

INSTITUTE  
OF LAW,  
POLITICS AND  
DEVELOPMENT



Sant'Anna  
School of Advanced Studies – Pisa

## Intro to R

Visualization and Programming in R

Prof. Kevin Koehler

[kevin.koehler@santannapisa.it](mailto:kevin.koehler@santannapisa.it)

# Schedule

We are missing two sessions. Do you prefer adding them to the back or do you want to schedule extra classes in the next weeks?



**Sant'Anna**

School of Advanced Studies – Pisa

# Today's class

## 1. Data classes and structures



Sant'Anna

School of Advanced Studies – Pisa

# Today's class

1. Data classes and structures
2. The grammar of graphics in `ggplot2`



Sant'Anna

School of Advanced Studies – Pisa

# Today's class

1. Data classes and structures
2. The grammar of graphics in `ggplot2`
3. Writing user-defined functions and automating repetitive tasks with `for` loops



Sant'Anna

School of Advanced Studies – Pisa

# Today's class

1. Data classes and structures
2. The grammar of graphics in `ggplot2`
3. Writing user-defined functions and automating repetitive tasks with `for` loops
4. Exercises



Sant'Anna

School of Advanced Studies – Pisa

# Why visualization?



Sant'Anna

School of Advanced Studies – Pisa

# Data visualization

1. Visualization allows you to better understand data and to efficiently communicate results



Sant'Anna

School of Advanced Studies – Pisa

# Data visualization

1. Visualization allows you to better understand data and to efficiently communicate results
2. Many important journals in Political Science require plots instead of tables for publication



Sant'Anna

School of Advanced Studies – Pisa

# Data visualization

1. Visualization allows you to better understand data and to efficiently communicate results
2. Many important journals in Political Science require plots instead of tables for publication
3. Visual representations are easier to read than full tables



Sant'Anna

School of Advanced Studies – Pisa

# The grammar of graphics



Sant'Anna

School of Advanced Studies – Pisa

# ggplot2



The `ggplot2` package is part of the `tidyverse`. It provides a specific system of constructing plots.



Sant'Anna

School of Advanced Studies – Pisa

First, some theory



Sant'Anna

School of Advanced Studies – Pisa

# Data classes and structures

There are three (four?) data classes in R

1. Numeric
2. Character
3. Factor
4. (Logical)

You can convert objects between classes with `as.numeric()`,  
`as.character()`, and `as.factor()`.

# Why factors are weird

factors have levels and labels. Internally, R represents factors as ordered sequences. This can lead to odd behavior such as:

```
factor <- as.factor(c(1,2,4,2,5,99))  
factor
```

```
[1] 1 2 4 2 5 99
```

```
Levels: 1 2 4 5 99
```

```
as.numeric(factor)
```

```
[1] 1 2 3 2 4 5
```

What R does here is to order the elements of the factor (in this case by size since the elements are numerical—if we had entered characters they would be ordered alphabetically). It then returns the position in this order as the numerical value.



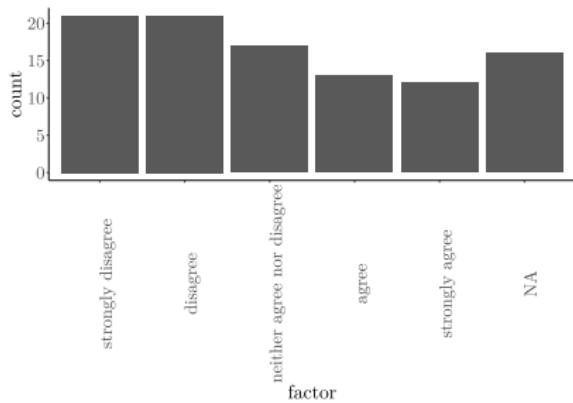
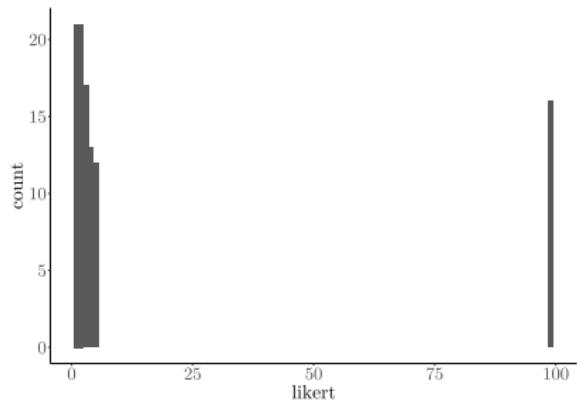
## Why factors are still useful

Despite this odd behavior, factors can be useful for plotting.

Say, we have a Likert scale response such as in the Tunisia survey.  
Plotting this variable as a vector allows us to control the labels:

```
likert <- sample(c(1:5, 99), size = 100, replace = TRUE)
factor <- factor(likert,
                  levels=c(1,2,3,4,5,99),
                  labels=c("strongly disagree",
                          "disagree",
                          "neither agree nor disagree",
                          "agree",
                          "strongly agree",
                          "NA"))
example <- data.frame(likert=likert,
                      factor=factor)
```

Plot both likert and factor to see the difference



This behavior is similar to what you might know from Stata  
(where you can assign value labels)

# Data structures

1. **Vectors** (numerical or character, only one class):  
`id <- c(1,2,3,4,5), gender <- c("m","f","m","f","m")`

# Data structures

1. **Vectors** (numerical or character, only one class): `id <- c(1,2,3,4,5), gender <- c("m","f","m","f","m")`
2. **Matrices** (two dimensional, only one class): `matrix1 <- matrix(c(1,2,3,4,5,23,34,19,46,38), ncol = 2)`



# Data structures

1. **Vectors** (numerical or character, only one class): `id <- c(1,2,3,4,5), gender <- c("m","f","m","f","m")`
2. **Matrices** (two dimensional, only one class): `matrix1 <- matrix(c(1,2,3,4,5,23,34,19,46,38), ncol = 2)`
3. **Lists** (collection of elements of different classes): `list1 <- list(id, gender, matrix1)`



# Data structures

1. **Vectors** (numerical or character, only one class):  
`id <- c(1,2,3,4,5), gender <- c("m","f","m","f","m")`
2. **Matrices** (two dimensional, only one class):  
`matrix1 <- matrix(c(1,2,3,4,5,23,34,19,46,38), ncol = 2)`
3. **Lists** (collection of elements of different classes):  
`list1 <- list(id, gender, matrix1)`
4. **Data frames** (two dimensional structure, different classes):  
`df1 <- data.frame(id = c(1,2,3,4,5), gender = c("m","f","m","f","m"))`



# Data structures

R saves the output of many more complex functions as lists. The following code creates random x and y variables and regresses y on x:

```
y <- runif(100, min = 0, max = 1)
x <- runif(100, min = 10, max = 20)
model1 <- lm(y ~ x)
names(model1)
```

```
[1] "coefficients"   "residuals"      "effects"        "rank"
[5] "fitted.values"  "assign"         "qr"             "df.residuals"
[9] "xlevels"        "call"          "terms"          "model"
```

# Data structures

```
summary(model1)
```

Call:

```
lm(formula = y ~ x)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.58155	-0.19851	0.06131	0.23177	0.47878

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )							
(Intercept)	0.738717	0.150656	4.903	3.74e-06 ***							
x	-0.013389	0.009752	-1.373	0.173							
---											
Signif. codes:	0	'***'	0.001	'**'	0.01	'*'	0.05	'.'	0.1	' '	1

Residual standard error: 0.2809 on 98 degrees of freedom

Multiple R-squared: 0.01887, Adjusted R-squared: 0.008857

F-statistic: 1

# Visualization with ggplot



Sant'Anna

School of Advanced Studies – Pisa

# ggplot2



In ggplot2, plots consist of three basic components:

1. data
2. aesthetics
3. coordinate system

You can combine as many different aesthetic layers as you like.



Sant'Anna

School of Advanced Studies – Pisa

## ggplot2

```
plot <- ggplot(data=[name of data],  
                aes(x=x,    → x-coordinates  
                     y=y)) + → y-coordinates (if applicable)  
                geom_[name of geom]()
```

Diagram illustrating the components of a ggplot2 command:

- Name of the plot (points to "ggplot")
- Declare ggplot (points to "ggplot")
- Name of the data set (points to "data")
- Name of the aesthetic (points to "aes")





## Plotting a single variable



Sant'Anna

School of Advanced Studies – Pisa

# Plotting a single variable

- We frequently need to plot single variables to visually inspect their distribution



# Plotting a single variable

- ▶ We frequently need to plot single variables to visually inspect their distribution
- ▶ There are various options for doing this:



# Plotting a single variable

- ▶ We frequently need to plot single variables to visually inspect their distribution
- ▶ There are various options for doing this:
  - ▶ Histograms (for interval-scaled variables)



Sant'Anna

School of Advanced Studies – Pisa

# Plotting a single variable

- ▶ We frequently need to plot single variables to visually inspect their distribution
- ▶ There are various options for doing this:
  - ▶ Histograms (for interval-scaled variables)
  - ▶ Density plots (for interval-scaled variables)



# Plotting a single variable

- ▶ We frequently need to plot single variables to visually inspect their distribution
- ▶ There are various options for doing this:
  - ▶ Histograms (for interval-scaled variables)
  - ▶ Density plots (for interval-scaled variables)
  - ▶ Box plots (for interval-scaled variables)



# Plotting a single variable

- ▶ We frequently need to plot single variables to visually inspect their distribution
- ▶ There are various options for doing this:
  - ▶ Histograms (for interval-scaled variables)
  - ▶ Density plots (for interval-scaled variables)
  - ▶ Box plots (for interval-scaled variables)
  - ▶ Bar charts (for categorical variables)



Let's plot the distribution of respondents' age  
in the Tunisia survey

# Histograms

- ▶ A **histogram** is a type of chart representing the distribution of a single numerical variable.
- ▶ In histograms, the x-axis is divided into **bins** which represent ranges of the variable.
- ▶ The height of the bars represent **how many observations** fall within a bin.

In `ggplot2`, histogram layers can be produced with `geom_histogram()`.

**Write code for a histogram of age.**



# Write code for a histogram of age

```
plot <- ggplot(data=[name of data],  
                aes(x=x, → x-coordinates  
                     y=y)) + → y-coordinates (if applicable)  
                geom_[name of geom]()
```

Name of the plot

Declare ggplot

Name of the data set

Name of the aesthetic

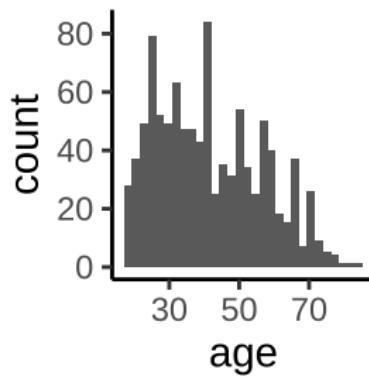


Sant'Anna

School of Advanced Studies – Pisa

# Histograms

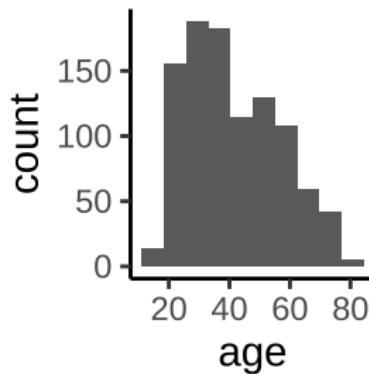
```
plot <- ggplot(data=tun22,  
                 aes(x=age)) +  
  geom_histogram() +  
  theme_classic()  
  
plot
```



# Histograms

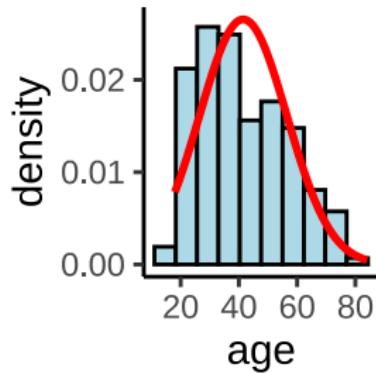
You can set the number of bins with the `bin` option (the default is 30):

```
plot1 <- ggplot(data=tun22,  
                  aes(x=age)) +  
  geom_histogram(bins = 10) +  
  theme_classic()  
  
plot1
```



# Histograms

We frequently want to check whether a variable aligns with the normal distribution. We can graphically check this by overlaying a normal distribution with the same mean and standard deviation.



# Histograms

Here is the code:

```
plot <- ggplot(data=tun22,  
                 aes(x=age)) +  
  geom_histogram(aes(y=..density..),  
                 bins = 10,  
                 fill="lightblue",  
                 color="black") +  
  stat_function(fun = dnorm,  
                args = list(mean = mean(tun22$age, na.rm=T),  
                            sd = sd(tun22$age, na.rm=T)),  
                color = "red", size = 1) +  
  theme_classic()  
plot
```

We now plot **density**, not counts

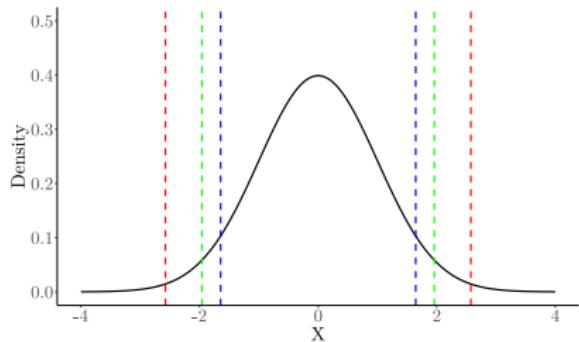


## Quick detour: What is density?

- ▶ Density refers to the **probability density function (PDF)**.
- ▶ It shows how likely a variable is to take a specific value.
- ▶ The total area under the density curve (or histogram) equals **1**.



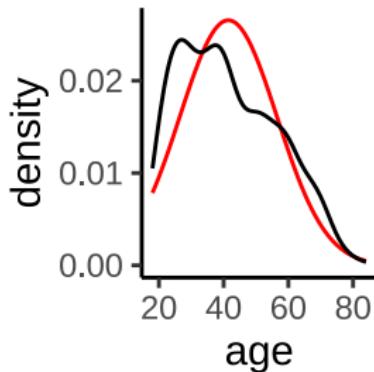
# The normal distribution



- ▶ 90% of observations fall into the area between the blue lines
- ▶ 95% of observations fall into the area between the green lines
- ▶ 99% of observations fall into the area between the red lines

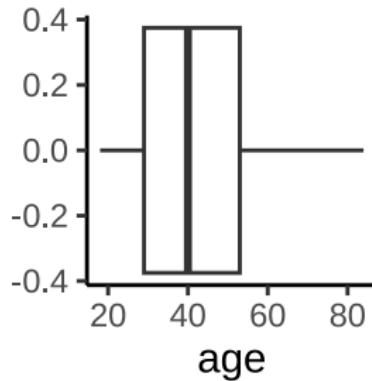
# Density plots

```
plot2 <- ggplot(data=tun22,
                  aes(x=age)) +
  stat_function(fun = dnorm,
                args = list(mean = mean(tun22$age, na.rm=T),
                            sd = sd(tun22$age, na.rm=T)),
                color = "red") +
  geom_density() +
  theme_classic()
plot2
```



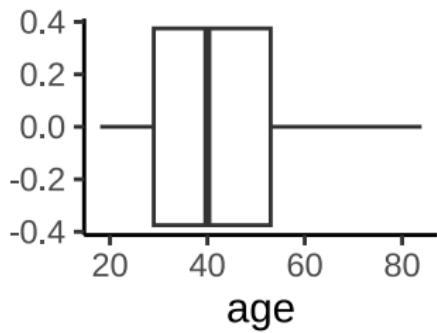
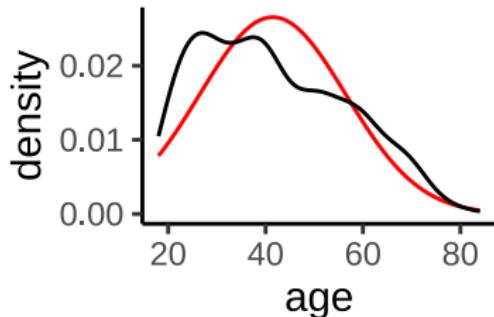
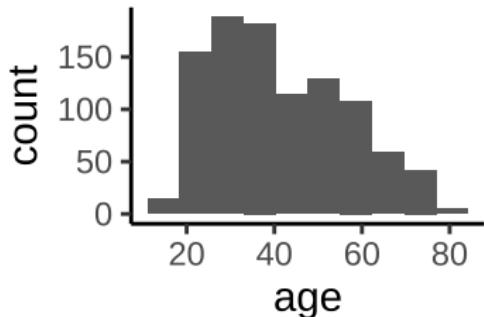
# Boxplots

```
plot3 <- ggplot(data=tun22,  
                  aes(x=age)) +  
  geom_boxplot() +  
  theme_classic()  
plot3
```



# Summary

## Respondent Age



# ggplot2

We can combine `ggplot` with `tidyverse` data manipulation:

```
tun22 %>%
  group_by(region) %>%
  summarize(n = n()) %>%
  ggplot() +
  geom_bar(
    aes(x = region, y = n),
    stat = "identity"
  ) +
  labs(x = "",
       y = "") +
  theme_classic() +
  theme(axis.text.x = element_text(angle = 90))
```



# ggplot2

We can combine `ggplot` with `tidyverse` data manipulation:

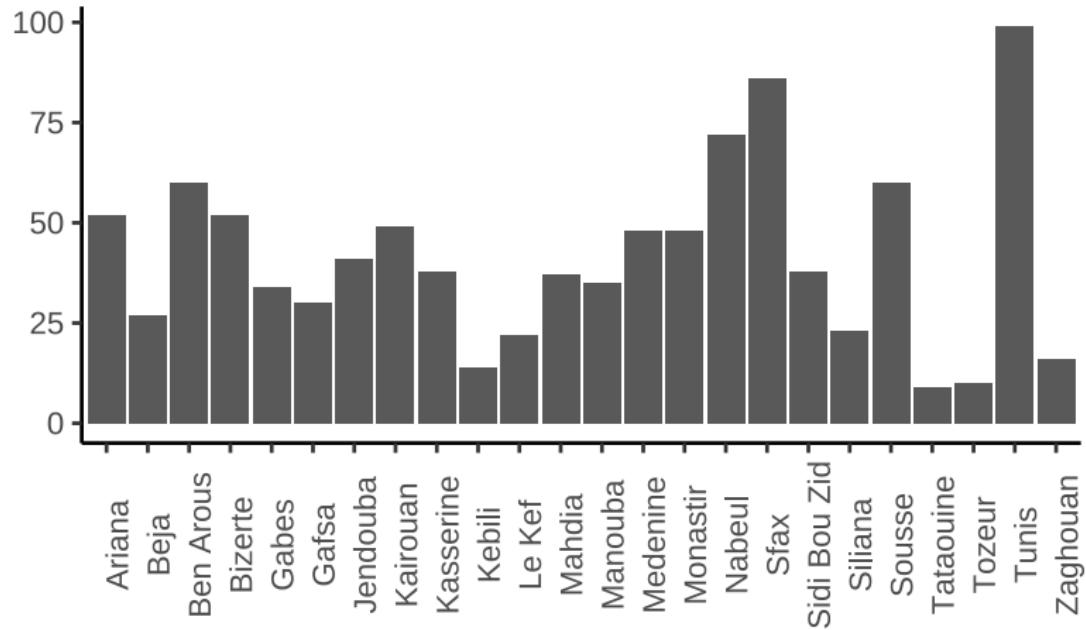
```
tun22 %>%
  group_by(region) %>%
  summarize(n = n()) %>%
  ggplot() +
  geom_bar(
    aes(x = region, y = n),
    stat = "identity"
  ) +
  labs(x="",
       y="",
       title="Number of respondents by region") +
  theme_classic() +
  theme(axis.text.x = element_text(angle = 90))
```

What will this graph look like?



# Bar charts

## Number of respondents by region



# ggplot2

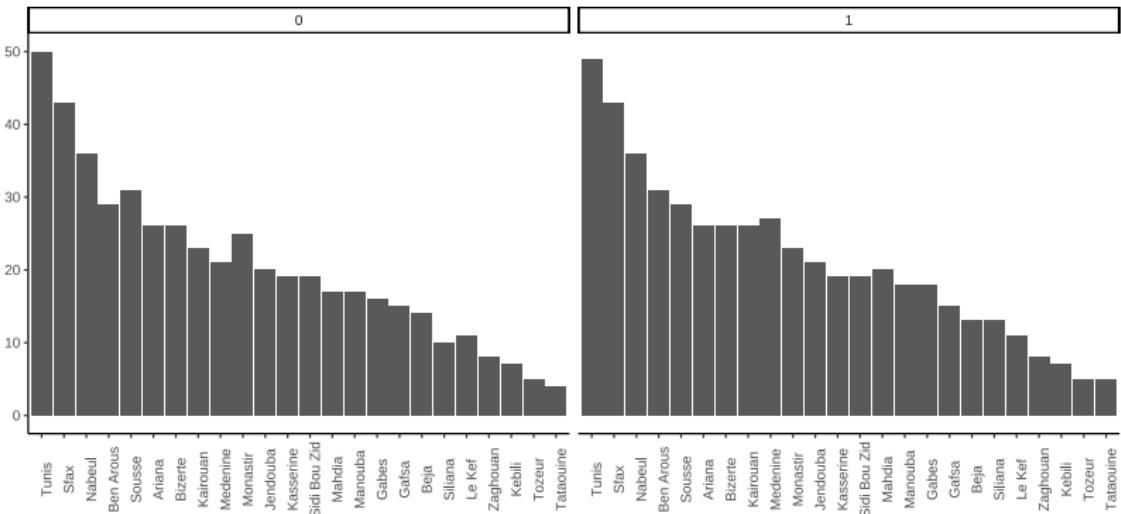
We can also produce different plots by group using the `facet()` layer

```
tun22 %>%
  group_by(region,female) %>%
  summarize(n = n()) %>%
  ggplot() +
  geom_bar(
    aes(x = reorder(region,-n), y = n), stat = "identity") +
  labs(x="",
       y="",
       title="Number of male and female respondents by region") +
  theme_classic() +
  theme(axis.text.x = element_text(angle = 90)) +
  facet_wrap(~female)
```



# ggplot2

Number of male and female respondents by region



## Plotting two variables



Sant'Anna

School of Advanced Studies – Pisa

# Scatterplots

- ▶ **Scatterplots** are plots displaying the relationship between two numerical variables.
- ▶ The variables are plotted on the x and y-axis, each dot represents one observation (or case).



# V-Dem data

Download the V-Dem data:

```
# First, you need to have the devtools package installed
install.packages("devtools")
# now, install the vdemdata package directly from GitHub
devtools::install_github("vdeminstitute/vdemdata")
# Now load the package
library(vdemdata)

data <- vdem

data <- data %>%
  dplyr::select(v2x_polyarchy,
                e_gdppc, year,
                country_name,
                e_regionpol)
```

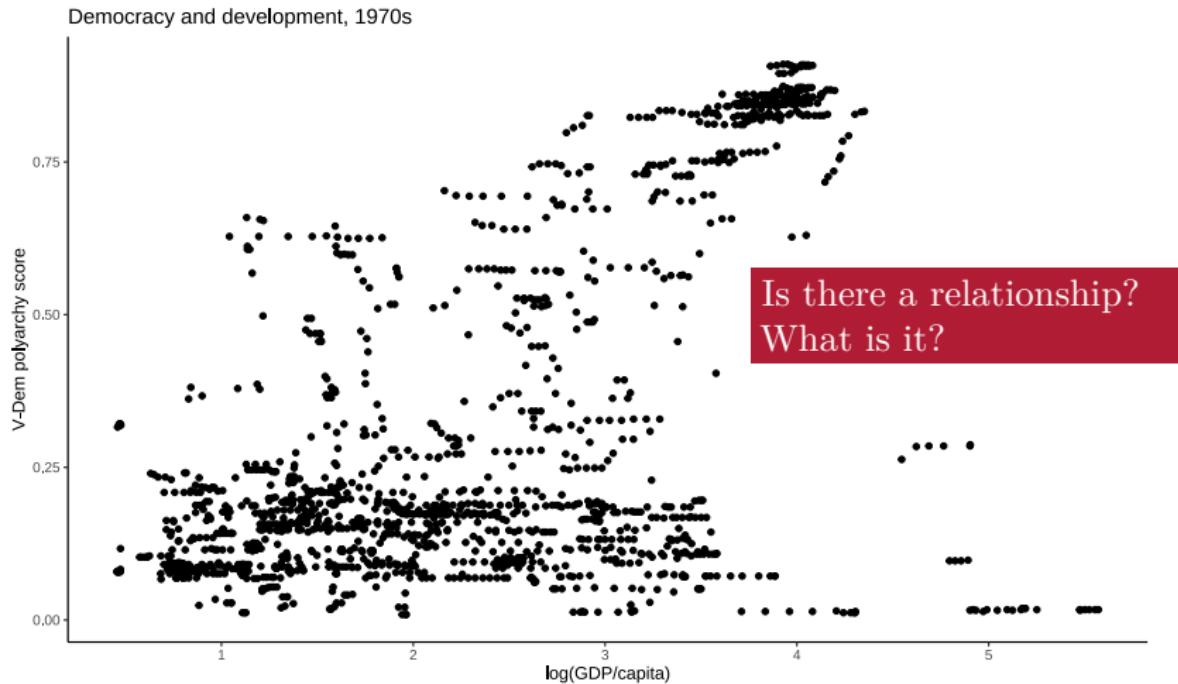


Let's create a scatterplot showing the relationship between log GDP/capita and the polyarchy score for the 1970s

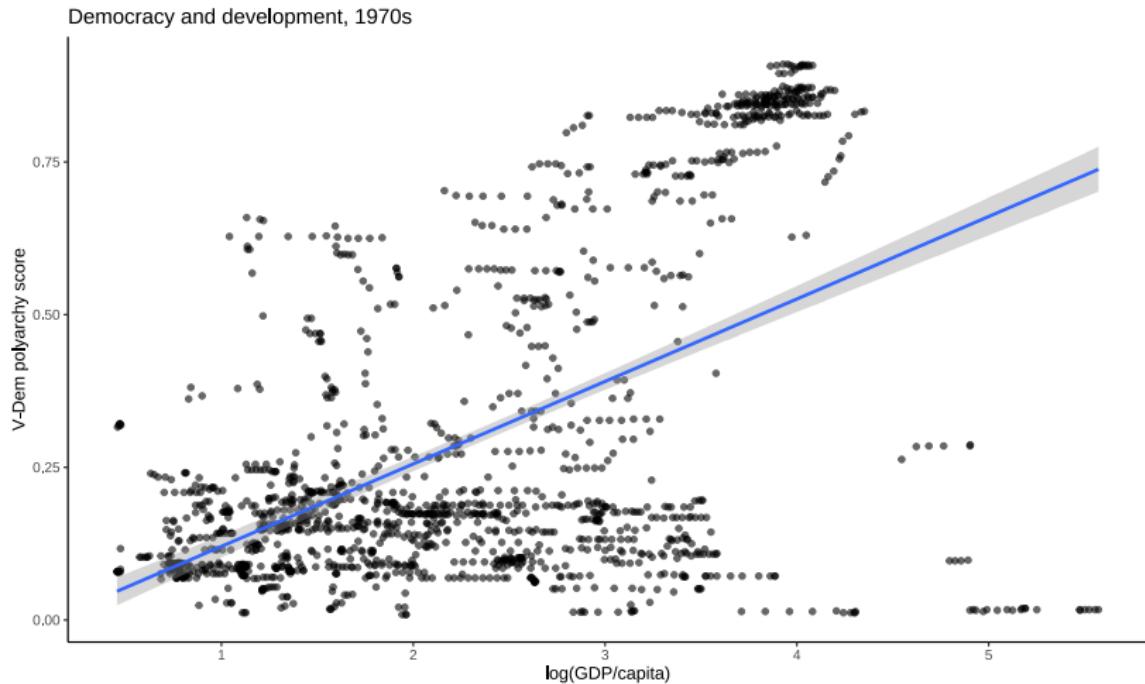
# Scatterplots



# Scatterplots



# Scatterplots

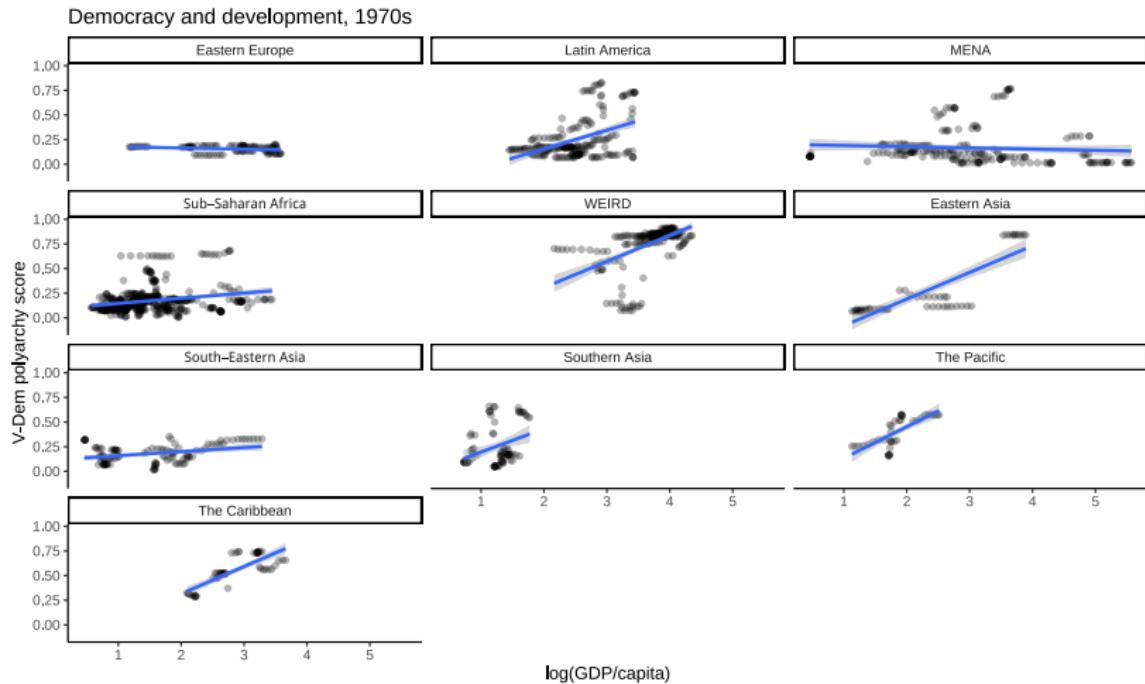


# Scatterplots

```
plot1 <- data %>%
  filter(year>1969 & year<1980) %>%
  rename(democracy=v2x_polyarchy) %>%
  mutate(gdp=log(e_gdppc)) %>%
  ggplot(aes(x=gdp,
             y=democracy)) +
  geom_point(alpha=0.6) +
  geom_smooth(method = "lm") +
  theme_classic() +
  labs(x="log(GDP/capita)",
       y="V-Dem polyarchy score",
       title = "Democracy and development, 1970s")
```



# By the way



# Combining graphs with the patchwork package

```
plot1 <- data %>%
  filter(year>1969 & year<1980) %>%
  rename(democracy=v2x_polyarchy) %>%
  mutate(gdp=log(e_gdppc)) %>%
  ggplot() +
  geom_point(aes(
    x=gdp,
    y=democracy
  ), alpha=0.6, size=1) +
  geom_smooth(aes(
    x=gdp,
    y=democracy),
    method = "lm") +
  theme_classic() +
  labs(x="log(GDP/capita)",
       y="V-Dem polyarchy score",
       title = "All countries")
```

Plot 1

```
plot2 <- data %>%
  mutate(region = factor(e_regionpol,
                        levels = 1:10,
                        labels = c("Eastern Europe",
                                  "Latin America",
                                  "MENA",
                                  "Sub-Saharan Africa",
                                  "WEIRD",
                                  "Eastern Asia",
                                  "South-Eastern Asia",
                                  "Southern Asia",
                                  "The Pacific",
                                  "The Caribbean")))) %>%
  filter(year > 1969 & year < 1980) %>%
  rename(democracy = v2x_polyarchy) %>%
  mutate(gdp = log(e_gdppc)) %>%
  ggplot(aes(x = gdp, y = democracy)) +
  geom_point(alpha=0.3, size=1, show.legend = F) +
  geom_smooth(method = "lm", show.legend = FALSE) +
  theme_classic() +
  labs(
    x = "log(GDP/capita)",
    y = "V-Dem polyarchy score",
    title = "By region"
  ) +
```

```
facet_wrap(~ region, ncol = 3)
```

Combine

```
library(patchwork)
all_plots <- plot1 + plot2
all_plots + plot_annotation(
  title = "Democracy and Development,
  1970s"
)
```



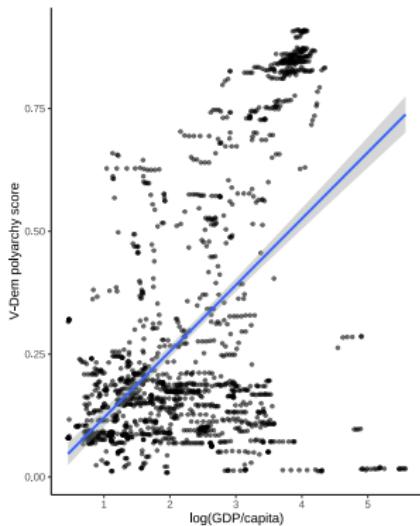
Sant'Anna

School of Advanced Studies – Pisa

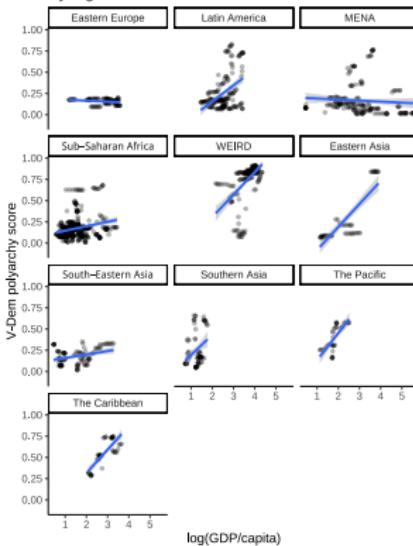
# Combining graphs with the patchwork package

Democracy and Development, 1970s

All countries



By region



Sant'Anna  
School of Advanced Studies – Pisa

# Basic programming



# Operators

## 1. Relational operators

1. < less than
2. > greater than
3. <= less than or equal to
4. >= greater than or equal to
5. == equal (identical) to
6. != not equal to

## 2. Logical operators

1. ! not ( $\neg$  in propositional logic)
2. & and ( $\wedge$  in propositional logic)
3. | or ( $\vee$  in propositional logic)

Operators return an object of the class “logical” which can have the values **TRUE** or **FALSE**.



# Operators

What will this evaluate to and why?

- ▶ `1 < 2 & "left"=="right"` ▶ ?
- ▶ `1 < 2 | "left"=="right"` ▶ ?
- ▶ `1 < 2 | !"left"=="right"` ▶ ?
- ▶ `!(1 < 2 | "left"=="right")` ▶ ?



# Operators

What will this evaluate to and why?

- |   |   |
|---|---|
| <ul style="list-style-type: none"><li>▶ <code>1 &lt; 2 &amp; "left"=="right"</code></li><li>▶ <code>1 &lt; 2   "left"=="right"</code></li><li>▶ <code>1 &lt; 2   !"left"=="right"</code></li><li>▶ <code>!(1 &lt; 2   "left"=="right")</code></li></ul> | <ul style="list-style-type: none"><li>▶ FALSE</li><li>▶ TRUE</li><li>▶ TRUE</li><li>▶ FALSE</li></ul> |
|---|---|



# Conditional statements

All programming languages have a way of making conditional statements and R is no exception:

```
if ([test]) {  
  [Action if test is passed]  
} else {  
  [Action if test is not passed]  
}
```

Try writing code which returns “positive” if a number is positive and “negative or zero” otherwise

# Conditional statements

```
x <- 0
if (x>0) {
  print("positive")
} else {
  print("negative or zero")
}
```

```
[1] "negative or zero"
```



# Conditional statements

You can nest if/else statements:

```
x <- 0
if (x>0) {
  print("positive")
} else if (x==0) {
  print("zero")
} else {
  print("negative")
}
```

```
[1] "zero"
```

# The `for` loop

`for` loops are useful if we need to repeat the same operation many times on different objects.

```
for (i in 1:10) {  
  print(i)  
}
```

```
[1] 1  
[1] 2  
[1] 3  
[1] 4  
[1] 5  
[1] 6  
[1] 7  
[1] 8  
[1] 9  
[1] 10
```

```
words <- c("this", "is",  
         "a", "for",  
         "loop")
```

```
for (w in words) {  
  print(w)  
}
```

```
[1] "this"  
[1] "is"  
[1] "a"  
[1] "for"  
[1] "loop"
```

# Task

Let's write a function to loop through all columns of a dataset and plot the variable if it is a Likert-scaled question.

