to the current working directory.

Linguagem, sintaxe e argumentos

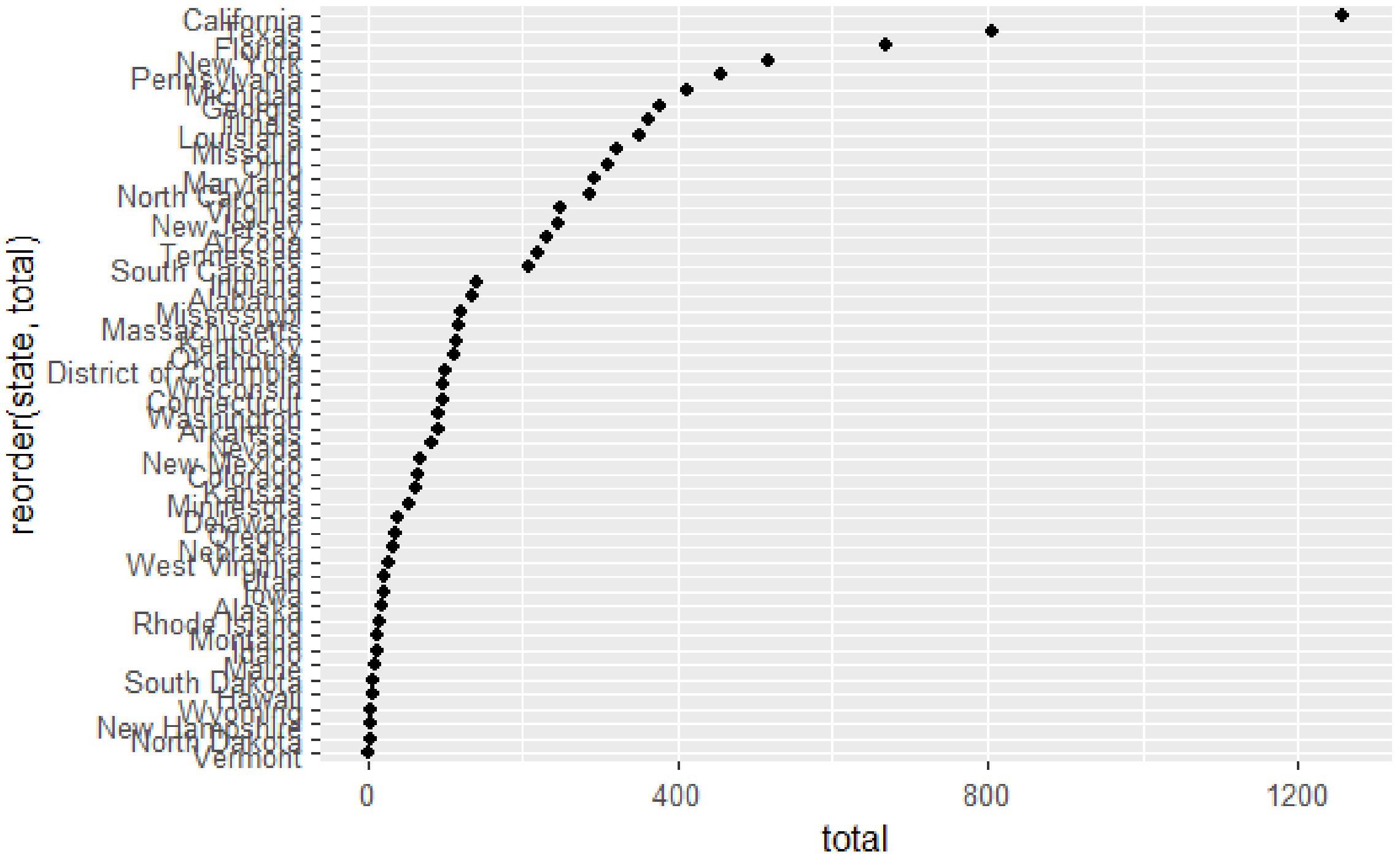
`funcao(arg_1 = var_1, arg_2 = var_2, ...)`

```
install.packages('tidyverse')
```

```
library(tidyverse)
```

Pacotes

```
read.csv('data.csv') %>%  
  ggplot() +  
    geom_point(aes(total, reorder(state,  
total)))
```



2. Criando, carregando e inspecionando dados

```
a <- 2
```

```
a
```

```
print(a)
```

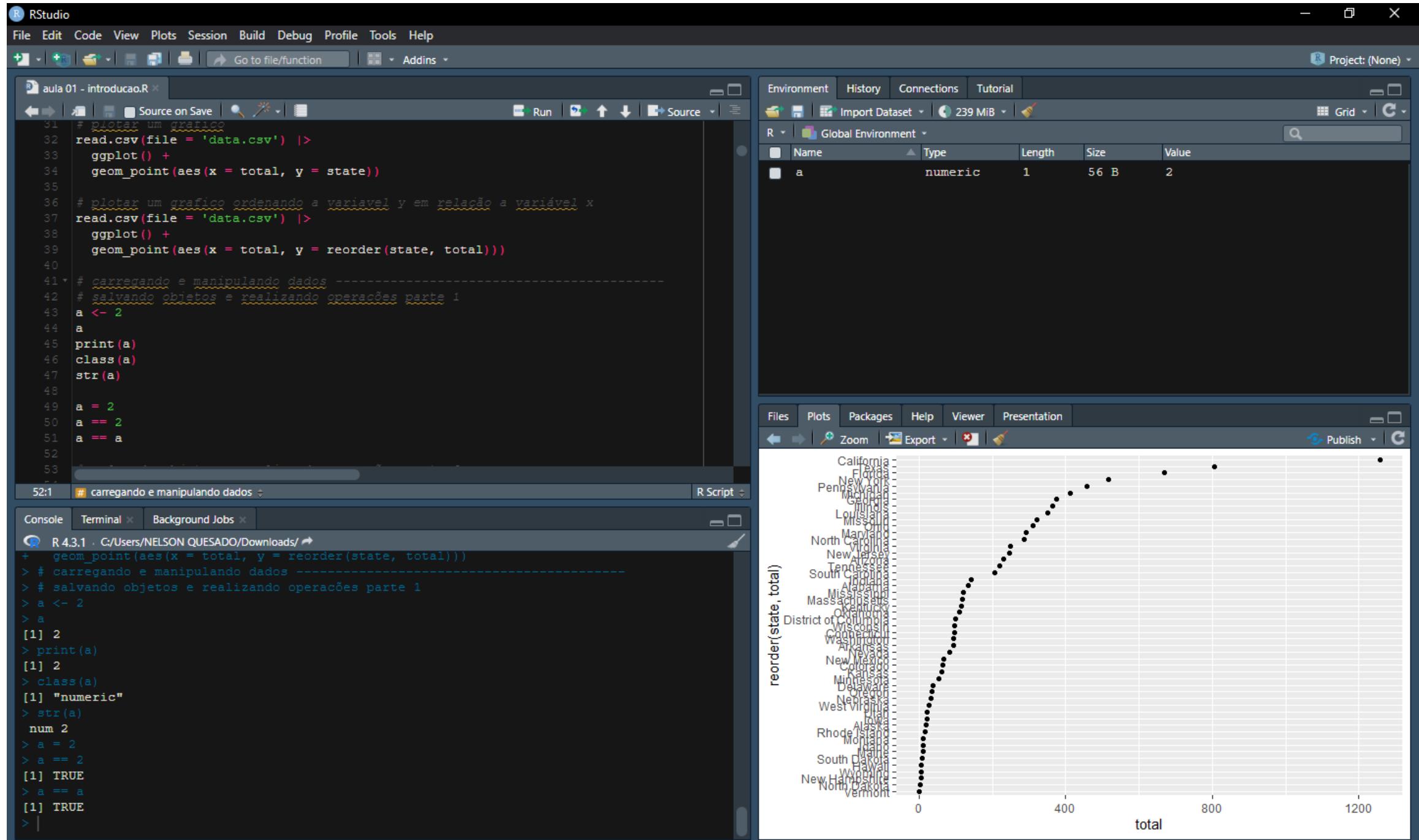
```
class(a)
```

```
str(a)
```

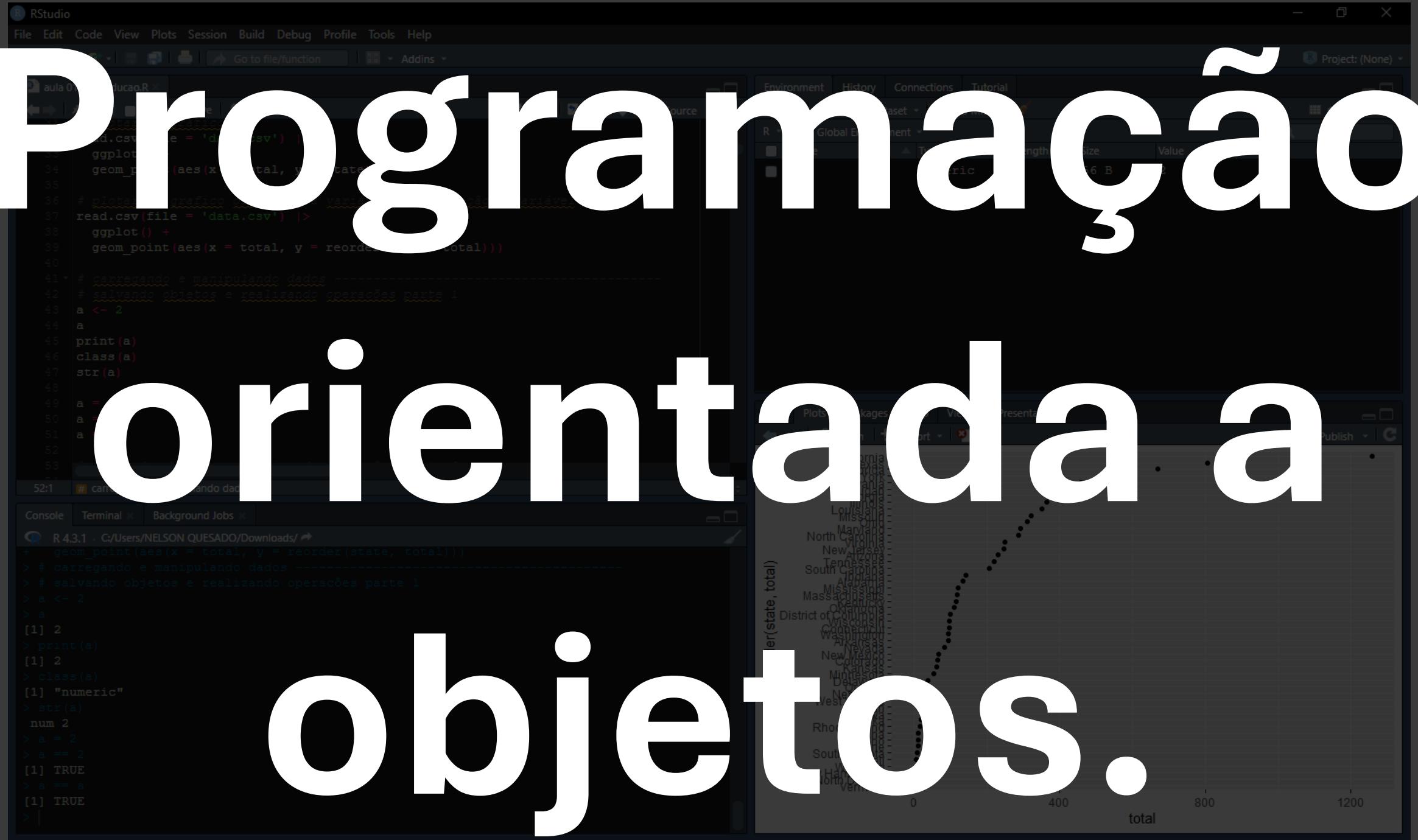
```
a = 2
```

```
a == 2
```

```
a == a
```



Programação orientada a objetos.

A screenshot of the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The left sidebar shows a file tree with 'aula 0' and 'data.csv'. The main workspace shows R code in the Source tab and a ggplot2 scatter plot in the Plots tab. The plot displays a positive correlation between state population and total value, with points for various US states and territories.

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Environment History Connections Tutorial
Project: (None)
R Global Environment Variables Value
source(aula 0.R)
## carregando e manipulando dados
## salvando objetos e realizando operações parte 1
a <- 2
a
print(a)
class(a)
str(a)
a
a = 2
a
[1] 2
> print(a)
[1] 2
> class(a)
[1] "numeric"
> str(a)
num 2
> a == 2
> a == 2
[1] TRUE
> a == a
[1] TRUE
>
ggplot(data = read.csv(file = 'data.csv')) +
  geom_point(aes(x = total, y = reorder(state, total)))
# carregando e manipulando dados
# salvando objetos e realizando operações parte 1
a <- 2
a
print(a)
class(a)
str(a)
a
a = 2
a
[1] 2
> print(a)
[1] 2
> class(a)
[1] "numeric"
> str(a)
num 2
> a == 2
> a == 2
[1] TRUE
> a == a
[1] TRUE
>

```

R 4.3.1 · C:/Users/NELSON QUESADO/Downloads/ ↵

```
+ geom_point(aes(x = total, y = reorder(state, total)))
> # carregando e manipulando dados
> # salvando objetos e realizando operações parte 1
> a <- 2
> a
[1] 2
> print(a)
[1] 2
> class(a)
[1] "numeric"
> str(a)
num 2
> a == 2
> a == 2
[1] TRUE
> a == a
[1] TRUE
>
```

Console Terminal Background Jobs

Plots Packages View Presentations Publish C

total

0 400 800 1200

state

Louisiana Mississippi Maryland North Carolina Virginia New Jersey Tennessee South Carolina Alabama Massachusetts District of Columbia Wisconsin Quebec Arkansas Nevada New Mexico Minnesota Delaware New Hampshire Rhode Island South Carolina Hawaii Vermont

```
Objeto <- funcao(arg_1 = var_1)
```

b <- 5

c <- -4

a == b

c == -4

class(a) == class(b)

a + b * c

(a + b) * c

```
d <- (a + b) * c
```

```
d
```

```
class(d)
```

```
A <- '2'
```

```
print(A)
```

```
class(A)
```

```
str(A)
```

```
class(A) == class(a)
```

```
class(a) <- 'character'
```

```
class(A) == class(a)
```

```
class(a) <- 'numeric'
```

```
class(A) == class(a)
```

```
ls()
```

```
rm(A, b, c, d)
```

```
c(1, 2, 3, 4)
```

```
1:4
```

```
c(1, 2, 3, 4) == 1:4
```

```
x <- 1:4
```

```
x
```

```
print(x)
```

```
class(x)
```

```
str(x)
```

x == a

y <- x == a

y

class(y)

as.numeric(y)

sum(as.numeric(y))

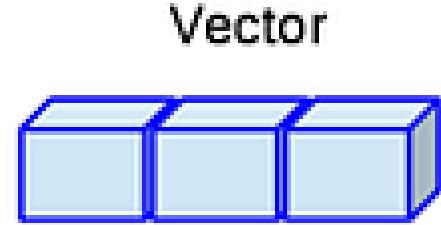
length(y)

sum(as.numeric(y)) / length(y)

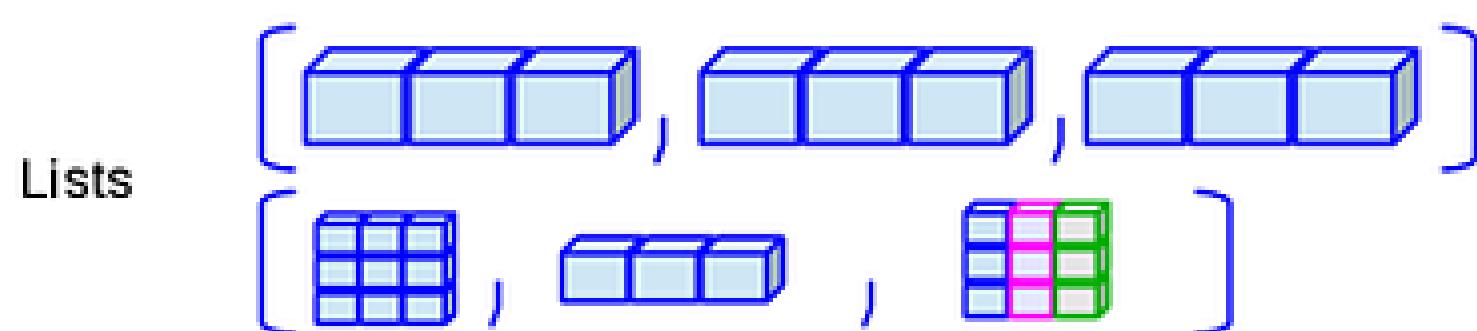
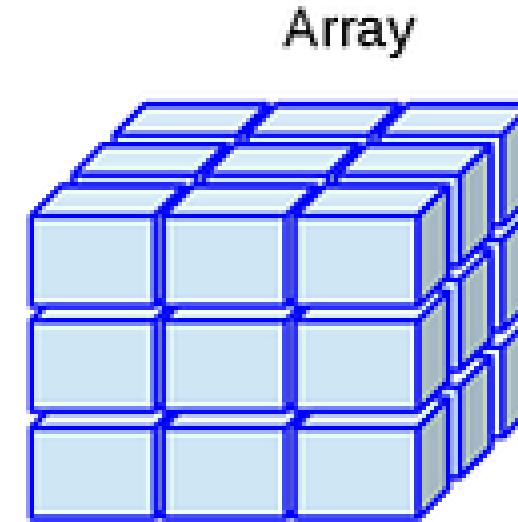
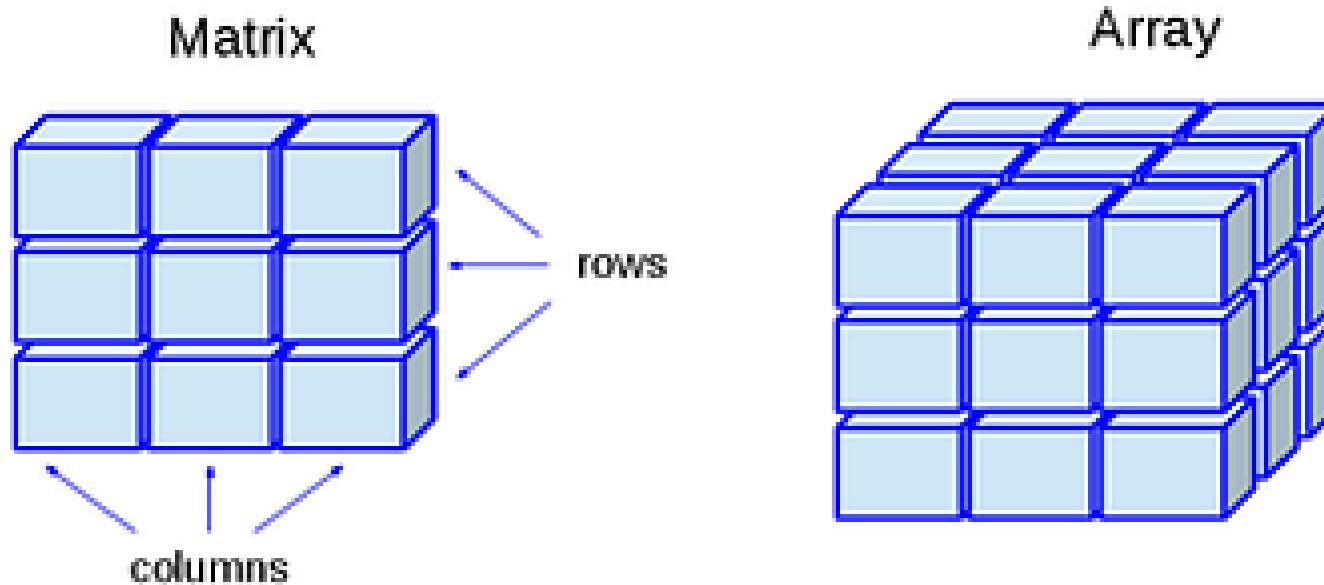
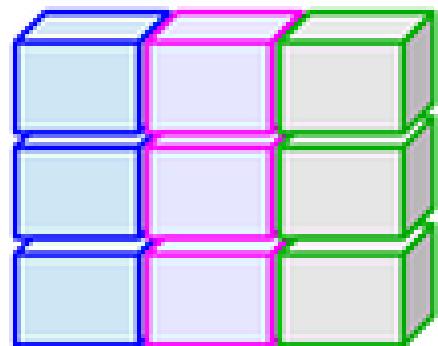
```
x<- c("a", "b", "c", "d", "e")
```

```
y<- c("a", "d", "f")
```

```
y %in% x
```



**Data Frame
(Table)**



```
data <- read.csv('data.csv')
data
print(data)
head(data)
names(data)
class(data)
str(data)
summary(data)
c("Boston", "Dakota", "Washington") %in% data$state
```

data frame

| | | |
|---|-----|-------|
| 1 | "R" | TRUE |
| 2 | "S" | FALSE |
| 3 | "T" | TRUE |

numeric character logical

A on one?



3. Manipulando dados

```
ls()  
rm(list = ls())  
tibble(data)  
data <- tibble(data)  
data  
data$state  
pull(data, state)  
data[, 1]  
data$state == pull(data, state)  
data[, 1] == pull(data, state)
```

```
identical(data$state, pull(data, state))
identical(data[, 1], pull(data, state))
data$state %>% class
pull(data, state) %>% class
data[, 1] %>% class
data[1, ]
data[1, 1]
data[c(1, 2, 3), ]
data[, 2:4]
```

```
data$state %>% class
```

```
pull(data, state) %>% class
```

```
data[, 1] %>% class
```

```
data[1, ]
```

```
data[1, 1]
```

```
data[c(1, 2, 3), ]
```

```
data[, 2:4]
```

0% > 0%

```
object %>% function()
```

```
function(object)
```

```
object %>% function(arg2)
```

```
function(object, arg2)
```

```
function(arg1, object)
```

```
object %>% function(arg1, .)
```

```
data$region  
table(data$region)  
data %>% mutate(region = as_factor(region))  
data %>% summary  
data %>% str  
data <- data %>%  
  mutate(region = as_factor(region))  
data %>% summary  
data %>% str
```

```
data %>%  
  mutate(dens = total/population)  
  
data %>%  
  mutate(dens = total/population * 100000)  
  
data %>%  
  mutate(dens = total/population * 100000,  
         seguro = dens < 1) %>%  
  filter(seguro == TRUE) %>%  
  arrange(dens)
```

```
data %>%  
  mutate(dens = total/population * 100000,  
        seguro = dens < 1) %>%  
  arrange(dens) %>%  
  select(state, dens, seguro)
```

```
data %>%  
  reframe(population = sum(population),  
          total = sum(total))
```

```
data %>%  
  group_by(region) %>%  
  reframe(n = n(),  
          pop.media = mean(population),  
          tot.population = sum(population),  
          tot.total = sum(total)) %>%  
  mutate(dens = tot.total/tot.population*100000,  
         seguro = if_else(dens < 3, 'seguro', 'inseguro'))
```

4. Dataviz

Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.
 Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite
 Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Séguir, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk en Mohilow et se rejoignaient vers Orscha en Wilebsk, avaient toujours marché avec l'armée.

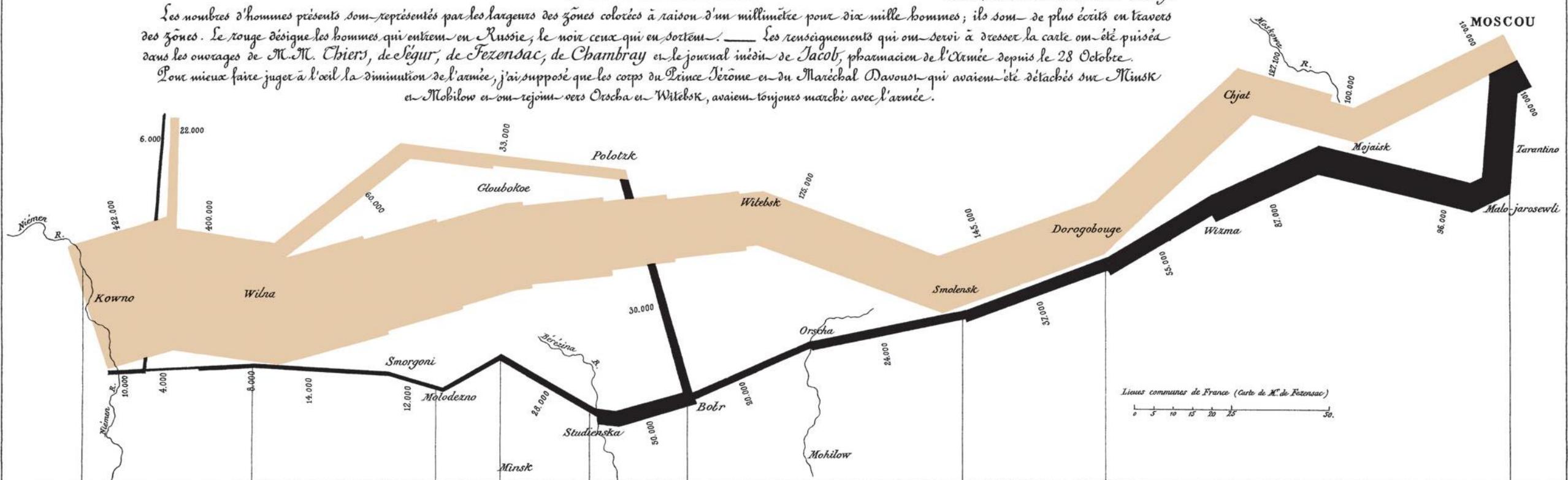
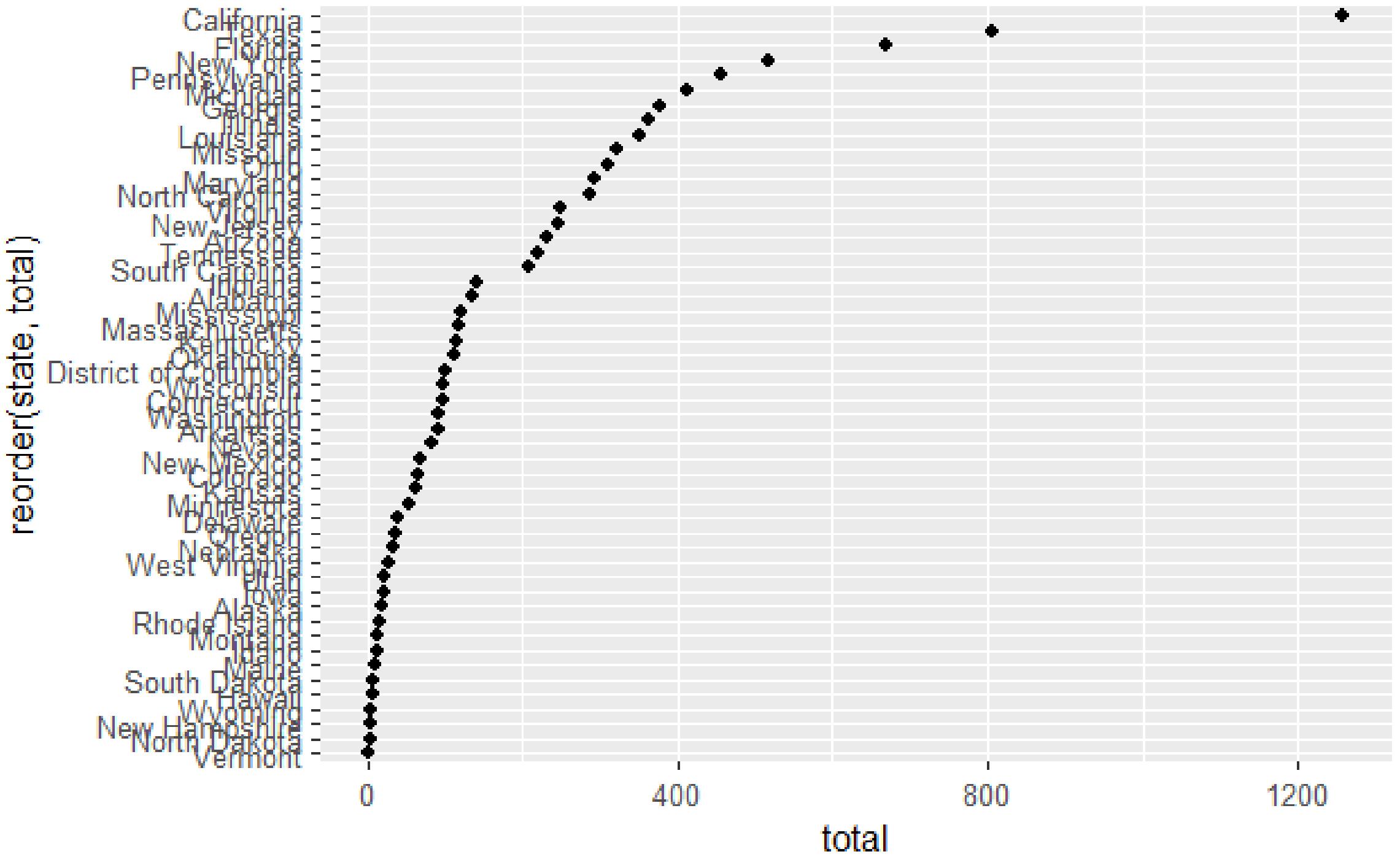


TABLEAU CRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les cosaques passent au galop
le Niemen gelé.

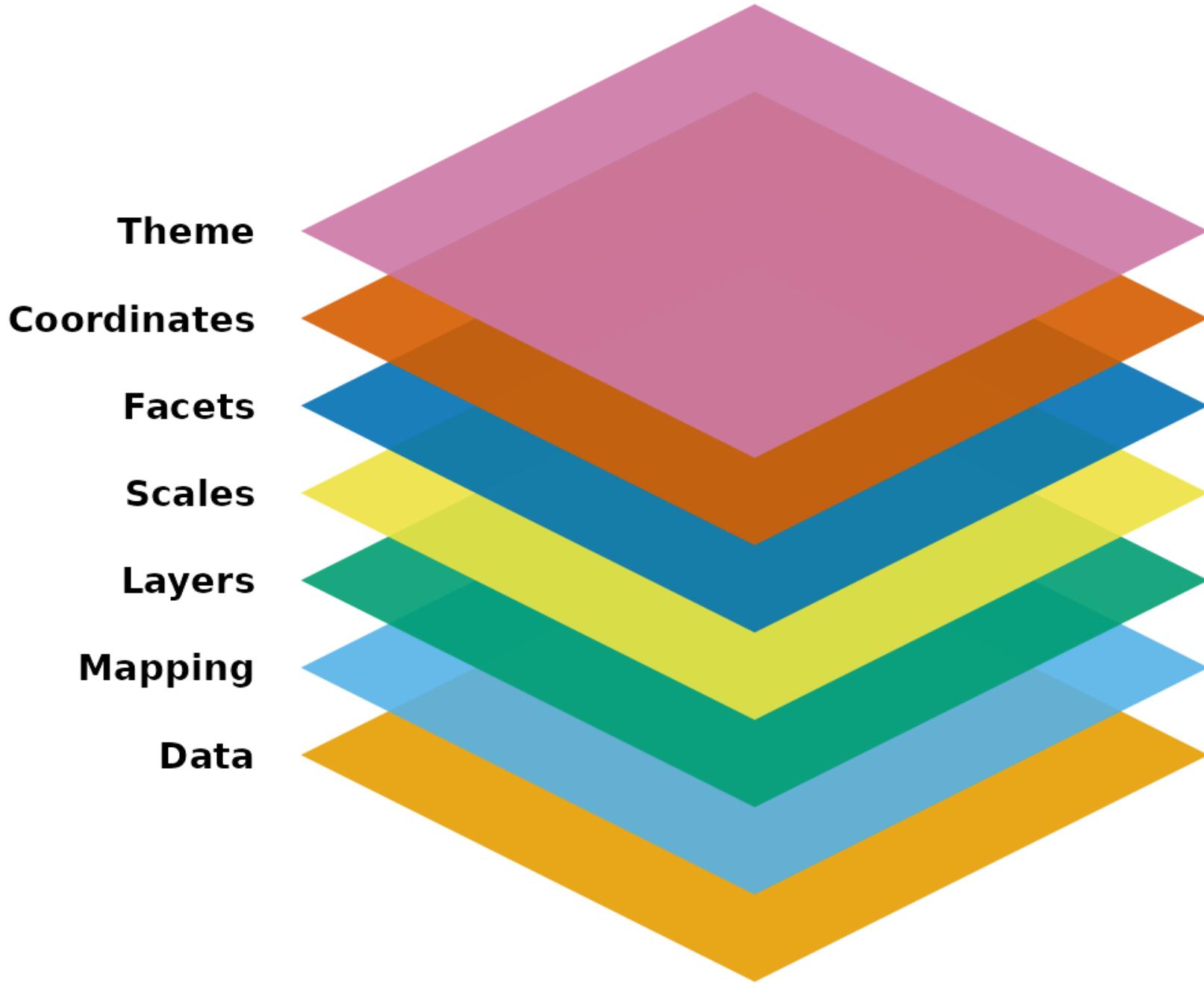


```
data %>%  
  ggplot() +  
    geom_point(aes(total, reorder(state,  
total)))
```



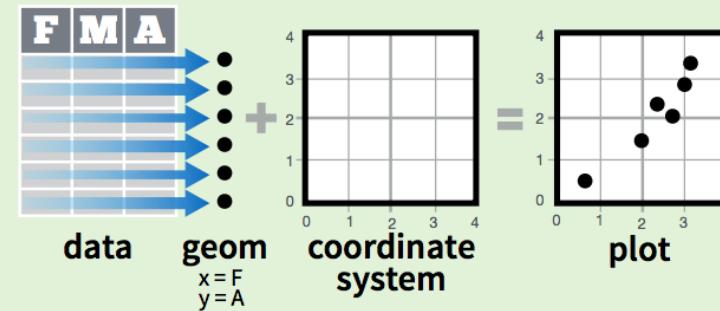
O **ggplot2** é um pacote R que usa uma estrutura conceitual baseada na gramática de gráficos para produzir plots.

Isso permite que você “descreva” um gráfico a partir de elementos combináveis.



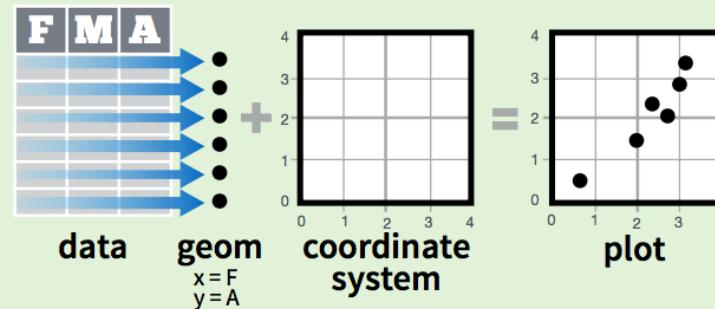
Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.

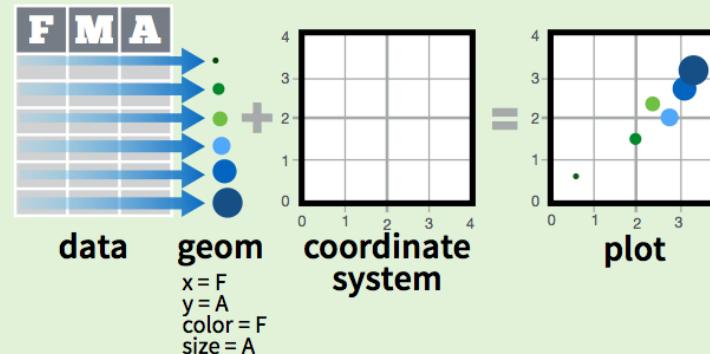


Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.

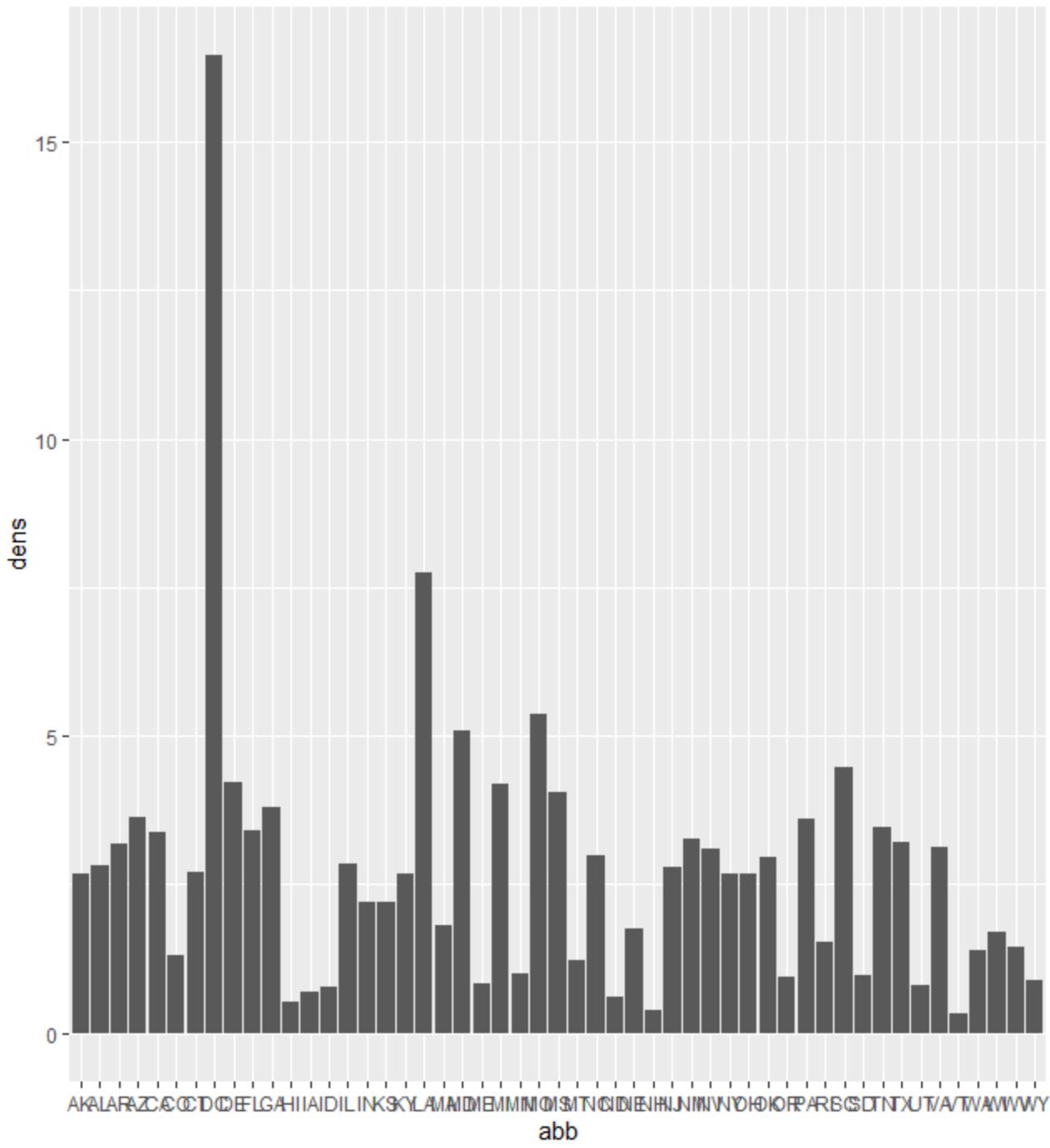


To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.

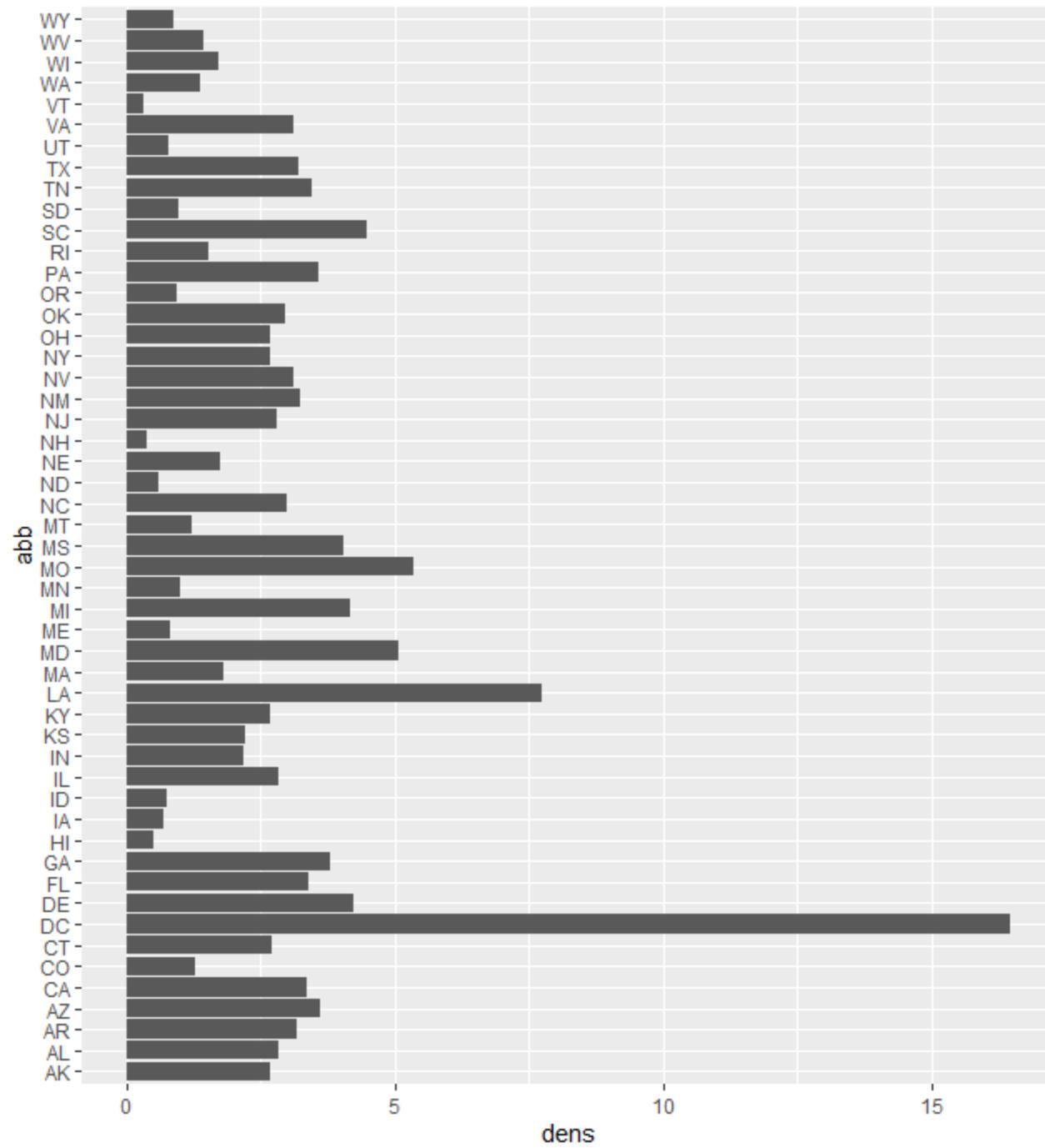


```
data %>%  
  ggplot()
```

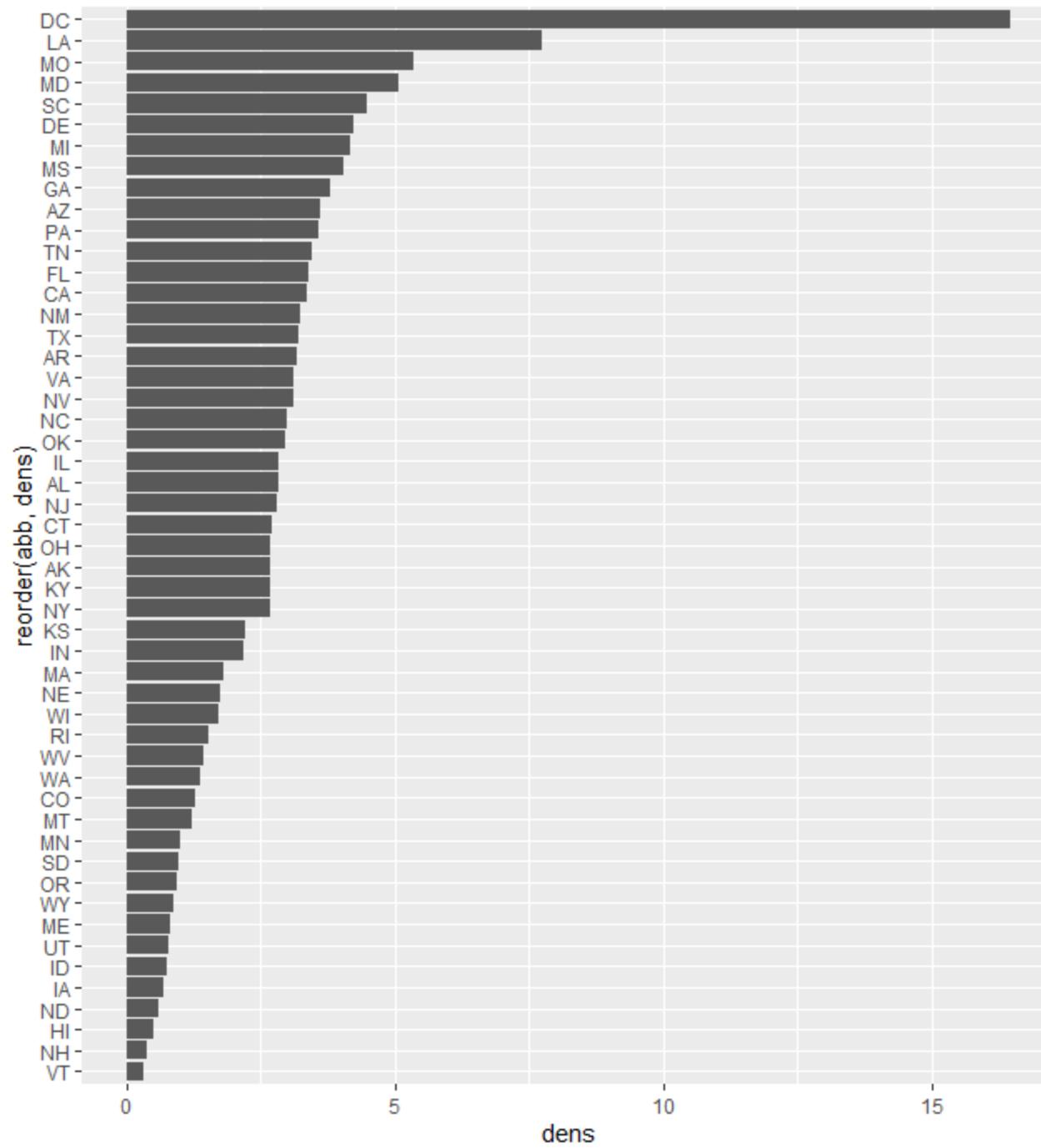
```
data %>%  
  ggplot() +  
  geom_col(aes(x = abb,  
               y = dens))
```



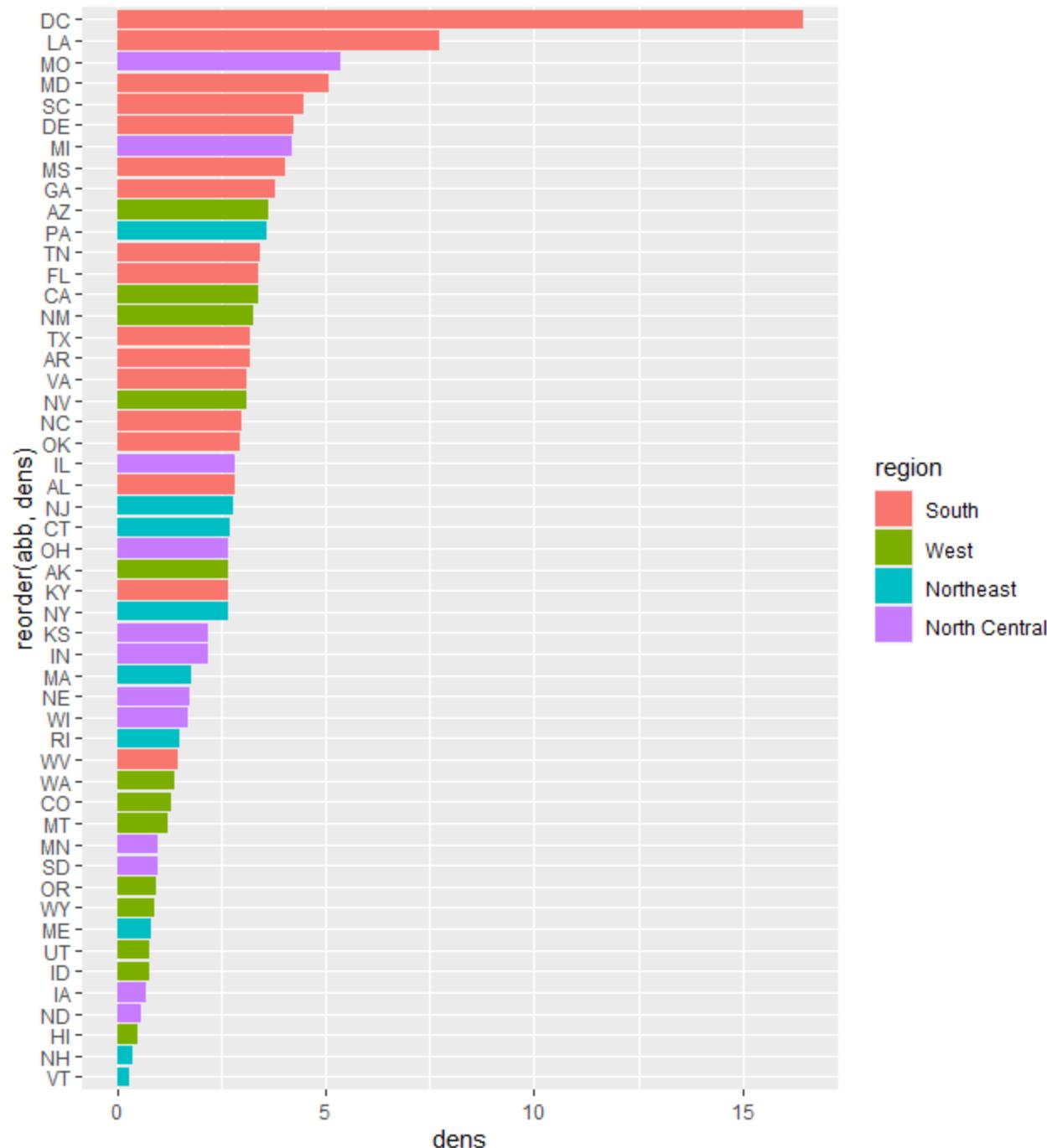
```
data %>%  
  ggplot() +  
  geom_col(aes(y = abb,  
               x = dens))
```



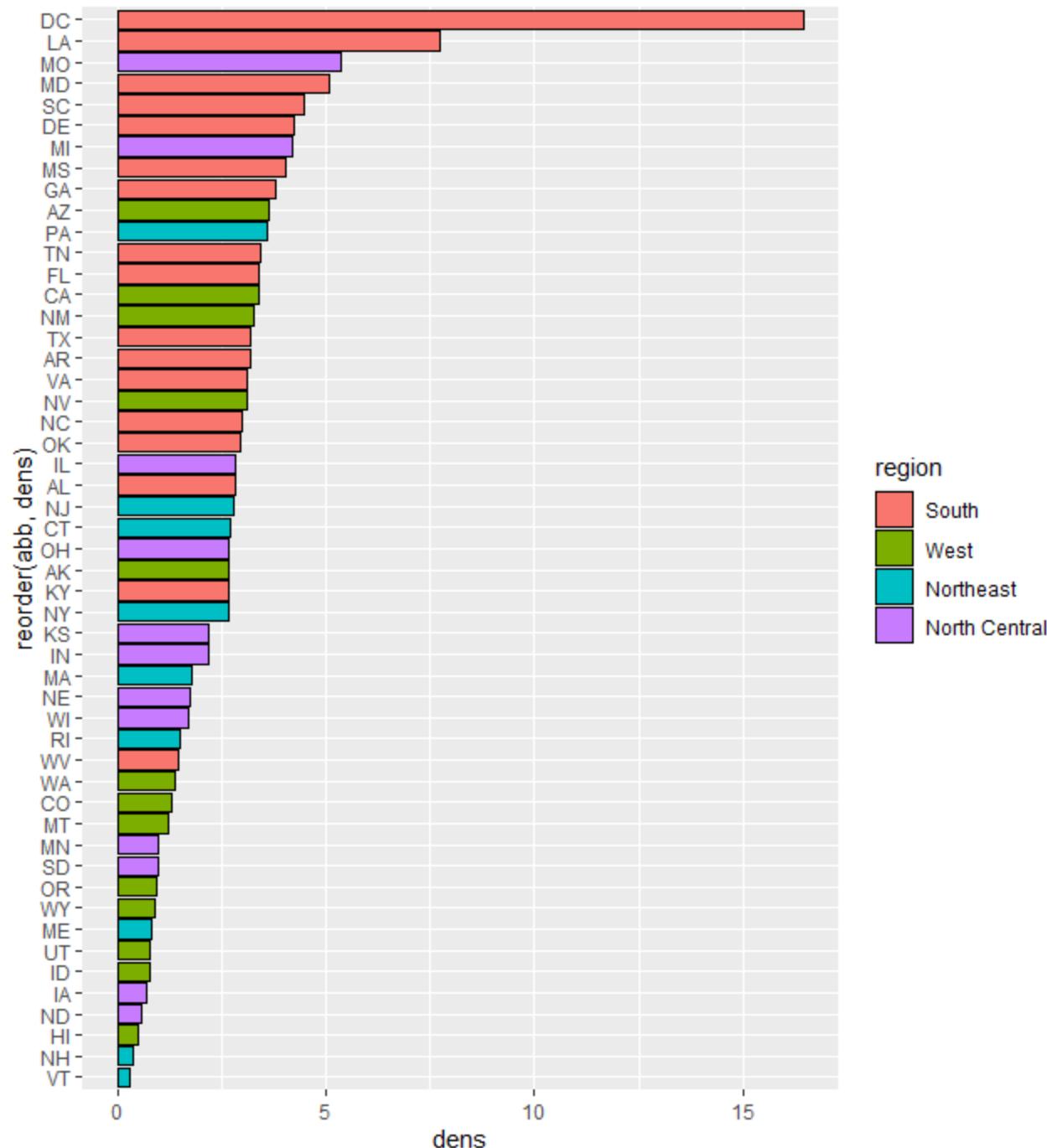
```
data %>%  
  ggplot() +  
  geom_col(aes(y = reorder(abb, dens),  
             x = dens))
```



```
data %>%  
  ggplot() +  
  geom_col(aes(y = reorder(abb, dens),  
               x = dens,  
               fill = region)
```

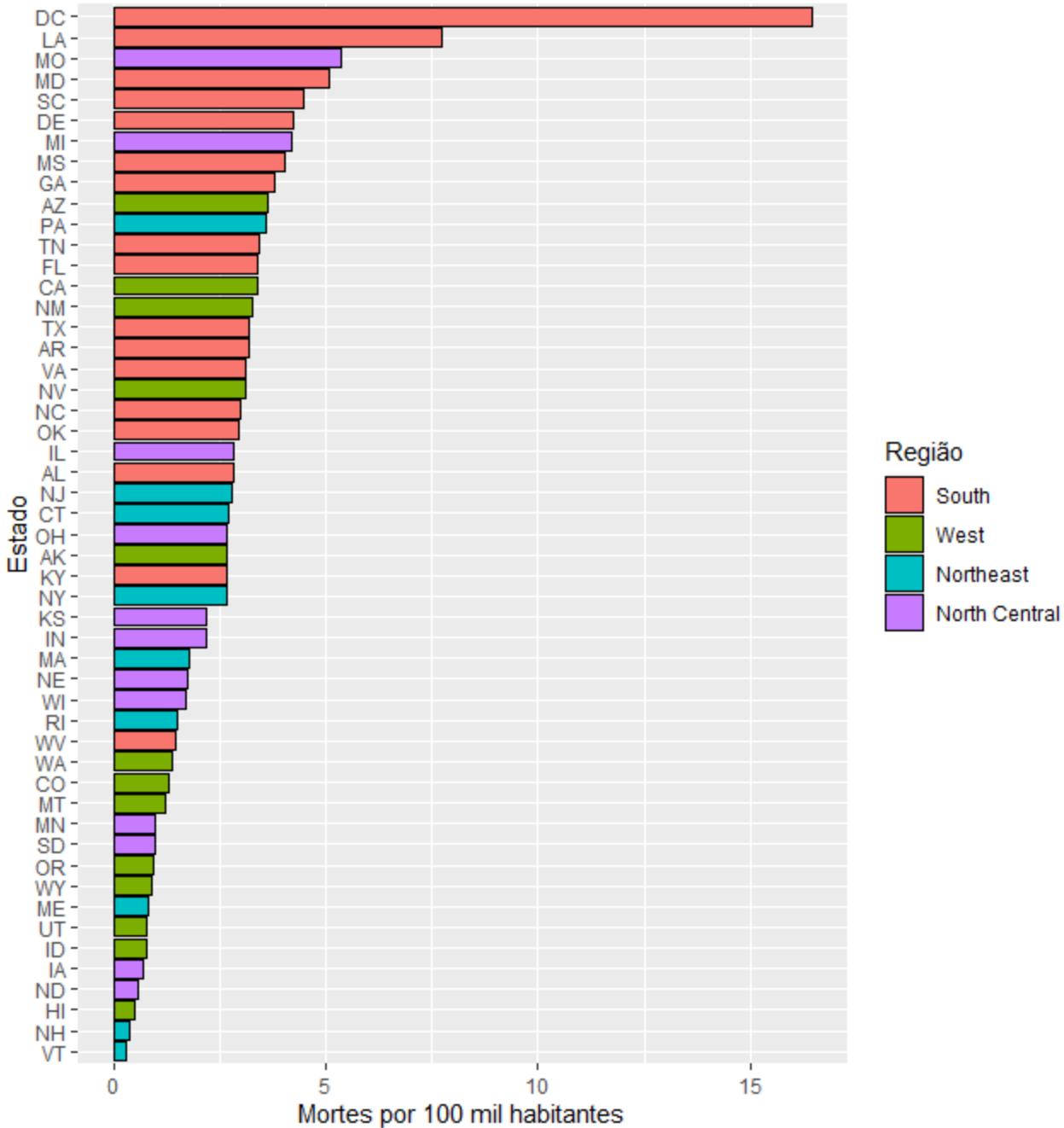


```
data %>%  
  ggplot() +  
  geom_col(aes(y = reorder(abb, dens),  
               x = dens,  
               fill = region),  
color = 'black')
```



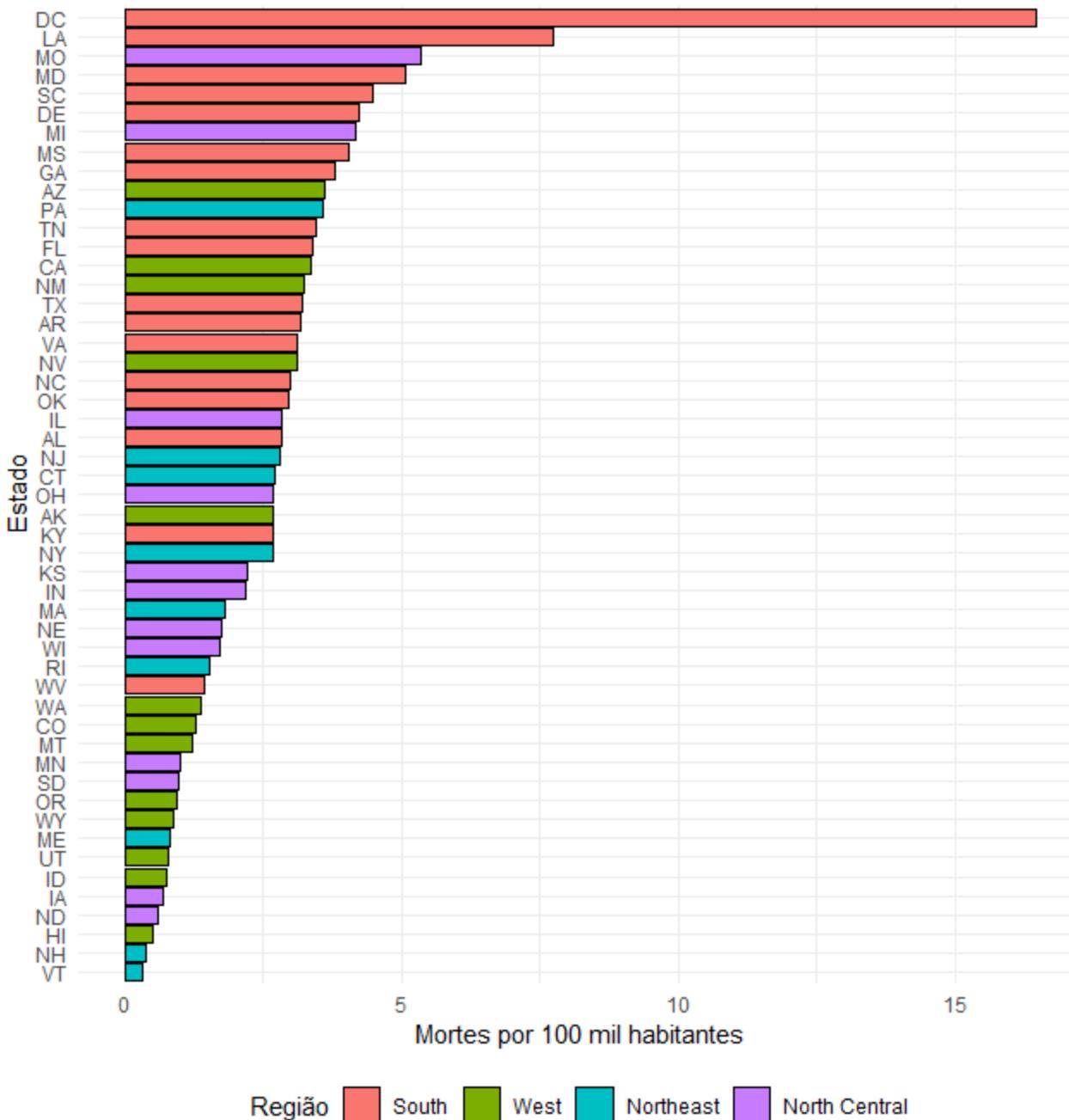
```
data %>%
  ggplot() +
  geom_col(aes(y = reorder(abb, dens),
                x = dens,
                fill = region),
            color = 'black') +
  labs(x = 'Mortes por 100 mil habitantes',
       y = 'Estado',
       title = 'Mortes por arma de fogo nos EUA em 2010',
       fill = 'Região)
```

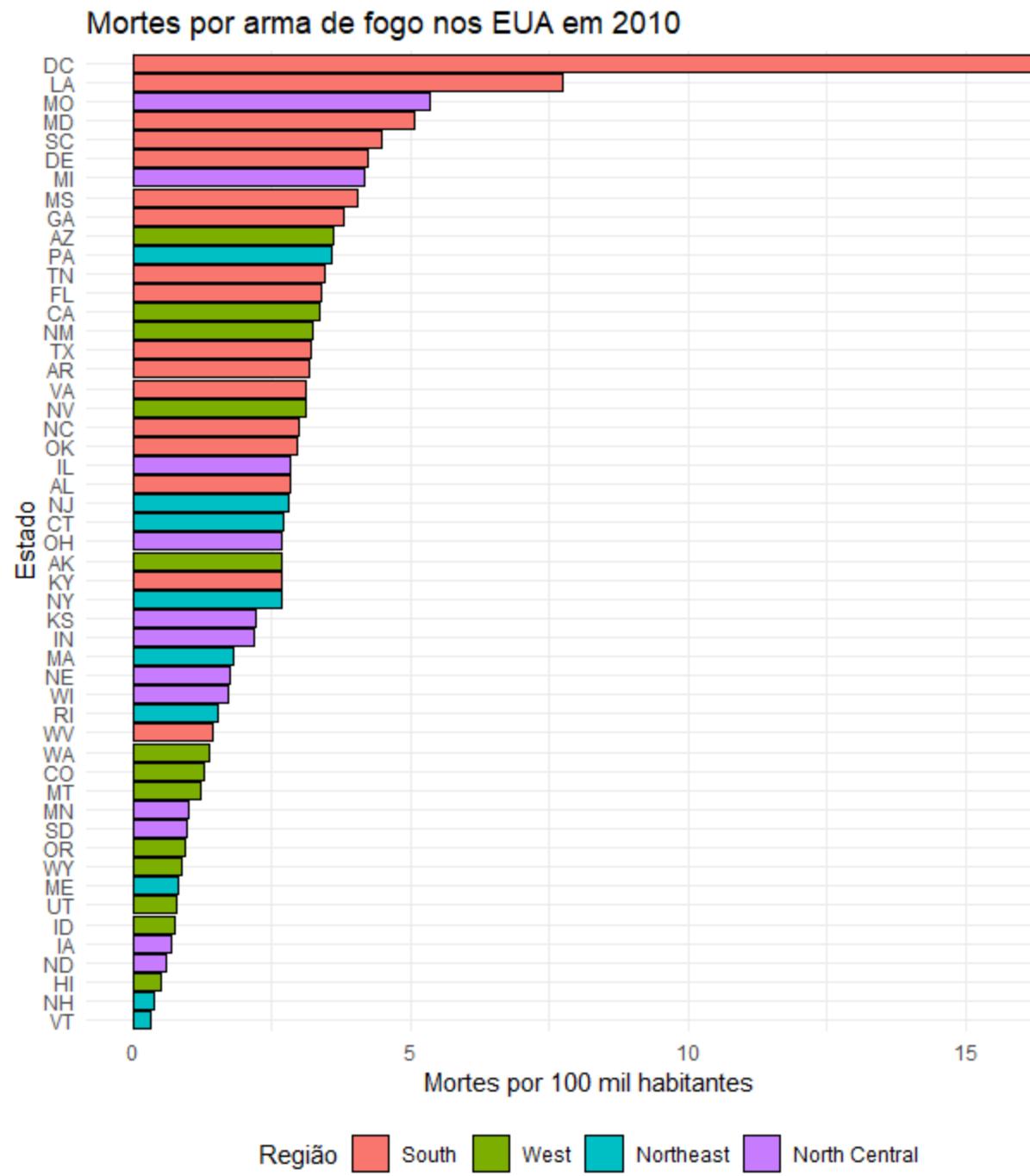
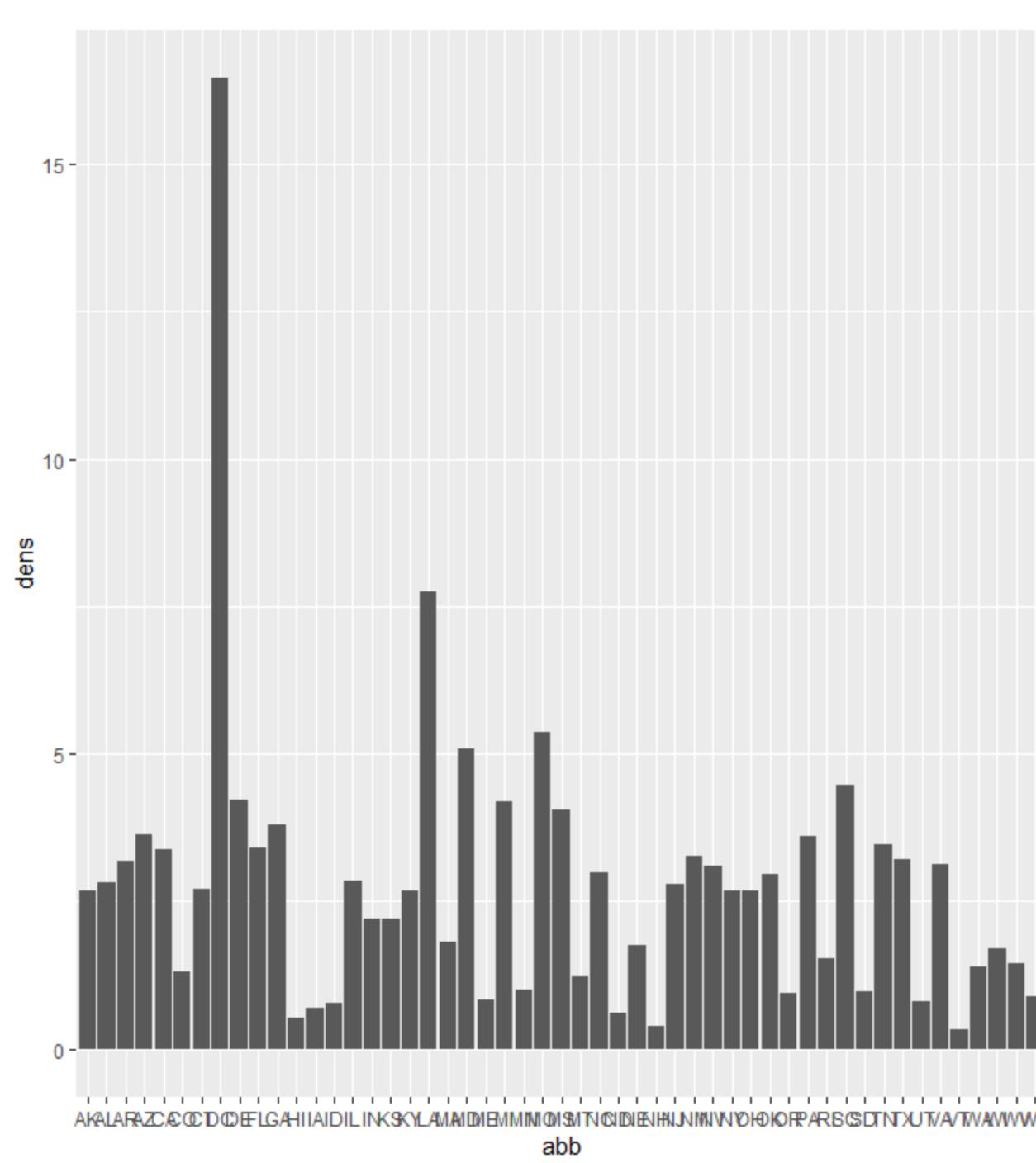
Mortes por arma de fogo nos EUA em 2010



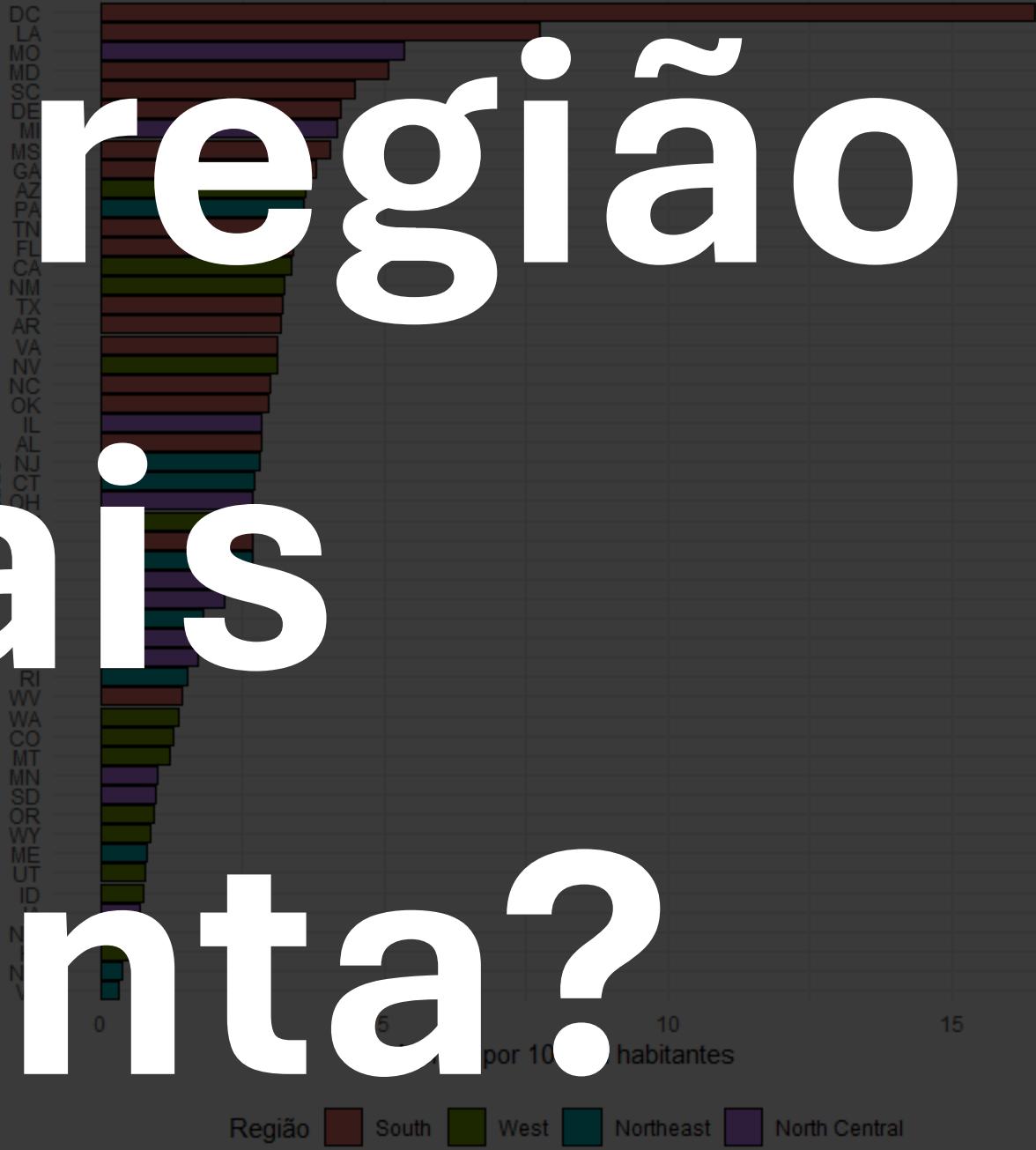
```
data %>%
ggplot() +
geom_col(aes(y = reorder(abb, dens),
             x = dens,
             fill = region),
color = 'black') +
labs(x = 'Mortes por 100 mil habitantes',
     y = 'Estado',
     title = 'Mortes por arma de fogo nos EUA em 2010',
     fill = 'Região') +
theme_minimal() +
theme(legend.position = 'bottom')
```

Mortes por arma de fogo nos EUA em 2010



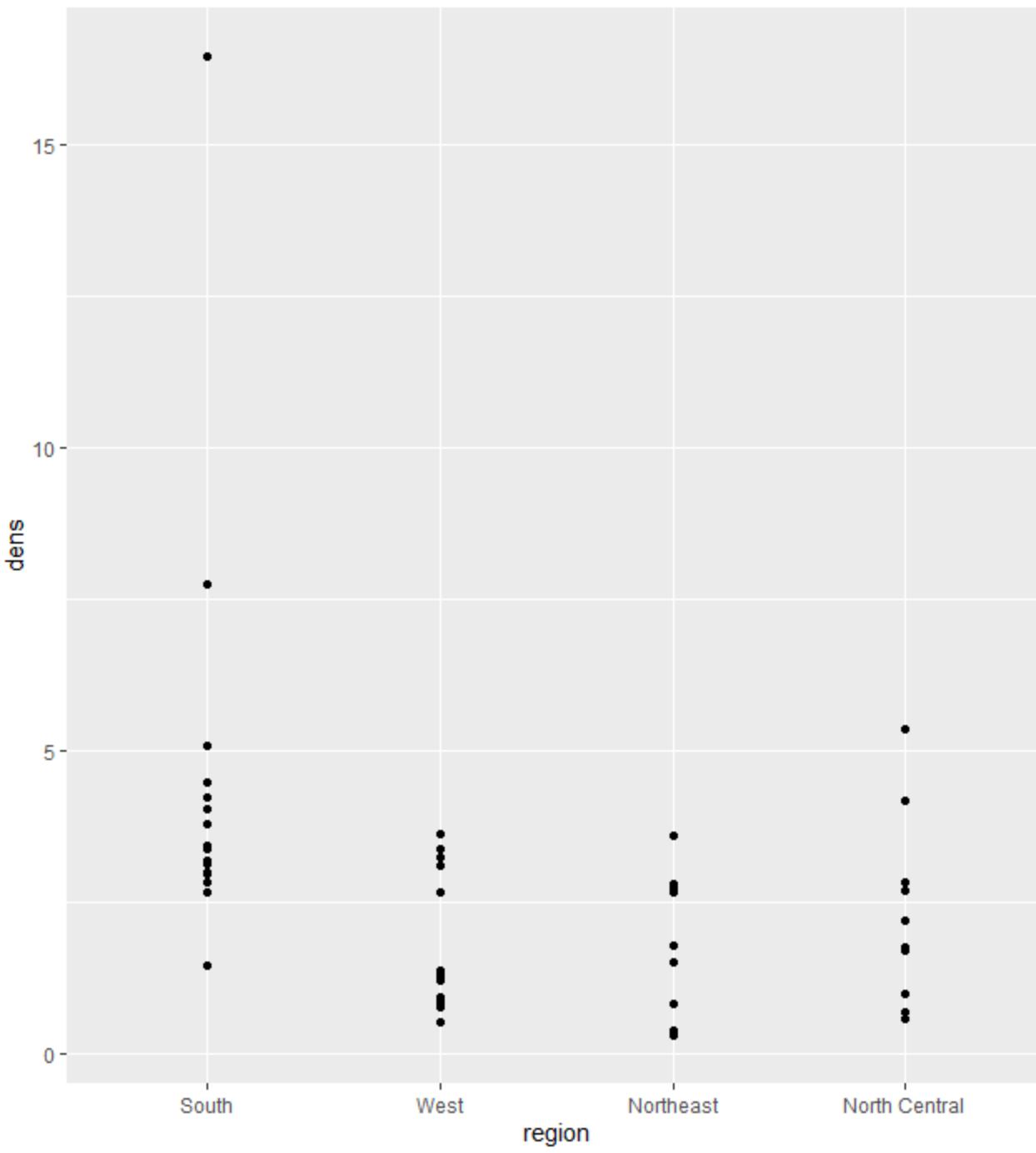


Mortes por arma de fogo nos EUA em 2010

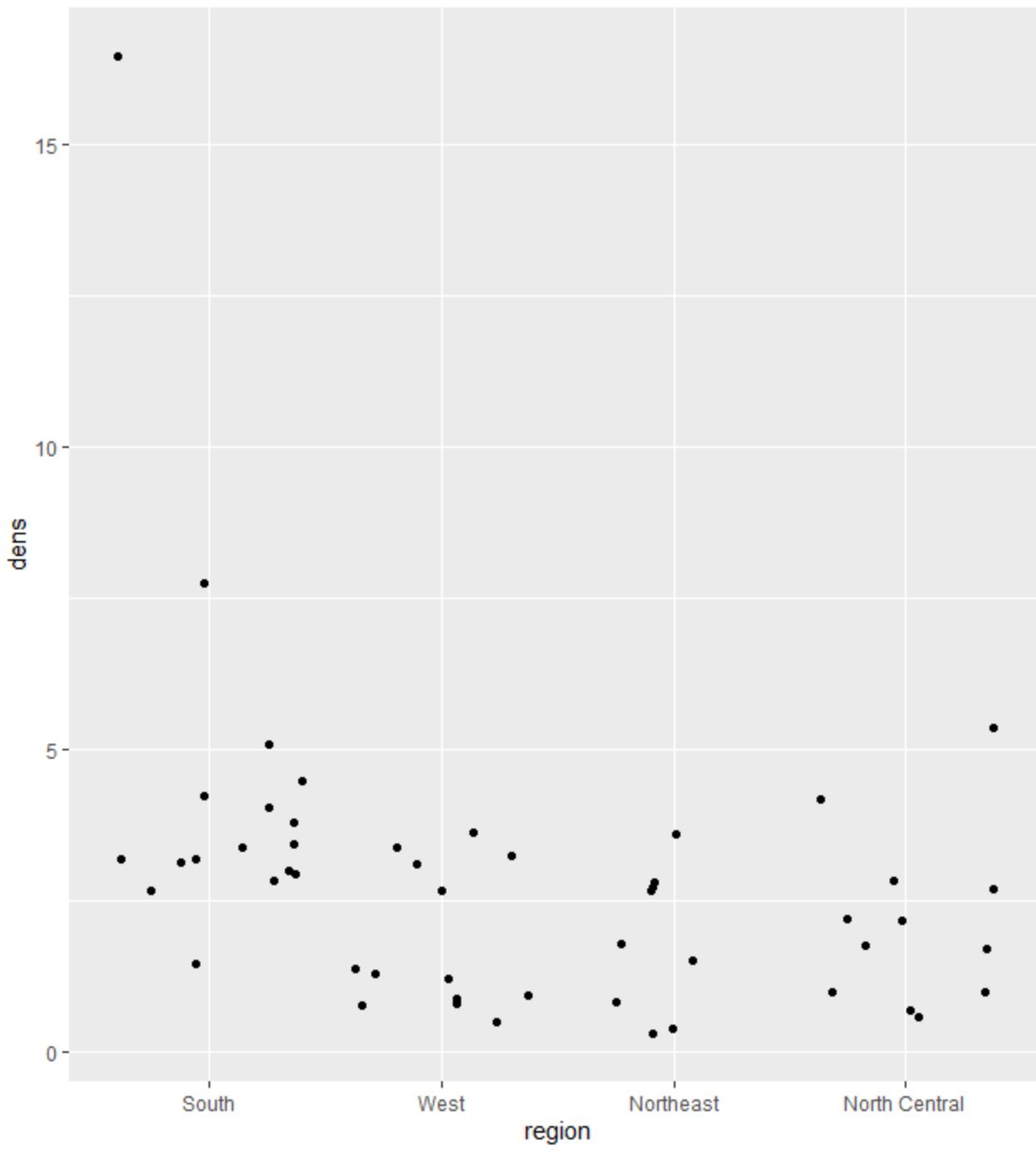


Qual a região
mais
violenta?

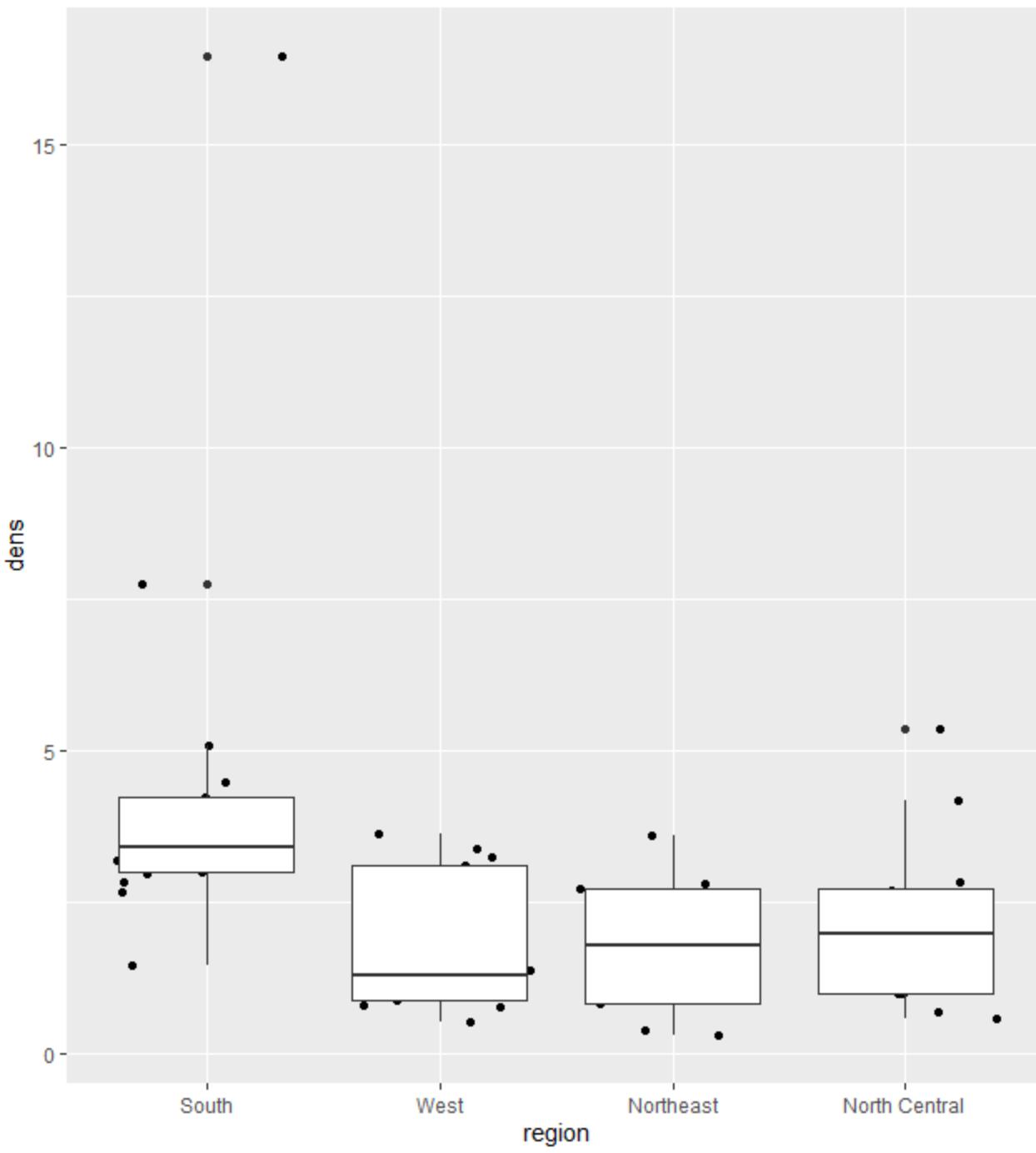
```
data %>%  
ggplot() +  
geom_point(aes(x = region, y = dens))
```



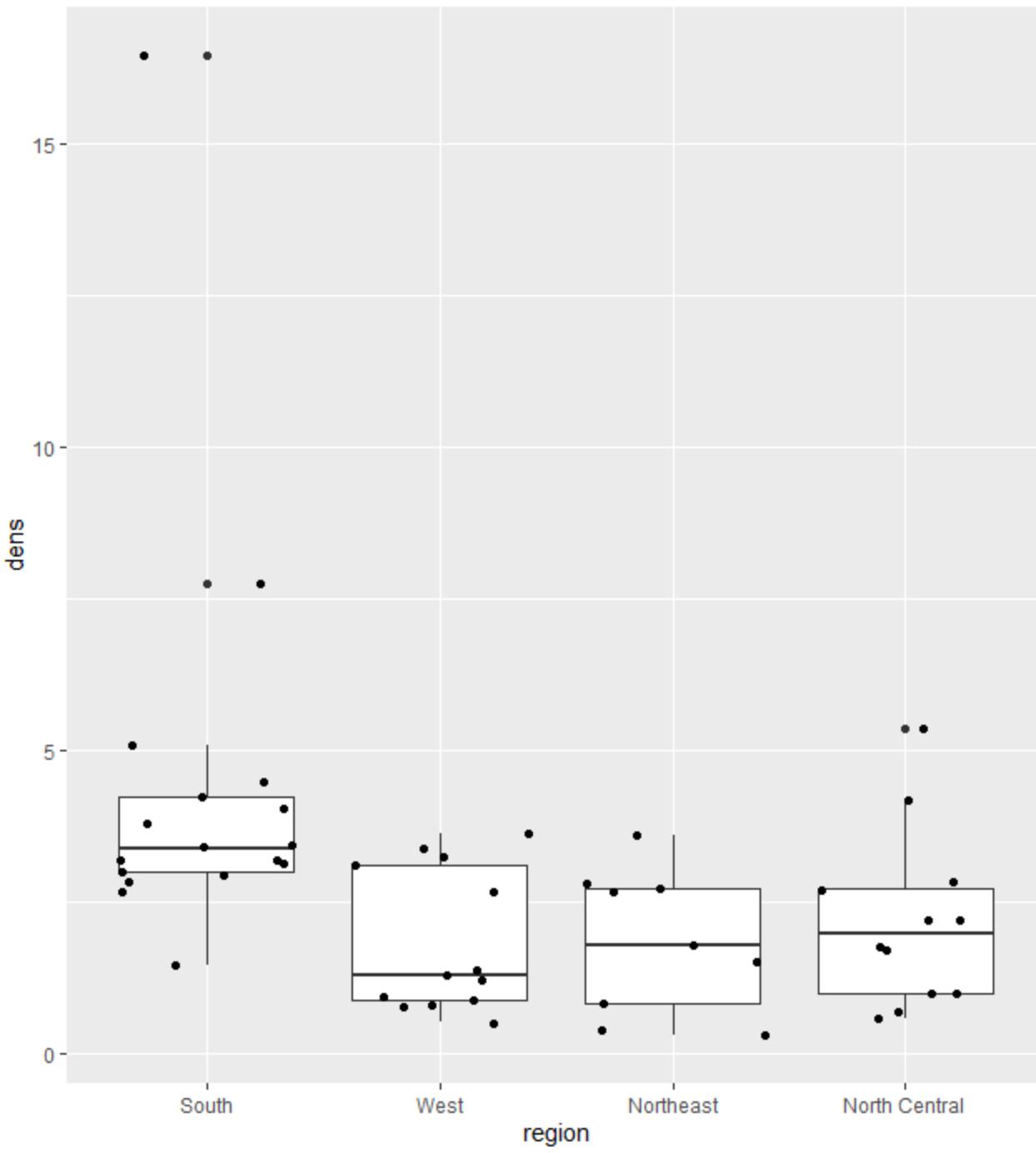
```
data %>%  
ggplot() +  
geom_jitter(aes(x = region, y = dens))
```



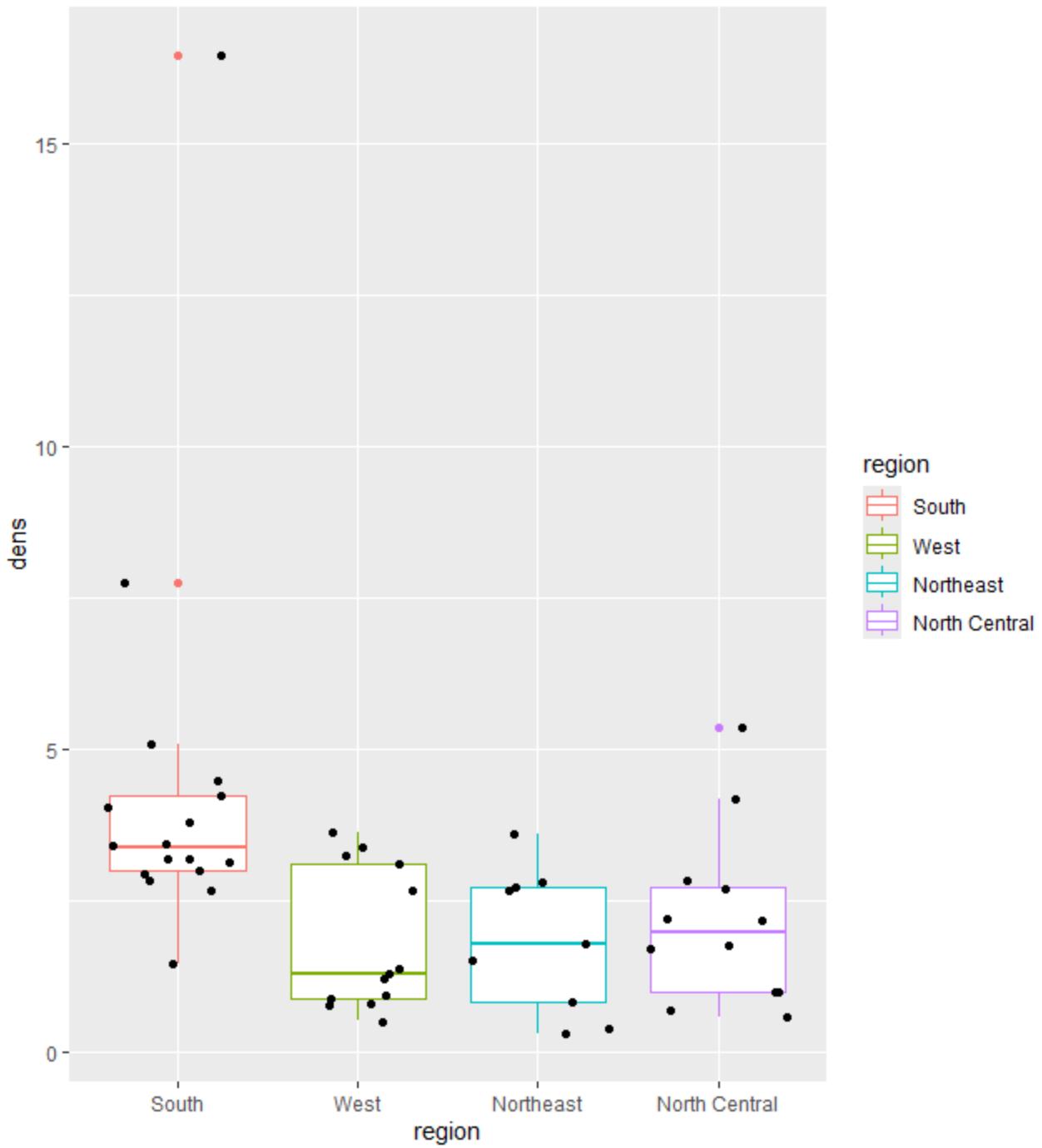
```
data %>%  
ggplot() +  
geom_jitter(aes(x = region, y = dens)) +  
geom_boxplot(aes(y = dens, x = region))
```



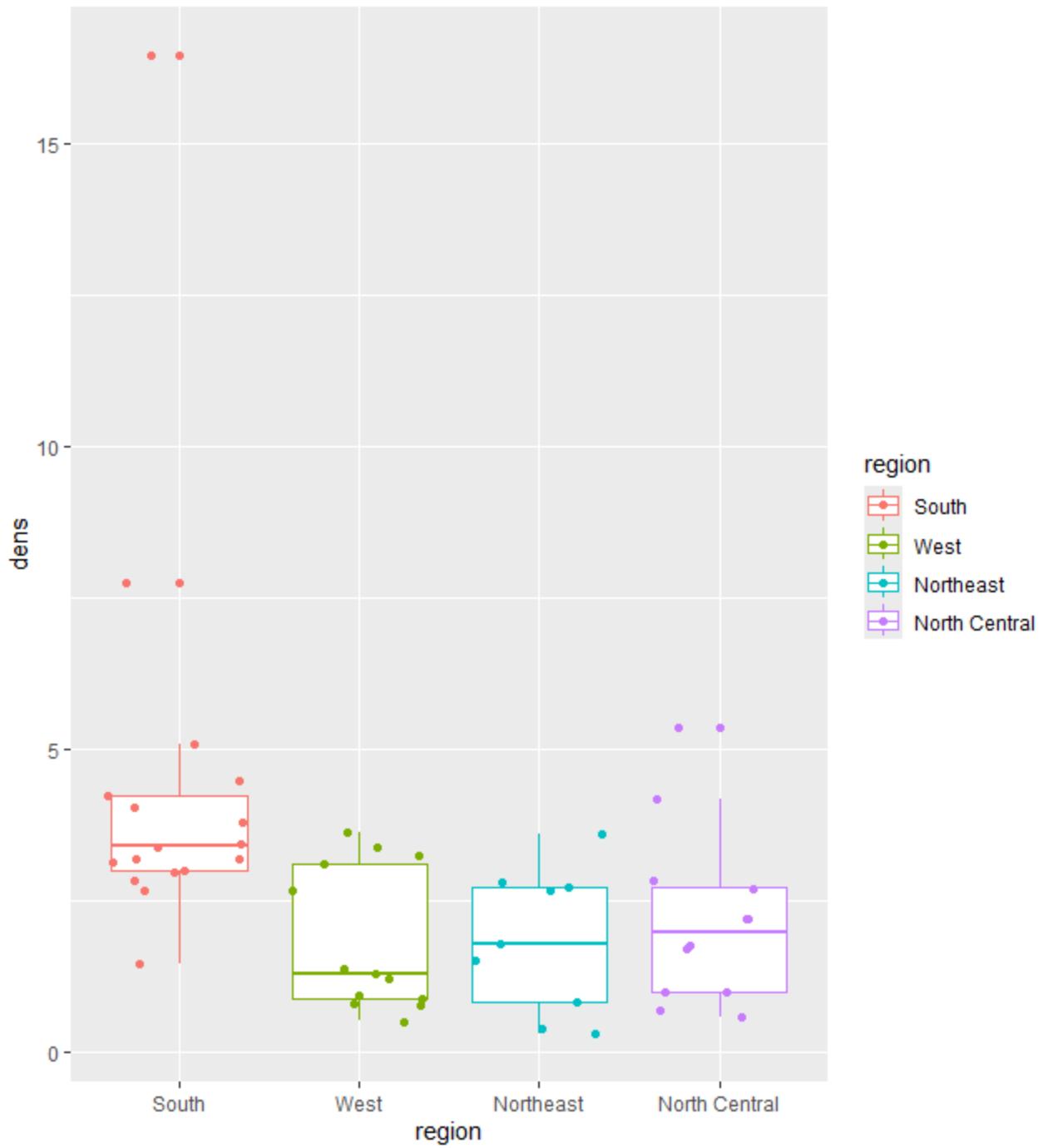
```
data %>%  
ggplot() +  
geom_boxplot(aes(y = dens, x = region)) +  
geom_jitter(aes(x = region, y = dens))
```



```
data %>%  
ggplot() +  
geom_boxplot(aes(y = dens, x = region, color = region))  
+  
geom_jitter(aes(x = region, y = dens))
```

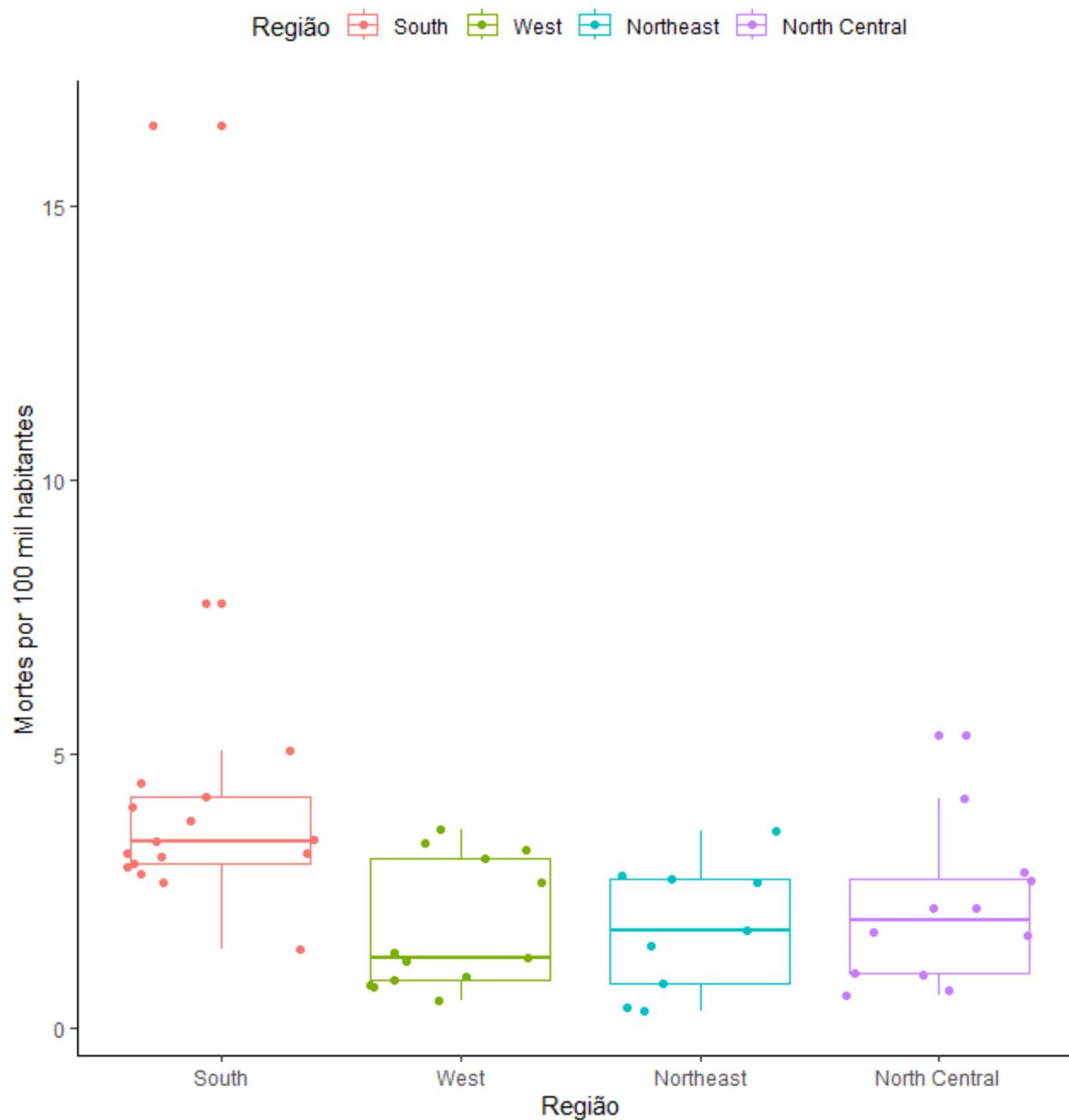


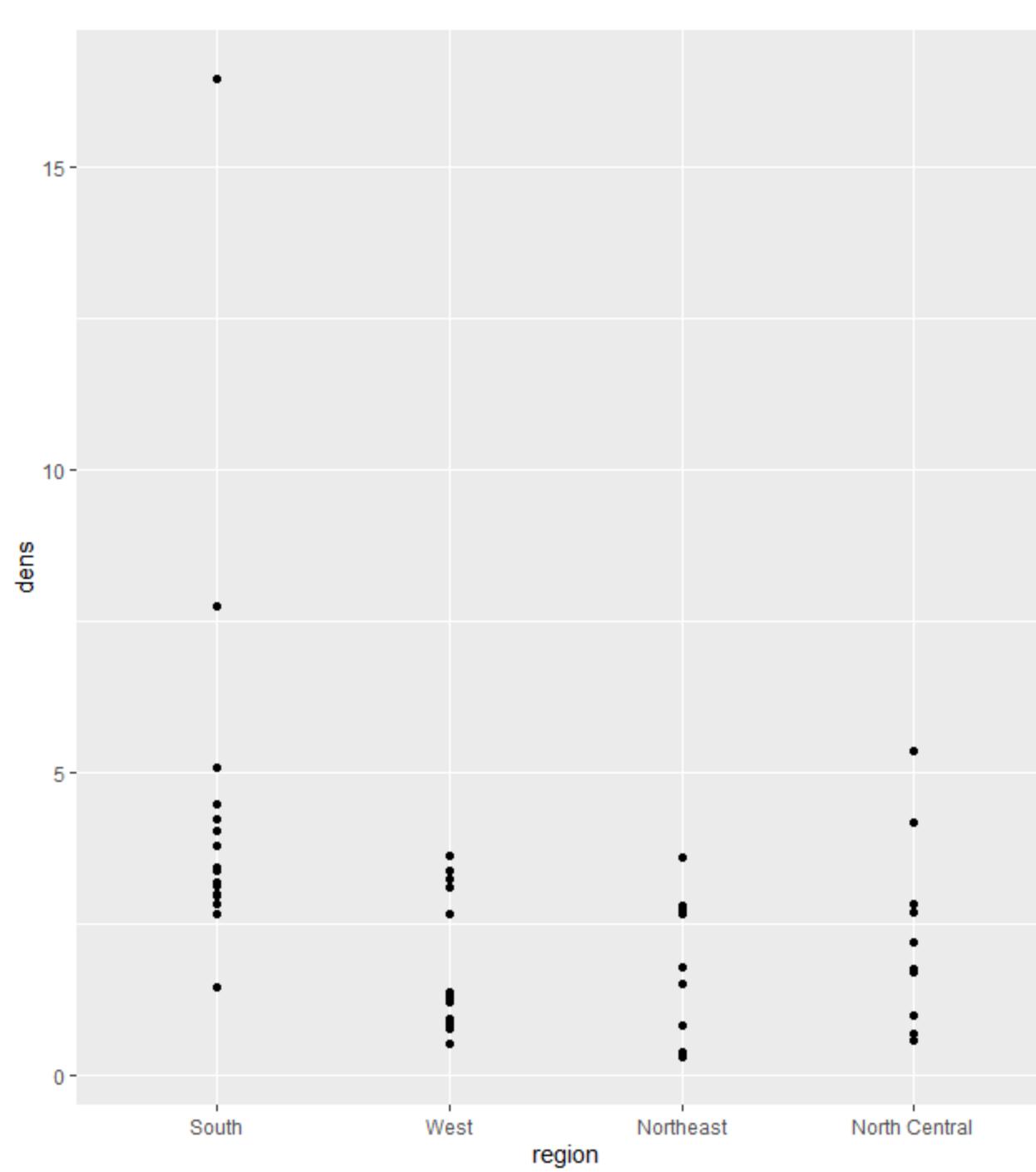
```
data %>%
ggplot() +
geom_boxplot(aes(y = dens, x = region, color = region))
+
geom_jitter(aes(x = region,
                 y = dens,
                 color = region))
```



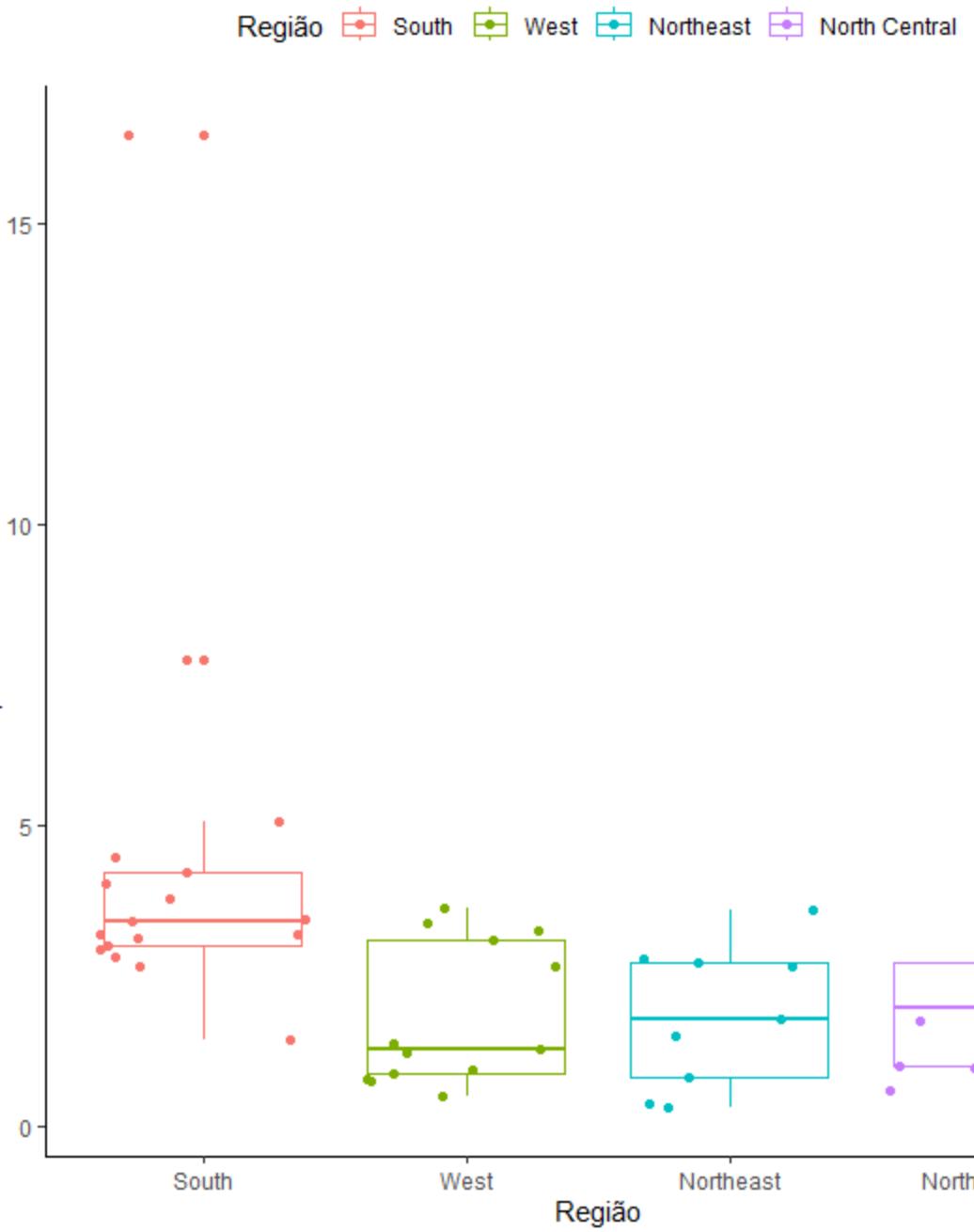
```
data %>%
ggplot() +
geom_boxplot(aes(y = dens, x = region, color = region)) +
geom_jitter(aes(x = region,
                 y = dens,
                 color = region)) +
labs(x = 'Região',
     y = 'Mortes por 100 mil habitantes',
     title = 'Mortes por arma de fogo nos EUA em 2010',
     color = 'Região') +
theme_classic() +
theme(legend.position = 'top')
```

Mortes por arma de fogo nos EUA em 2010





Mortes por arma de fogo nos EUA em 2010



Podemos salvar
o gráfico como
um objeto!

region

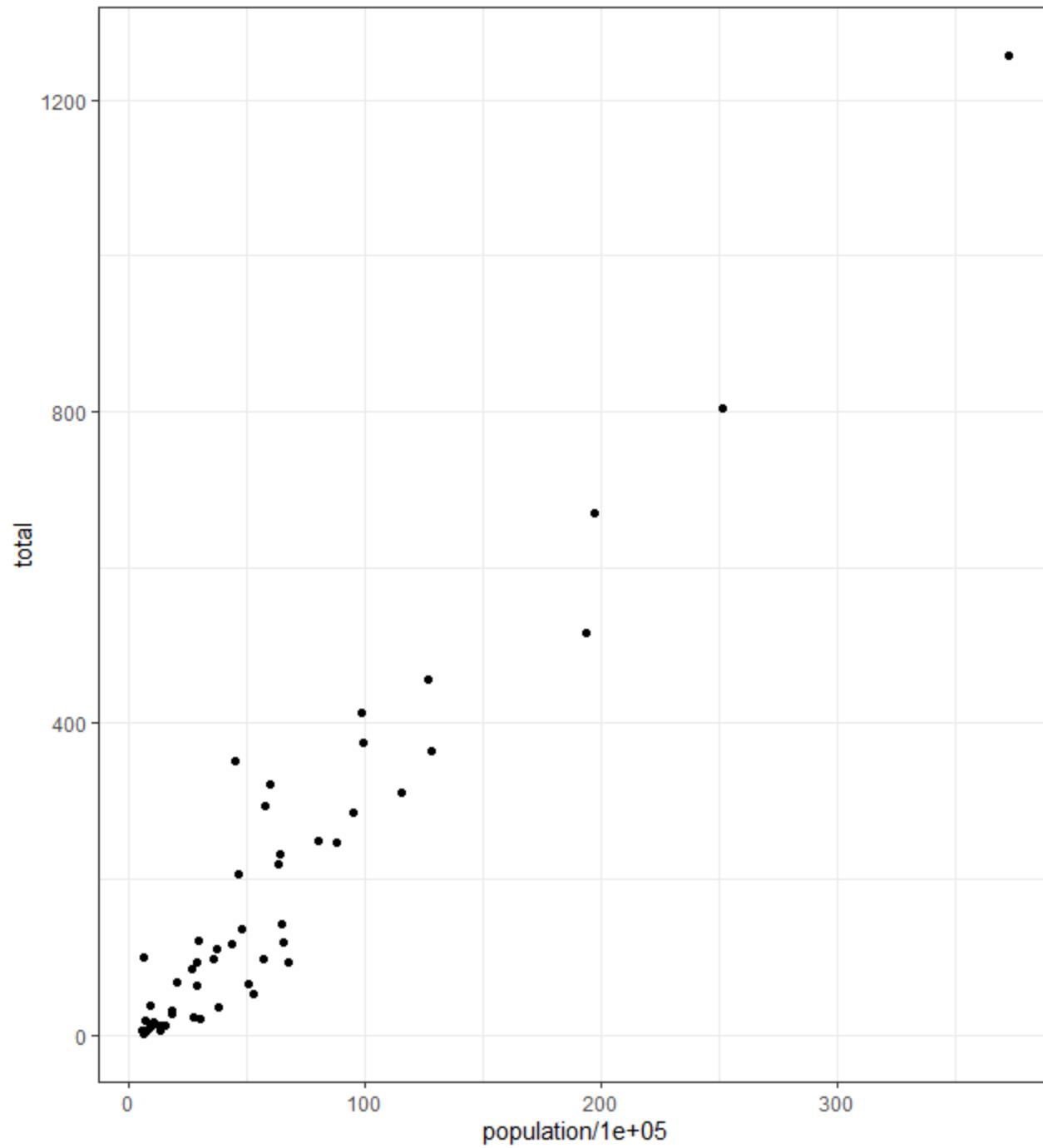
Região

```
p <- ggplot(data = data)
```

```
class(p)
```

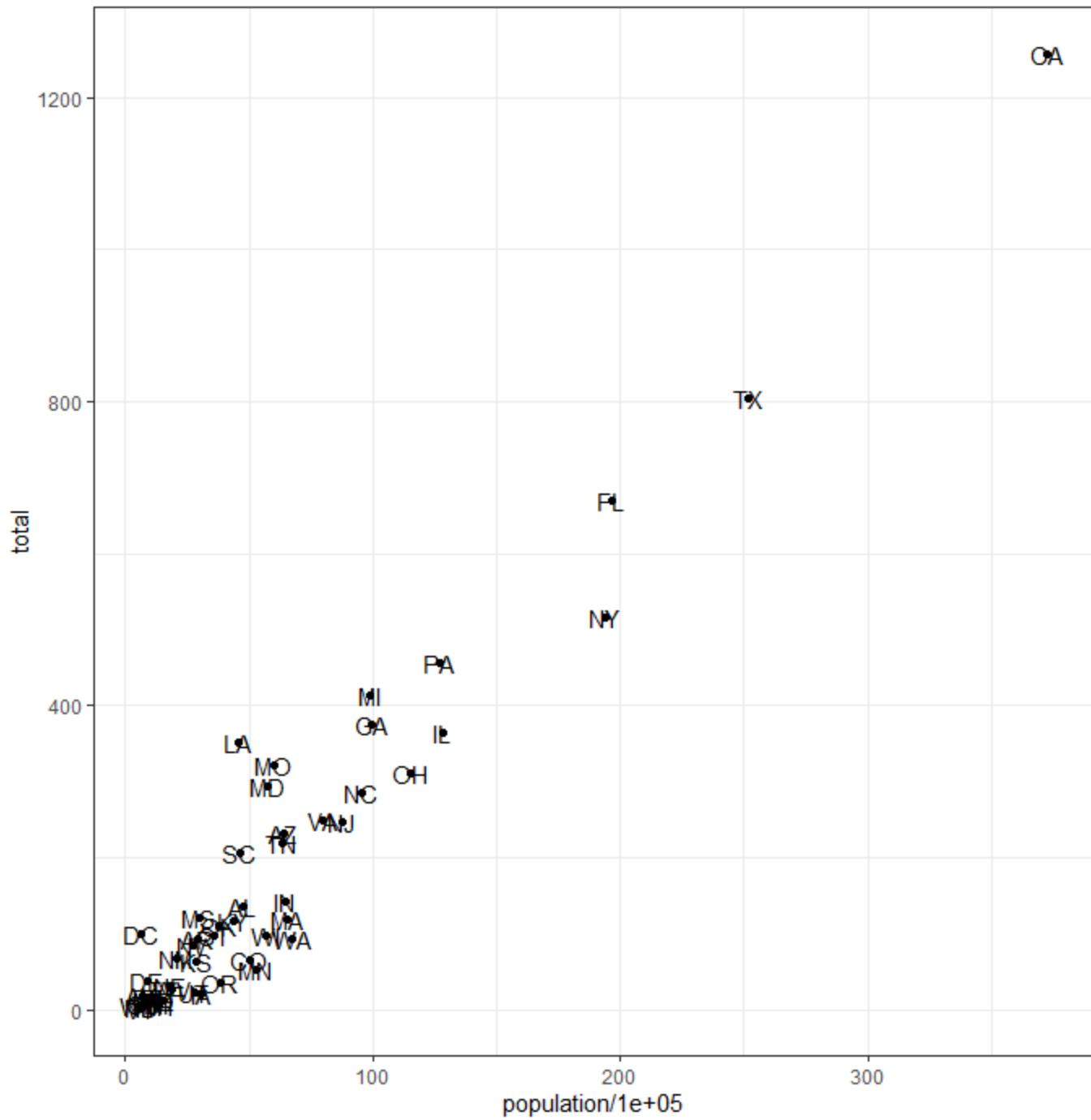
```
p +
```

```
geom_point(aes(x = population/100000, y = total))
```



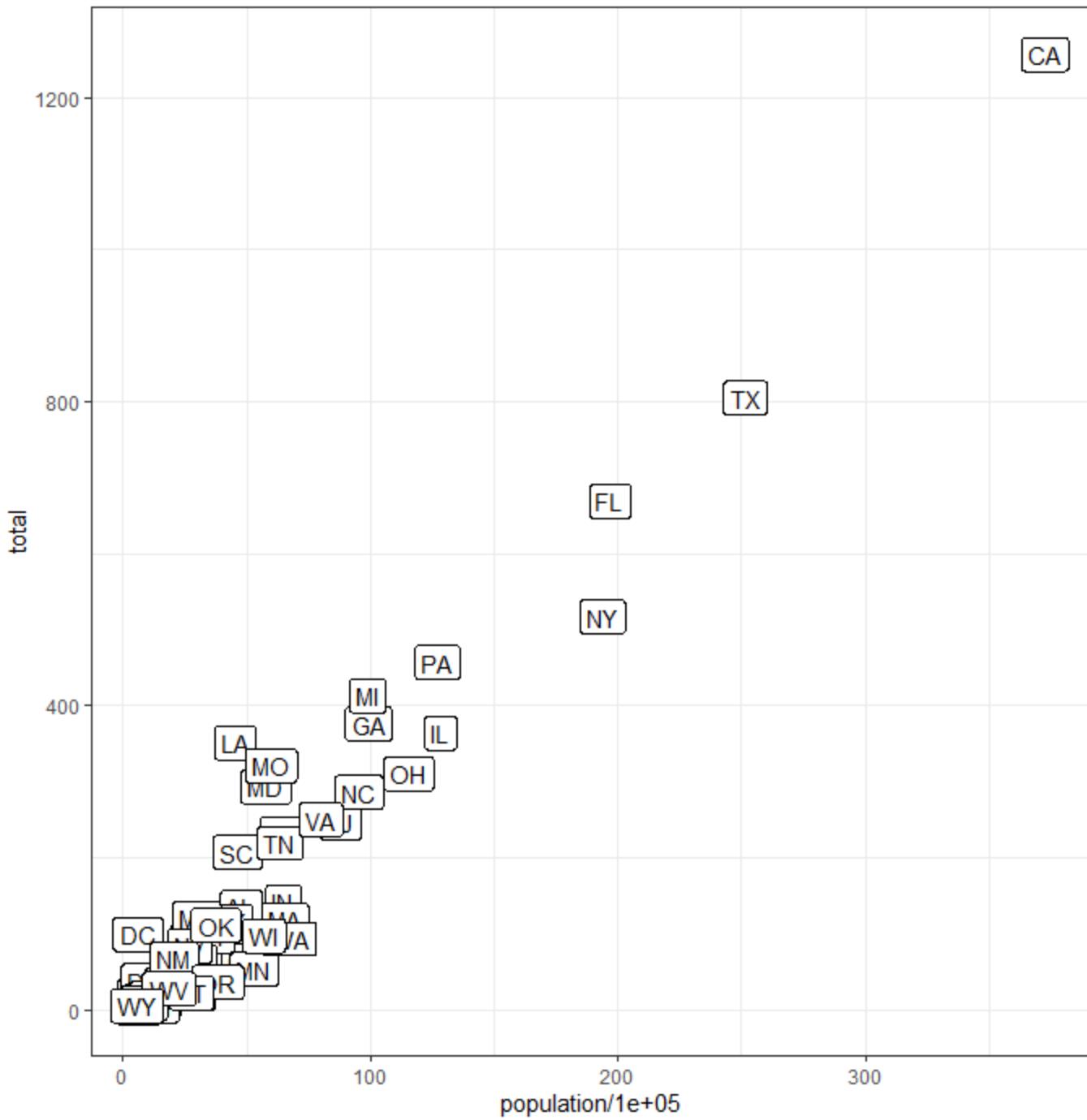
p +

```
geom_point(aes(x = population/100000, y = total), size = 3) +  
geom_text(aes(x = population/100000, y = total, label = abb))
```



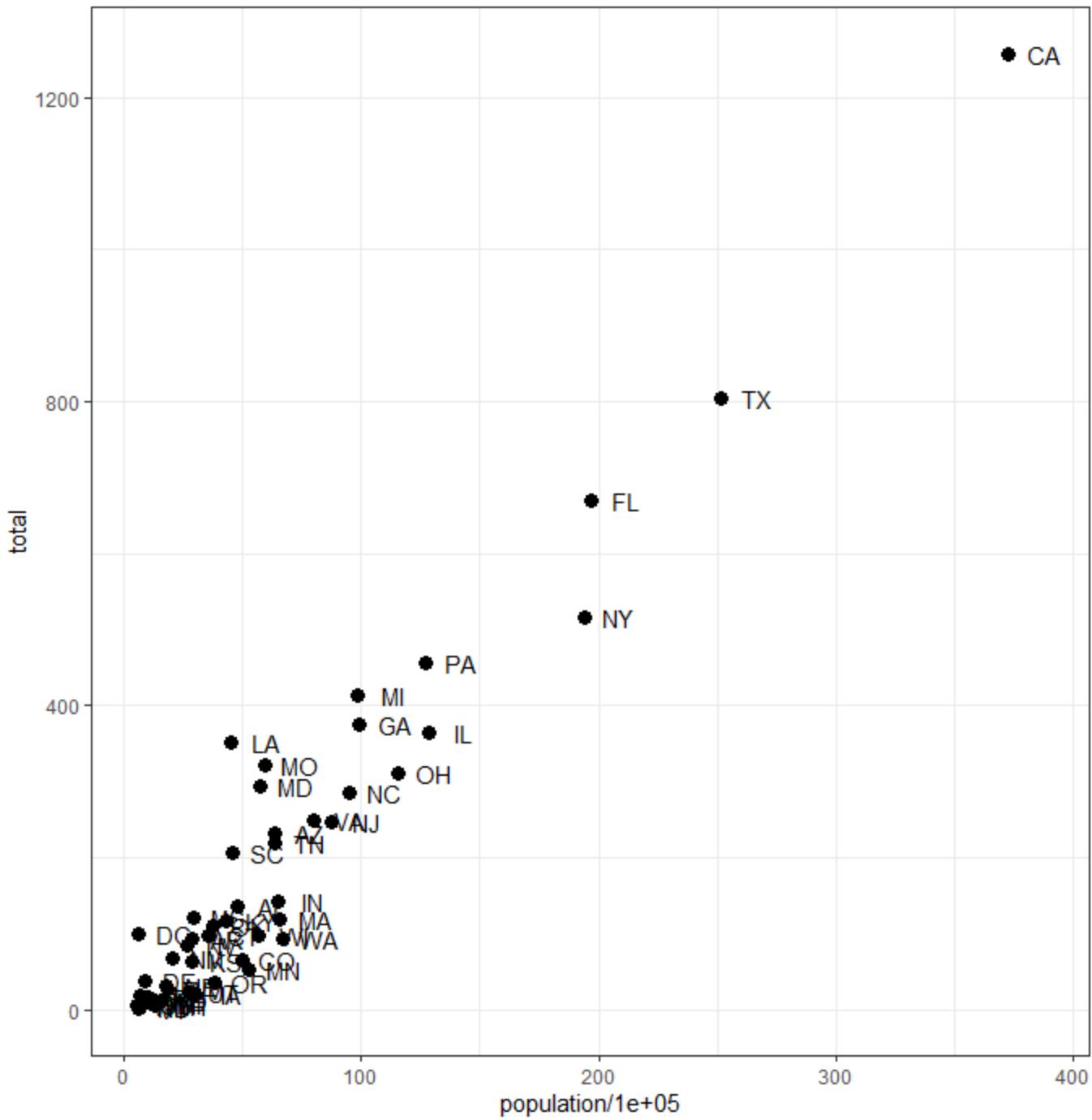
p +

```
geom_point(aes(x = population/100000, y = total), size = 3) +  
geom_label(aes(x = population/100000, y = total, label = abb))
```



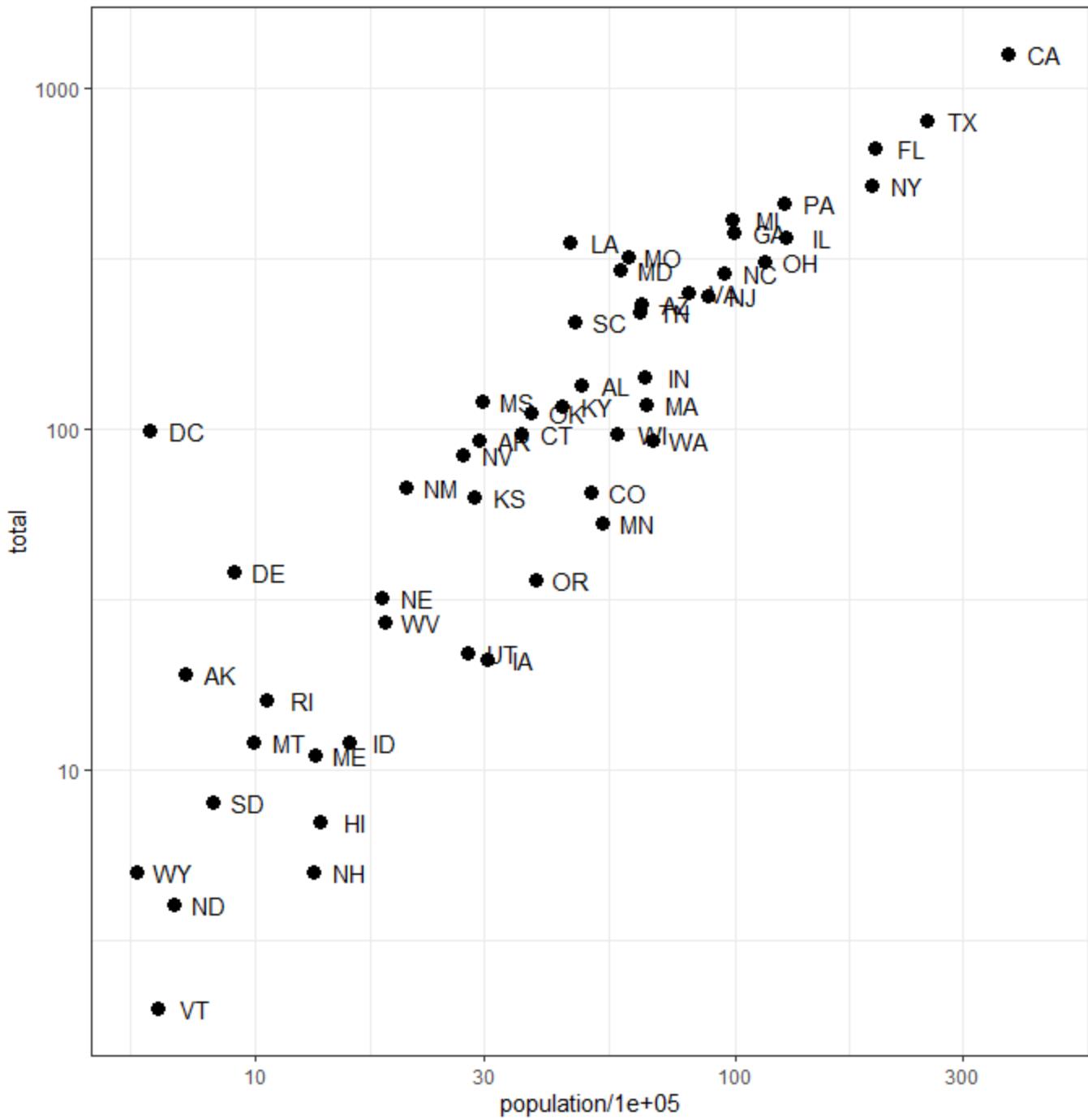
p +

```
geom_point(aes(x = population/100000, y = total), size = 3) +  
geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 15)
```



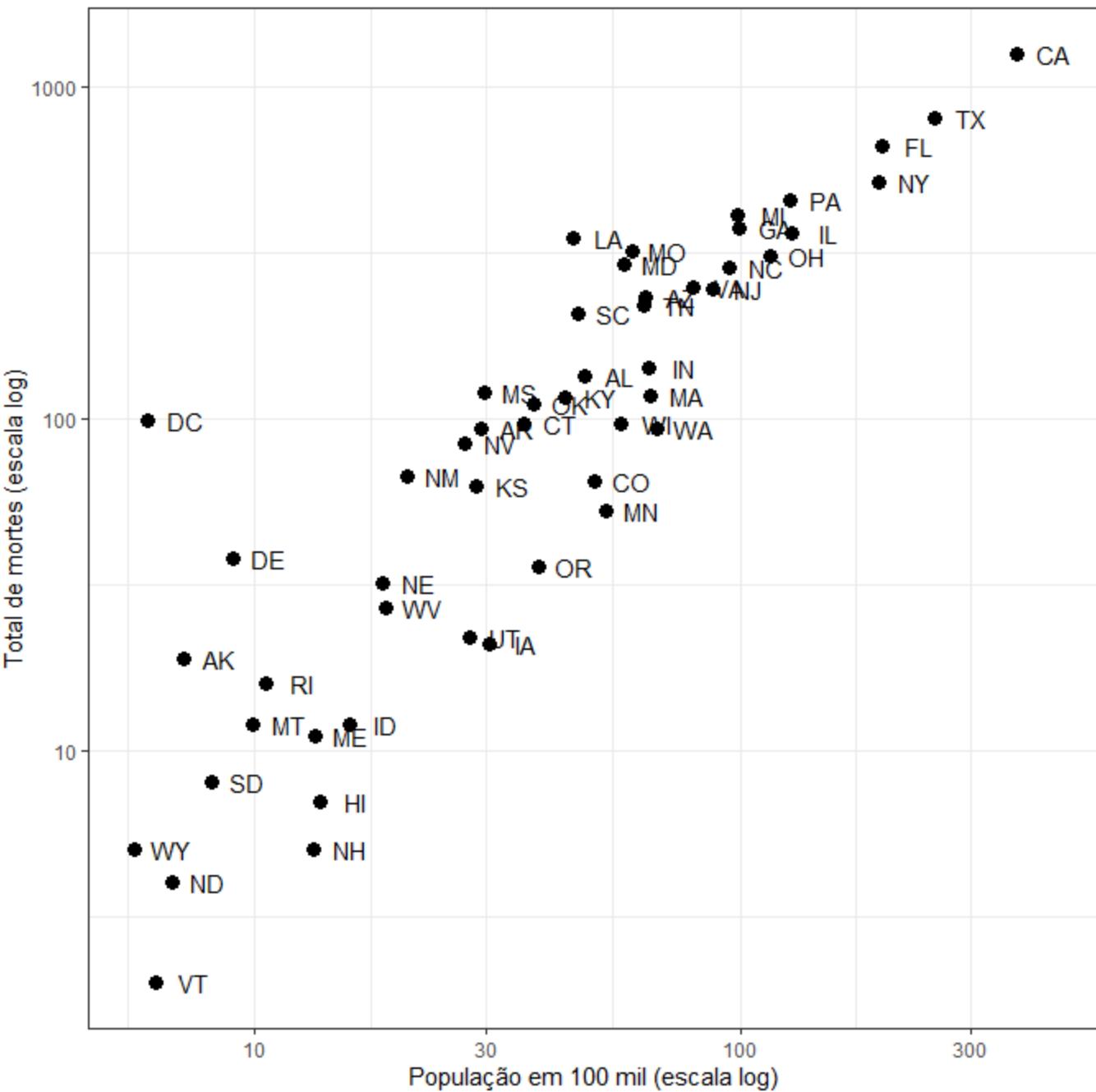
p +

```
geom_point(aes(x = population/100000, y = total), size = 3) +  
geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +  
scale_x_log10() +  
scale_y_log10()
```



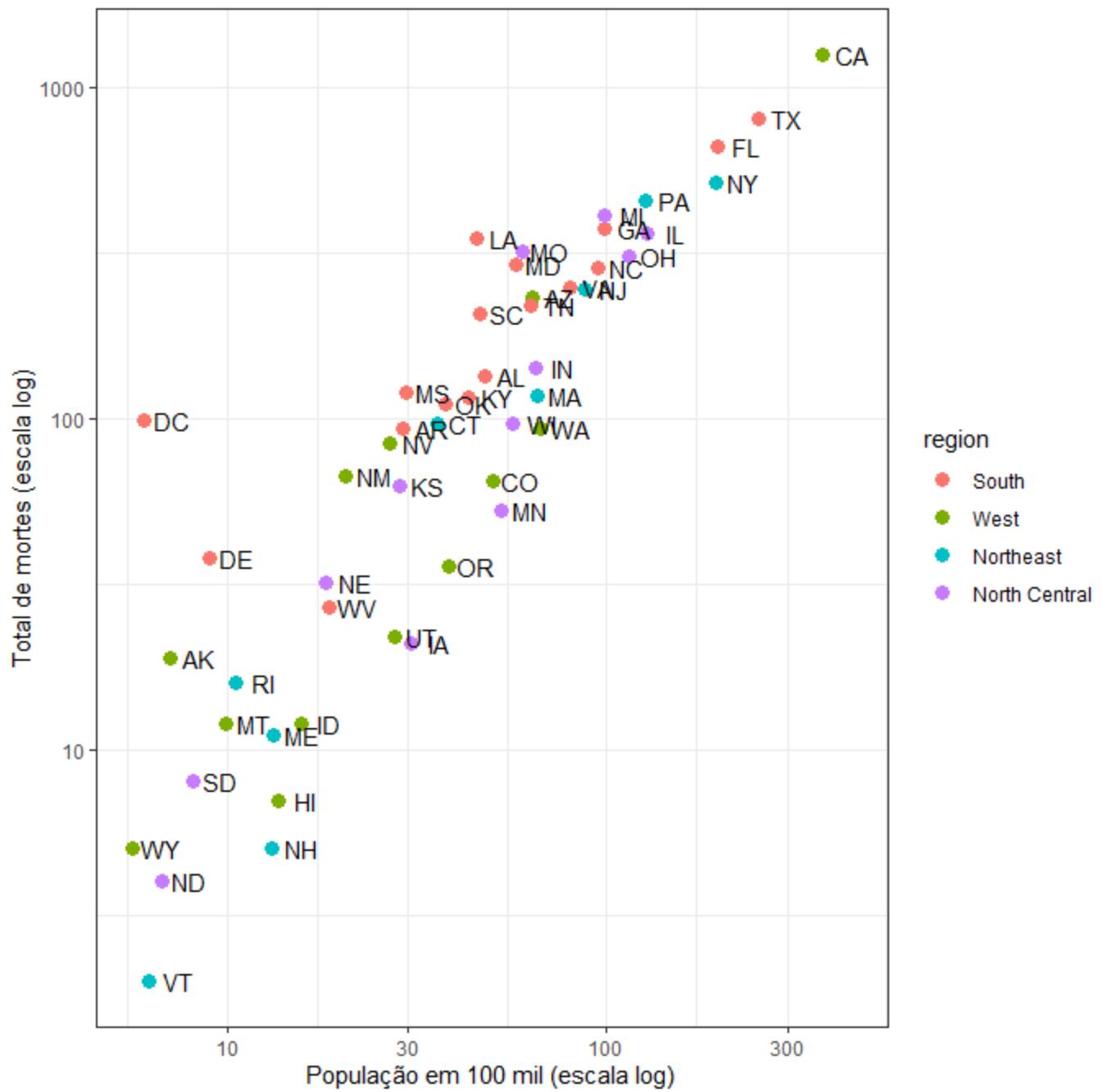
```
p +  
  geom_point(aes(x = population/100000, y = total), size = 3) +  
  geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +  
  scale_x_log10() +  
  scale_y_log10() +  
  xlab("População em 100 mil (escala log)") +  
  ylab("Total de mortes (escala log)") +  
  ggtitle("Mortes por arma de fogo nos EUA em 2010")
```

Mortes por arma de fogo nos EUA em 2010



```
p +  
  geom_point(aes(x = population/100000, y = total, color = region), size = 3) +  
  geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +  
  scale_x_log10() +  
  scale_y_log10() +  
  xlab("População em 100 mil (escala log)") +  
  ylab("Total de mortes (escala log)") +  
  ggtitle("Mortes por arma de fogo nos EUA em 2010")
```

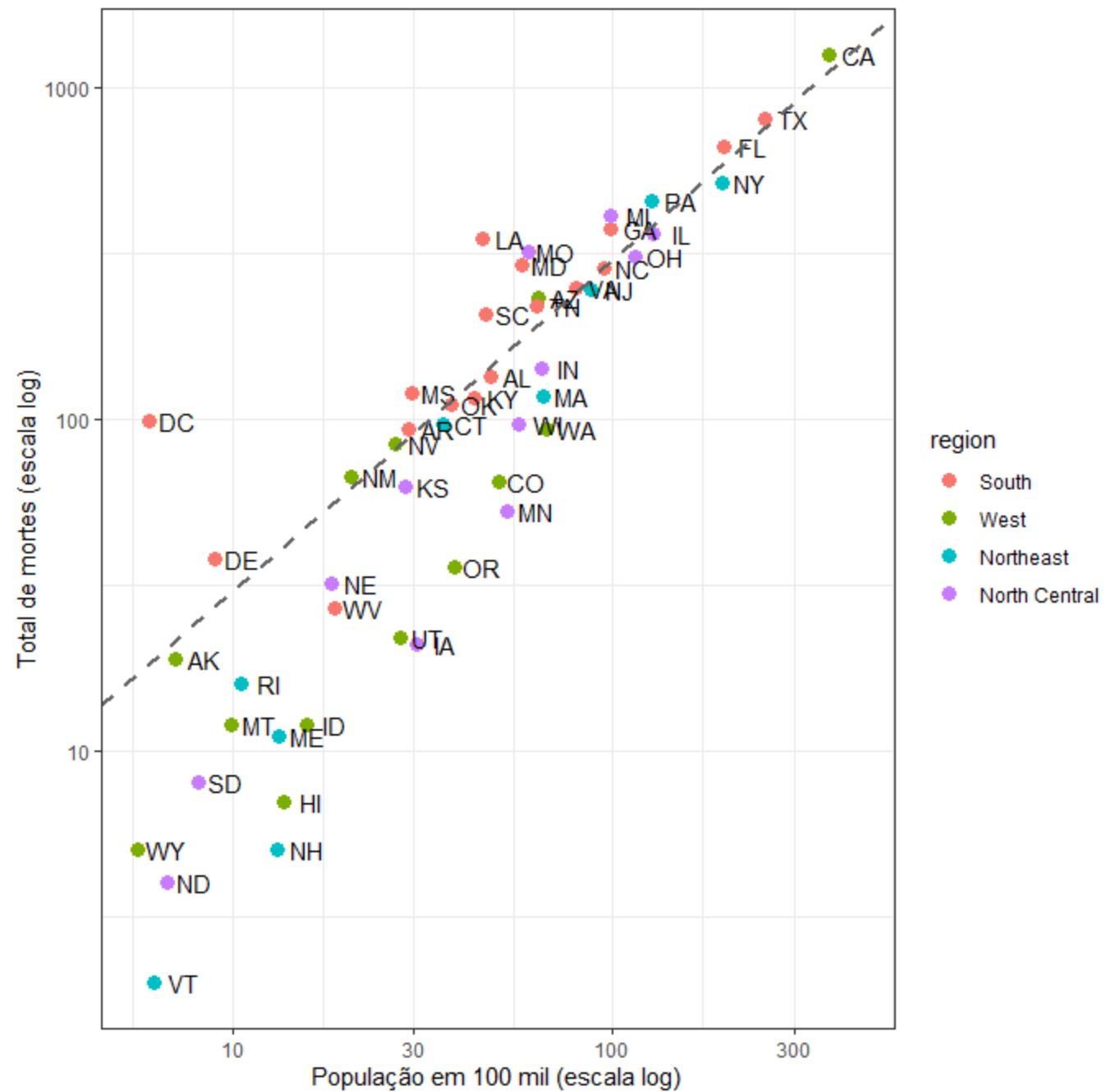
Mortes por arma de fogo nos EUA em 2010



```
tx <- sum(data$total) / sum(data$population) * 100000

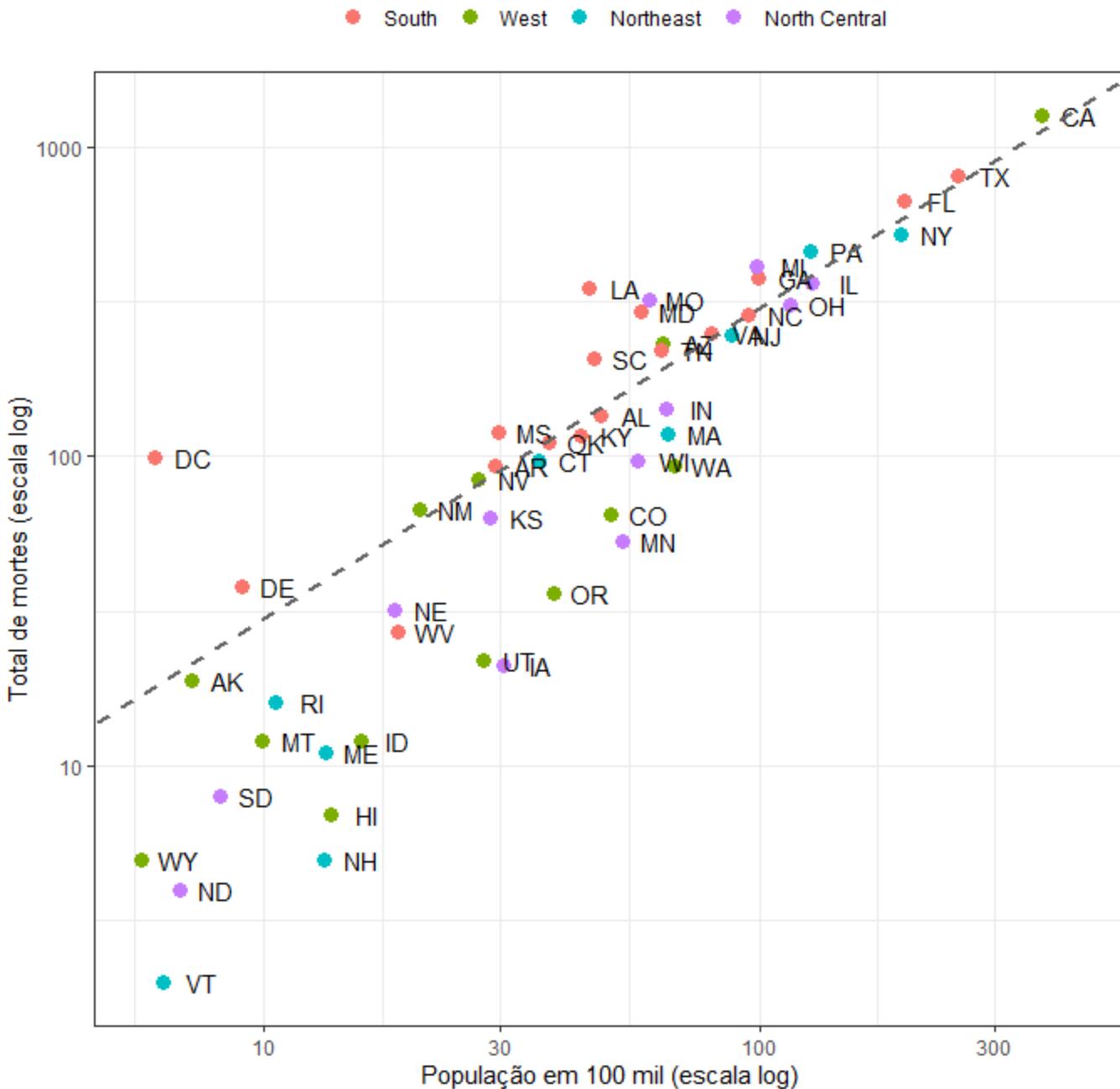
p +
  geom_point(aes(x = population/100000, y = total, color = region), size = 3) +
  geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +
  geom_abline(intercept = log10(tx), slope = 1,
    linetype = 'dashed', linewidth = 1, color = 'gray40') +
  scale_x_log10() +
  scale_y_log10() +
  xlab("População em 100 mil (escala log)") +
  ylab("Total de mortes (escala log)") +
  ggtitle("Mortes por arma de fogo nos EUA em 2010")
```

Mortes por arma de fogo nos EUA em 2010



```
p +  
  geom_point(aes(x = population/100000, y = total, color = region), size = 3) +  
  geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +  
  geom_abline(intercept = log10(tx), slope = 1,  
              linetype = 'dashed', linewidth = 1, color = 'gray40') +  
  scale_x_log10() +  
  scale_y_log10() +  
  xlab("População em 100 mil (escala log)") +  
  ylab("Total de mortes (escala log)") +  
  ggtitle("Mortes por arma de fogo nos EUA em 2010") +  
  theme(legend.title = element_blank(), legend.position = 'top')
```

Mortes por arma de fogo nos EUA em 2010



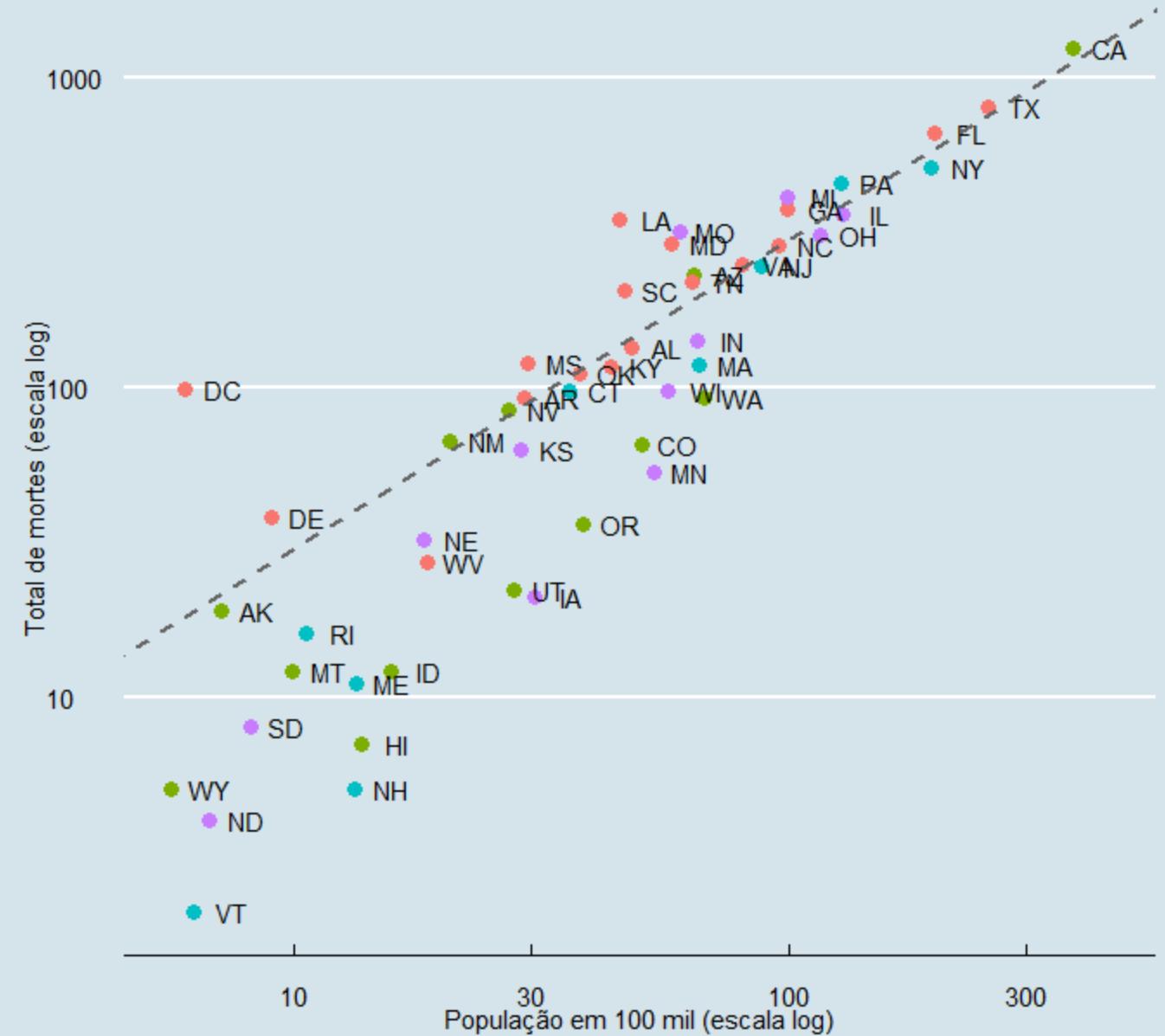
```
install.packages("ggthemes")
```

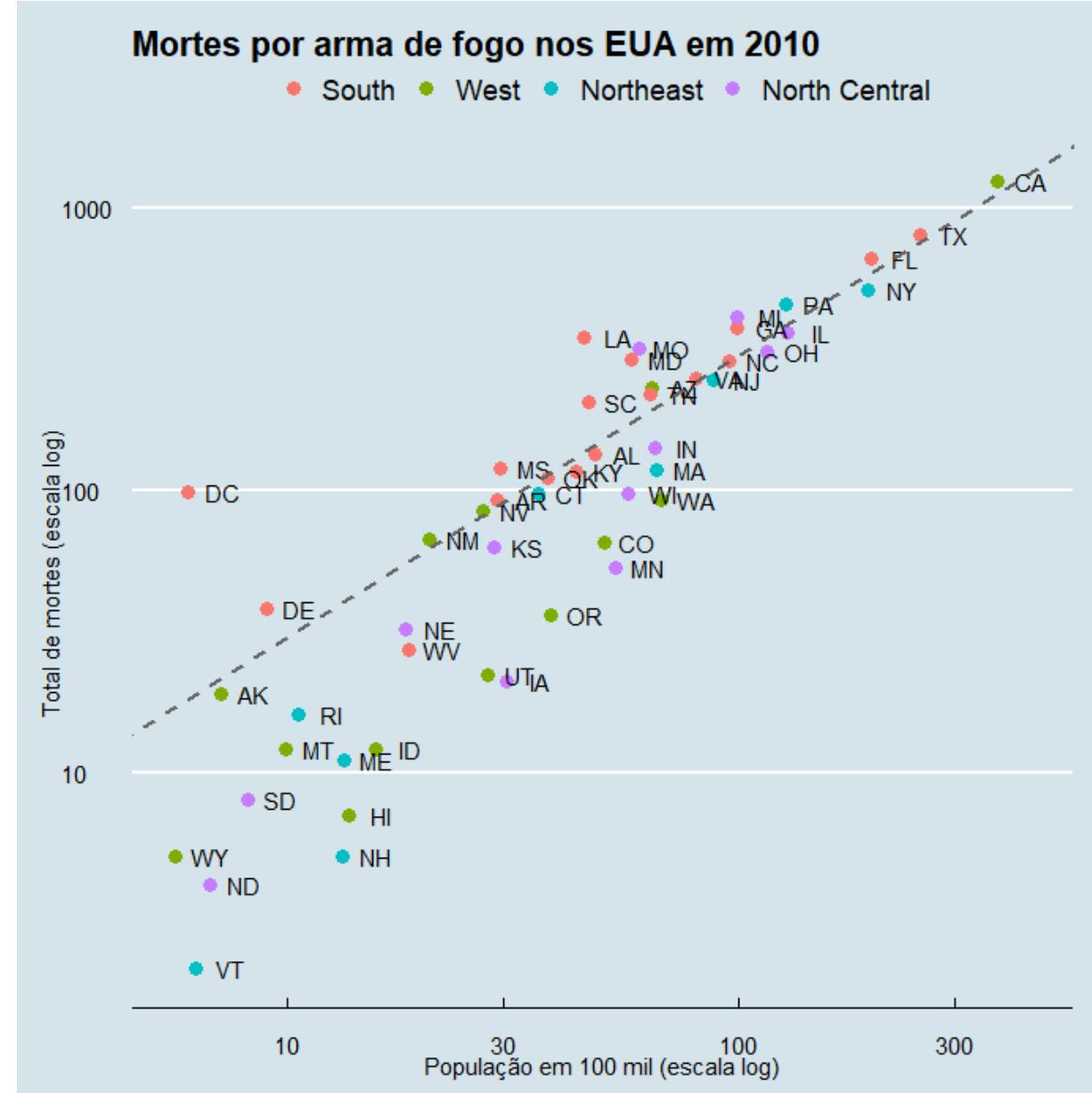
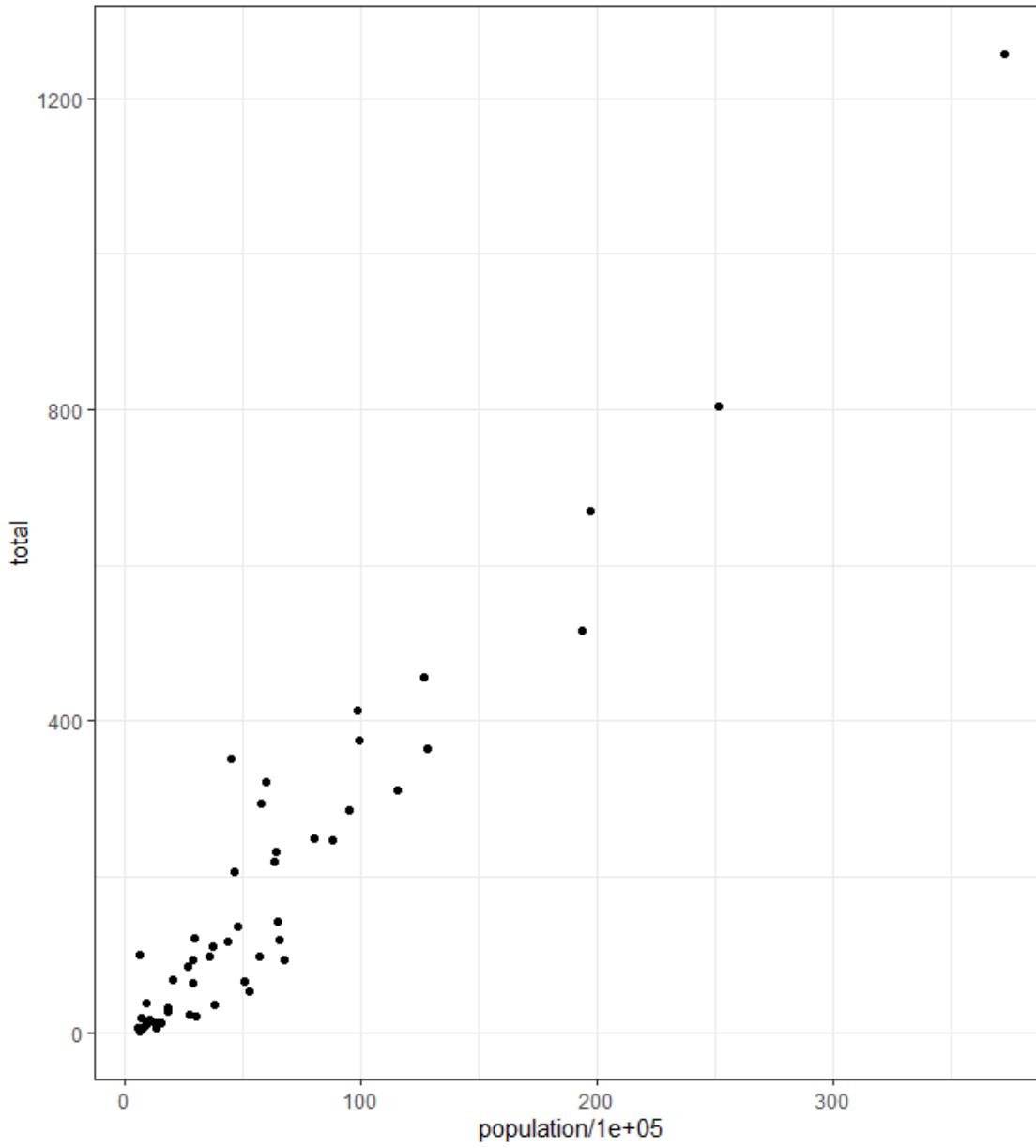
```
library(ggthemes)
```

```
p +  
  geom_point(aes(x = population/100000, y = total, color = region), size = 3) +  
  geom_text(aes(x = population/100000, y = total, label = abb), nudge_x = 0.075) +  
  geom_abline(intercept = log10(tx), slope = 1,  
              linetype = 'dashed', linewidth = 1, color = 'gray40') +  
  scale_x_log10() +  
  scale_y_log10() +  
  xlab("População em 100 mil (escala log)") +  
  ylab("Total de mortes (escala log)") +  
  ggtitle("Mortes por arma de fogo nos EUA em 2010") +  
  theme_economist() +  
  theme(legend.title = element_blank(), legend.position = 'top')
```

Mortes por arma de fogo nos EUA em 2010

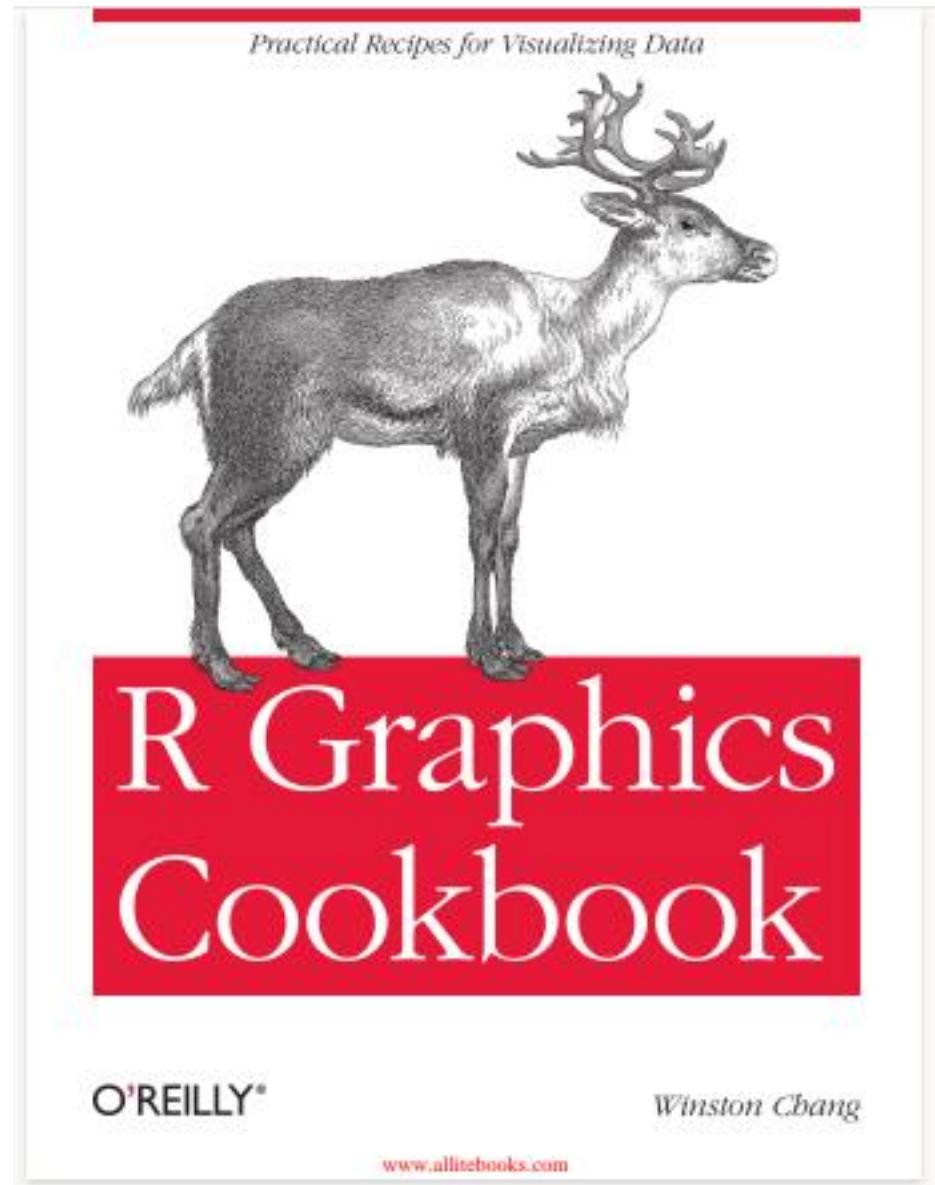
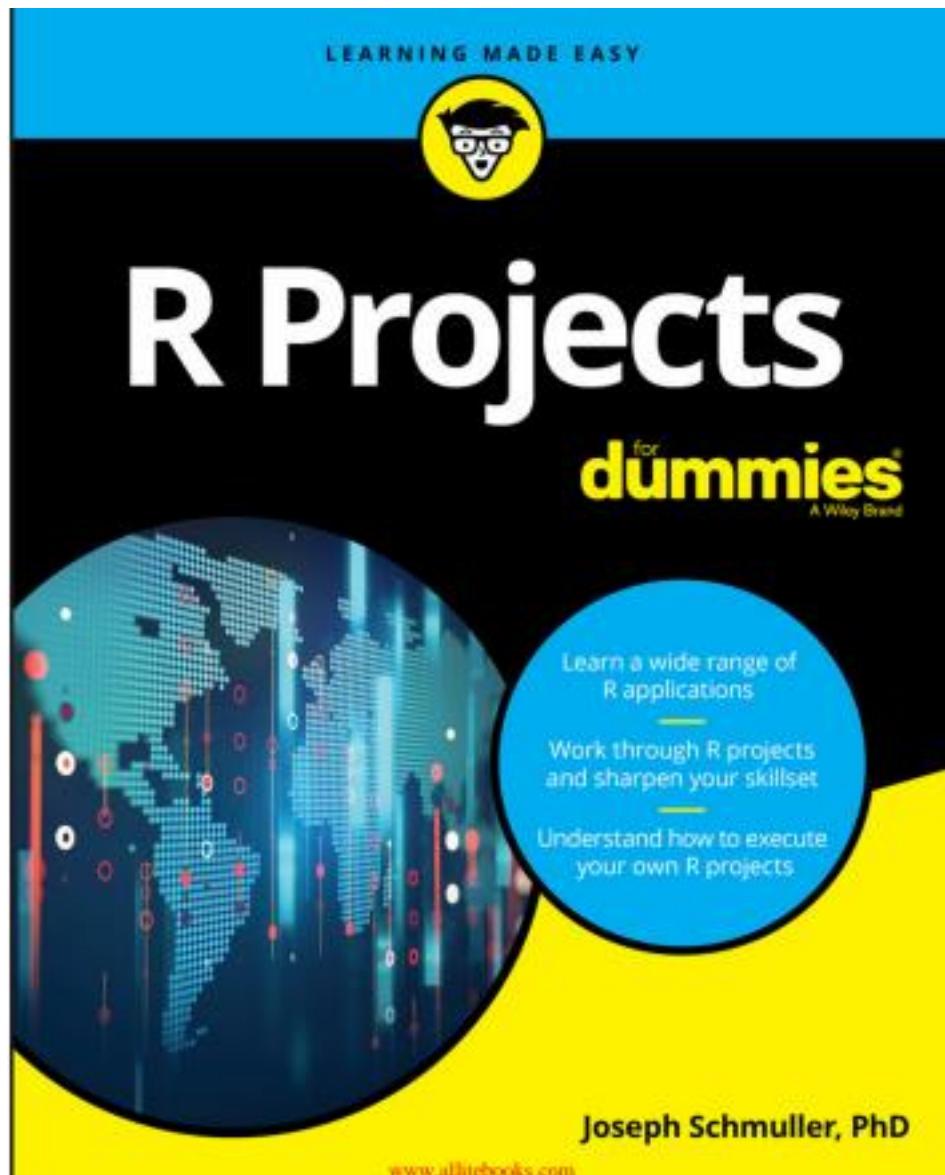
- South
- West
- Northeast
- North Central

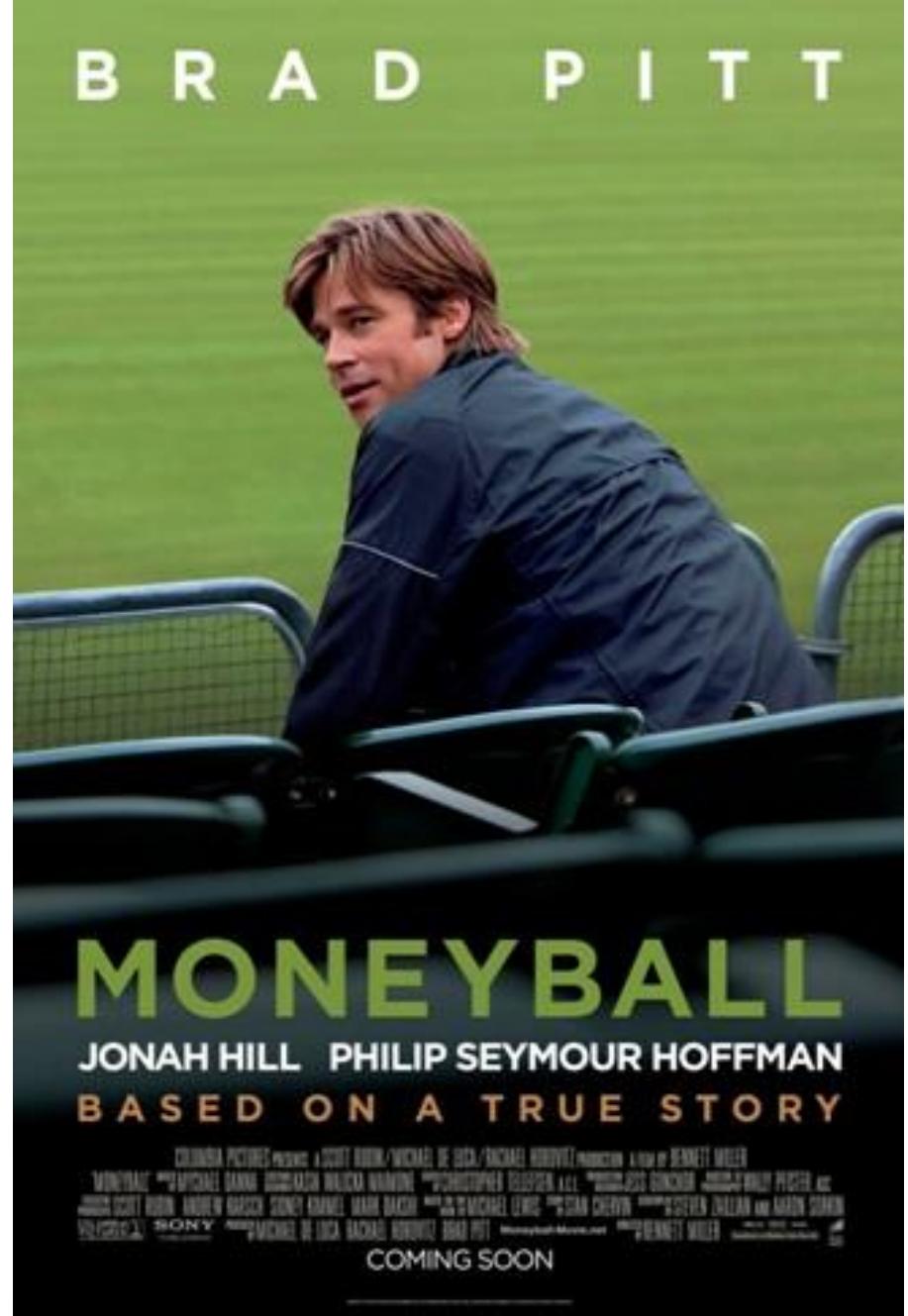
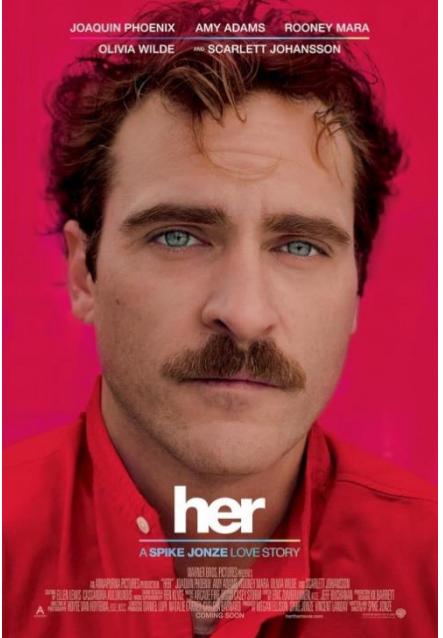




```
ggsave('meu.primeiro.plot.png')
```

**Outras fontes de
inspiração...**





OBRIGADO!

Nelson Quesado
Caio Gustavo

Introdução à Ciência de Dados com R

