

INSTITUTE  
OF LAW,  
POLITICS AND  
DEVELOPMENT



Sant'Anna  
School of Advanced Studies – Pisa

# Data Analysis for the Social Sciences with R

## Data Visualization

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

# Why visualization?



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

# Today's class

1. The grammar of graphics in `ggplot2`



Sant'Anna

School of Advanced Studies – Pisa

# Today's class

1. The grammar of graphics in `ggplot2`
2. Histograms, bar charts, and scatterplots



Sant'Anna

School of Advanced Studies – Pisa

# Today's class

1. The grammar of graphics in `ggplot2`
2. Histograms, bar charts, and scatterplots
3. Exercises



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

# 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")

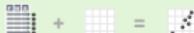


# ggplot2

## Data visualization with ggplot2 :: CHEATSHEET

### Basics

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



To do **layer**s, many variables in the data to visual properties of the geom (aesthetics), like size, color, and x and y locations.



Complete the template below to build a graph.

```
ggplot(data = DATA) +  
  GEOM_FUNCTIONS: mapping = MAPPINGS +  
  stat = STAT_FUNCTIONS +  
  position = POSITION_FUNCTIONS +  
  coordinate = COORDINATE_FUNCTIONS +  
  facet = FACET_FUNCTIONS +  
  scale = SCALE_FUNCTIONS +  
  theme = THEME_FUNCTIONS
```

ggplot(data = mpg, aes(x = cyl, y = hwy))  
plots that you finish by adding layers. To add one geom function per layer.

last\_plot() returns the last plot.

ggname("plot.png", width = 5, height = 5) saves last plot as 5" x 5" file named "plot.png" in working directory.  
Matches file type to file extension.

### Aes

Common aesthetics values.

color and fill - string ("red", "VRBGGBB")  
linetype - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "longdash", 6 = "twohead")

size - integer (in mm for size of points and text)

linewidth - integer (in mm for widths of lines)

shape - integer (shape name or a single character ("\*"))

stroke - integer (stroke width)

strokeDash - vector of length 2 (in mm)

strokeWidth - integer (stroke width)

strokeWidthDash - vector of length 2 (in mm)

strokeWidthDashOffset - integer (stroke width offset)

strokeWidthDashPhase - integer (stroke width phase)

strokeWidthDashType - integer (stroke width type)

strokeWidthDashUnit - character ("mm", "px", "pt")

strokeWidthDashValue - vector of length 2 (in px, pt, mm)

strokeWidthDashX - integer (stroke width x offset)

strokeWidthDashY - integer (stroke width y offset)



### Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables.  
Each function returns a layer.

#### GRAPHICAL PRIMITIVES

a <- ggplot(economics, aes(date, unemploy))

b <- ggplot(heads, aes(x = long, y = lat))

a + geom\_blank() and a + expand\_limits()

Ensure limits provide values across all plots

b + geom\_rect() and b + expand\_limits()

same as long + 1, curvature = 0, x = send, y = end,

alpha, angle, color, curvature, x, linetype, size

a + geom\_path() and a + facet

lines = "round", linestreng = 1)

same as long + 1, curvature = 0, x = send, y = end,

alpha, angle, color, curvature, x, linetype, size

a + geom\_polygon(alpha = 0.5) x, y, alpha,

color, fill, group, subgroup, linetype, size

b + geom\_rect() and b + expand\_limits()

same as long + 1, ymax = lat + 1 - smax, smin,

ymax, yend, alpha, color, fill, linetype, size

a + geom\_raster() and a + expand\_limits()

width = 1000, height = 1000, x, y, alpha, color, fill, group, linetype, size

b + geom\_rect() and b + expand\_limits()

width = 1000, height = 1000, x, y, alpha, color, fill, group, linetype, size

#### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

b + geom\_abline(intercept = 0, slope = 1)

b + geom\_abline(intercept = 0, slope = 0)

b + geom\_segment(aesend = 1, vend = long + 1)

b + geom\_spoke(aesangle = 11.25, radius = 2)

continuous

ONE VARIABLE continuous

c <- ggplot(mpg, aes(wt))

c + geom\_area() and c + geom\_ribbon()

c + geom\_density(kernel = "gaussian")

c + geom\_dotplot(binaxis = "y", stackdir = "center")

c + geom\_hex()

c + geom\_jitter()

c + geom\_linerer()

c + geom\_point()

c + geom\_rect()

c + geom\_raster()

c + geom\_rect()

## 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)



# 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

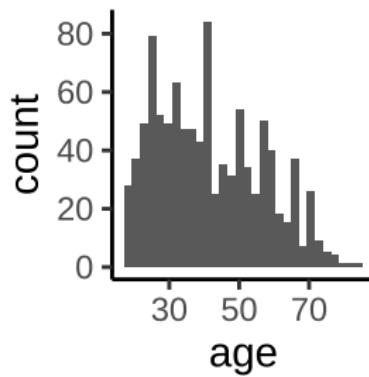
- ▶ A **histogram** is a type of chart representing the distribution of a single numerical variable.
- ▶ In histograms, the x-axis is divided into **bins** which represent ranges of the variable.
- ▶ The height of the bars represent **how many observations** fall within a bin.

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

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

# Histograms

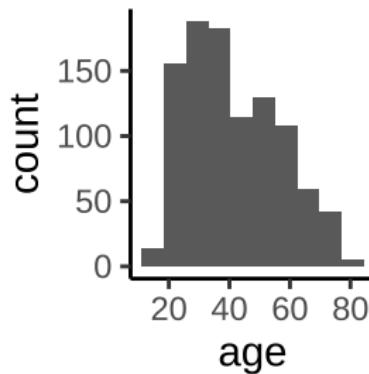
```
plot <- ggplot(data=tun22,  
                 aes(x=age)) +  
  geom_histogram() +  
  theme_classic()  
  
plot
```



# Histograms

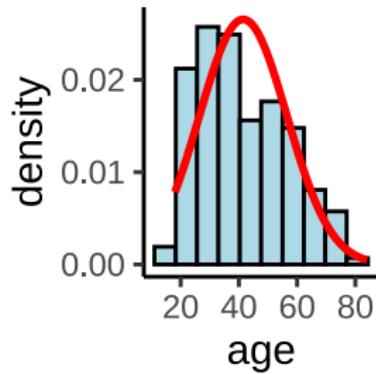
You can set the number of bins with the `bin` option (the default is 30):

```
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

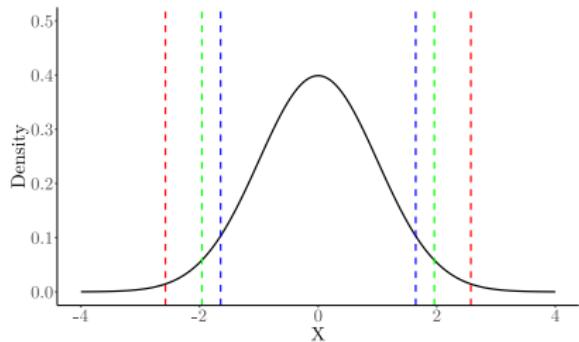


## Quick detour: What is density?

- ▶ 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**.



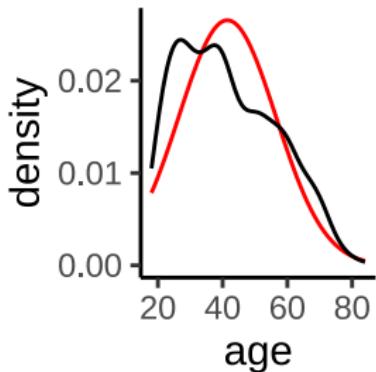
# The normal distribution



- ▶ 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

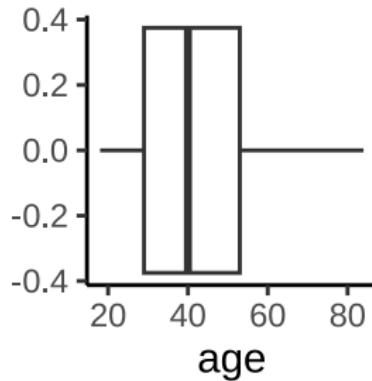
# Density plots

```
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
```



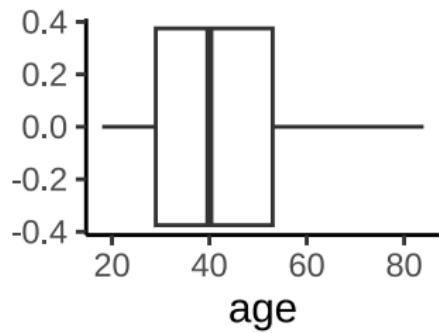
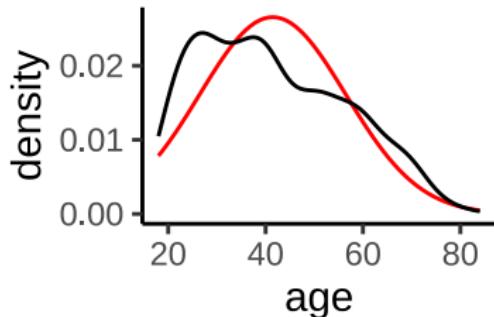
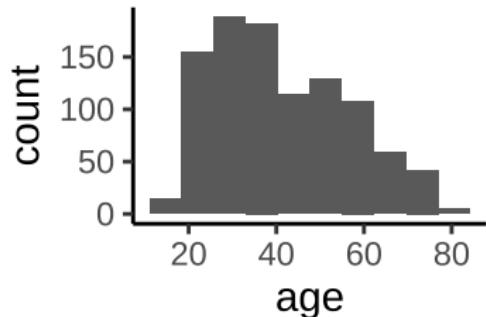
# Boxplots

```
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:

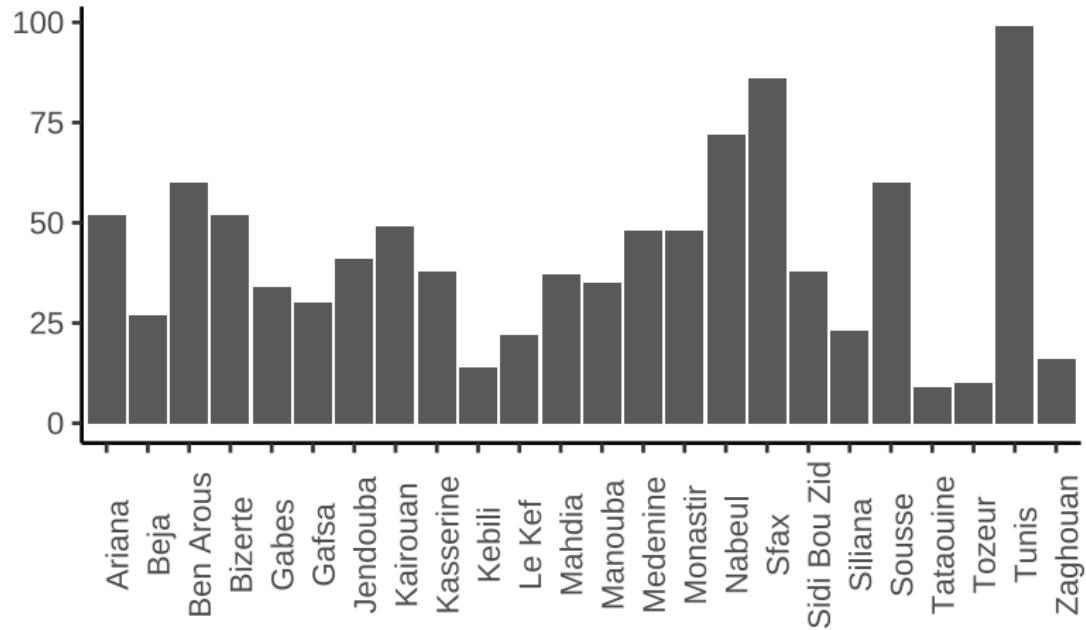
```
tun22 %>%
  group_by(region) %>%
  summarize(n = n()) %>%
  ggplot() +
  geom_bar(
    aes(x = region, y = n),
    stat = "identity"
  ) +
  labs(x="",
       y="",
       title="Number of respondents by region") +
  theme_classic() +
  theme(axis.text.x = element_text(angle = 90))
```

What will this graph look like?



# Bar charts

## Number of respondents by region



# ggplot2

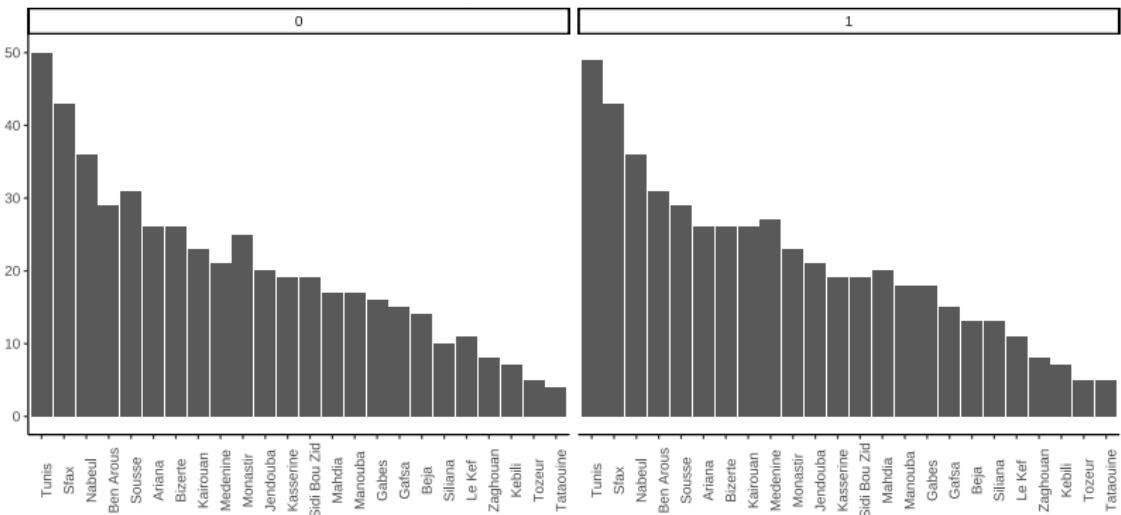
We can also produce different plots by group using the `facet()` layer

```
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)
```

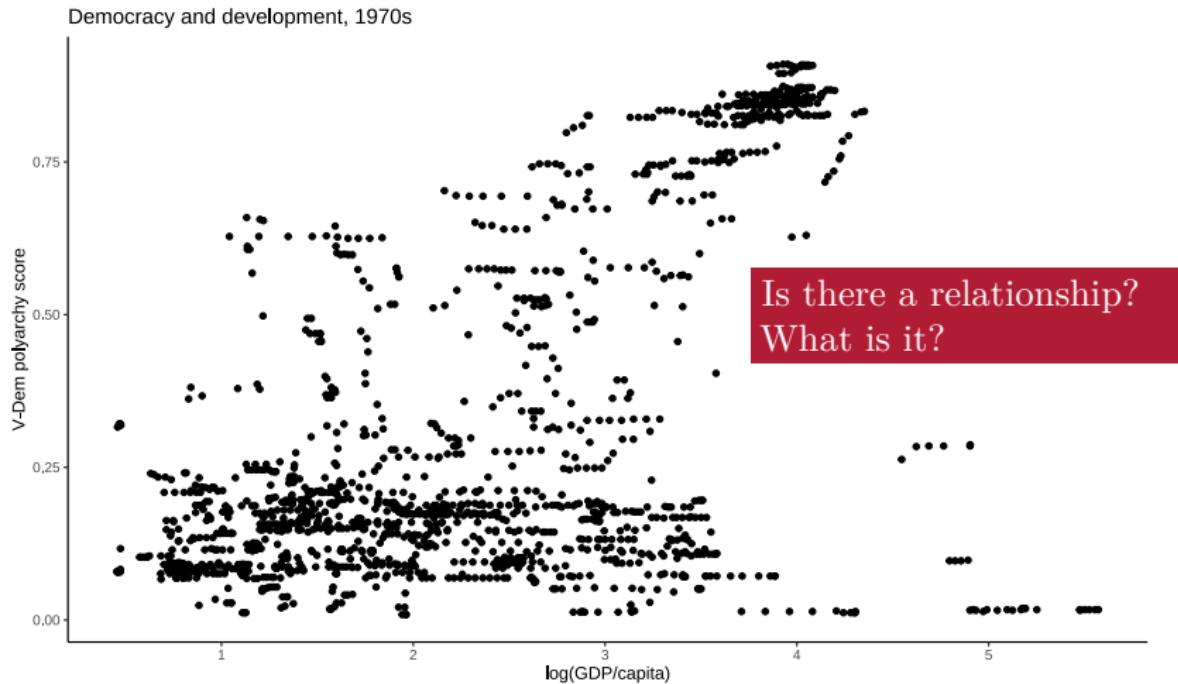


Let's create a scatterplot showing the relationship between log GDP/capita and the polyarchy score for the 1970s

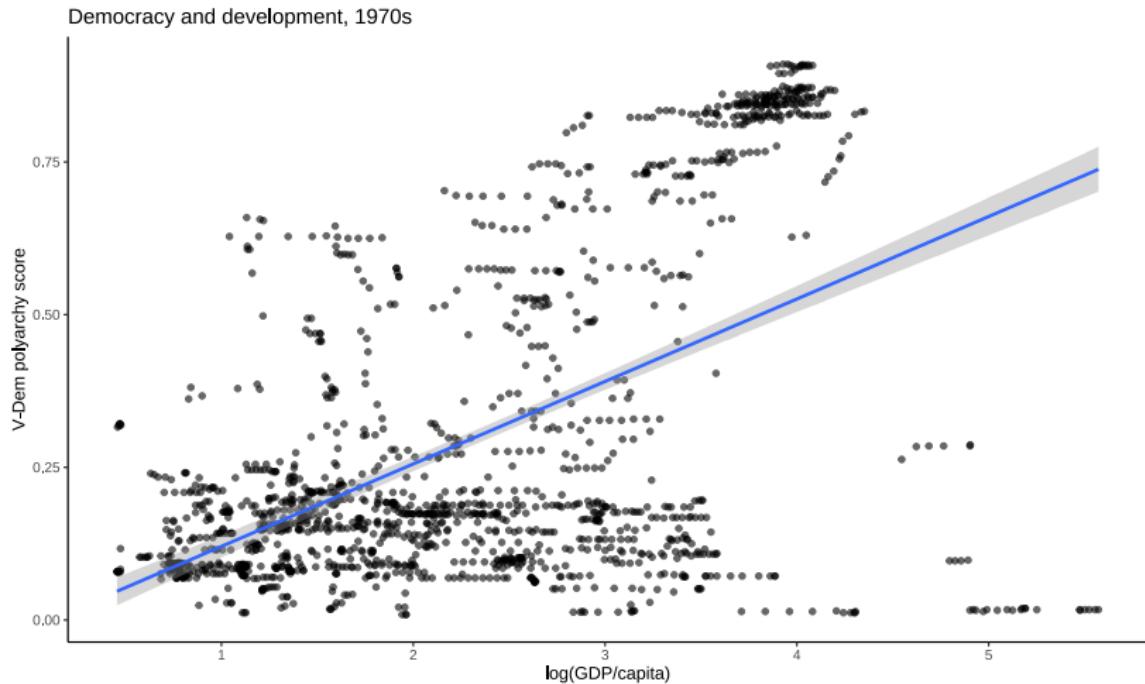
# Scatterplots



# Scatterplots



# Scatterplots

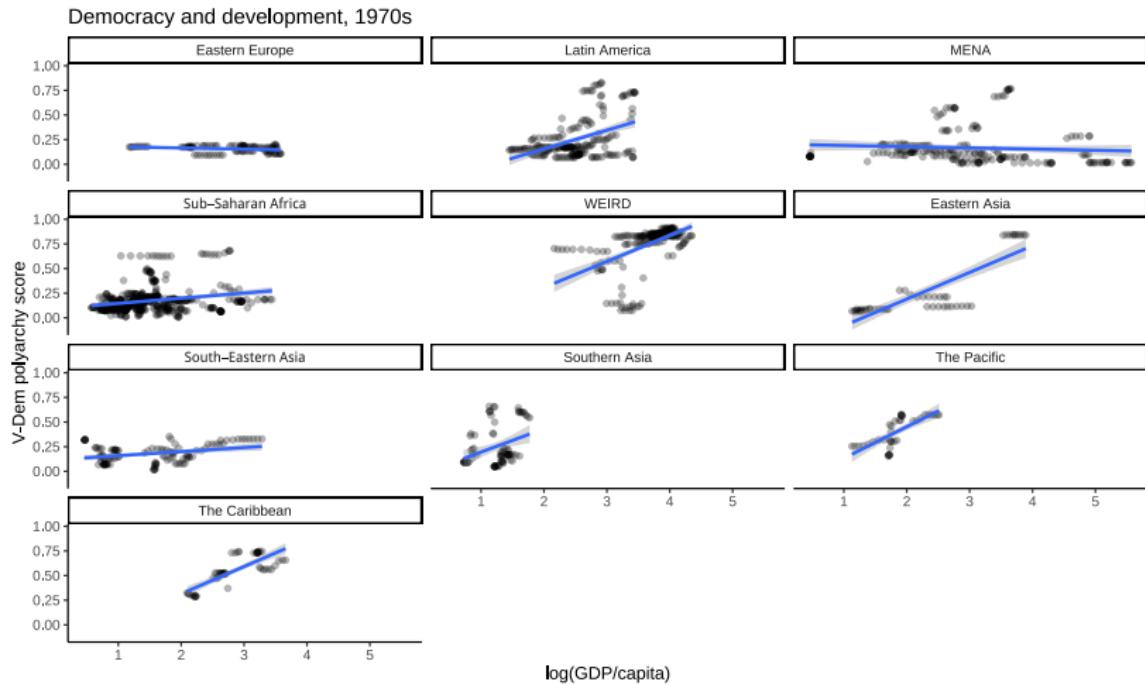


# Scatterplots

```
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"
)
```

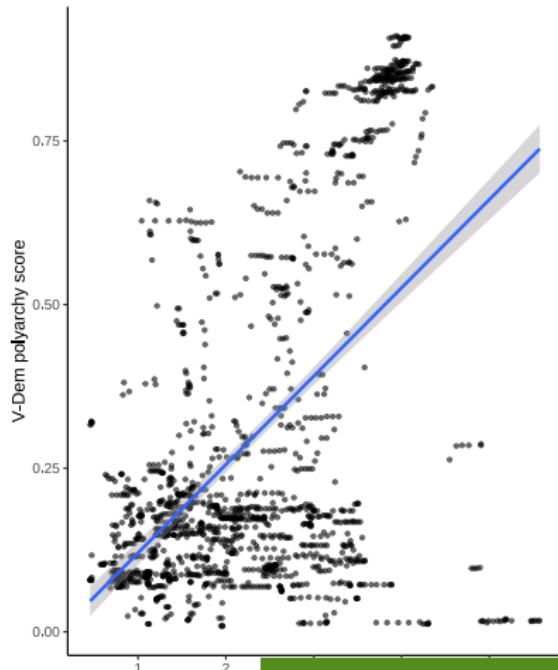


Sant'Anna  
School of Advanced Studies – Pisa

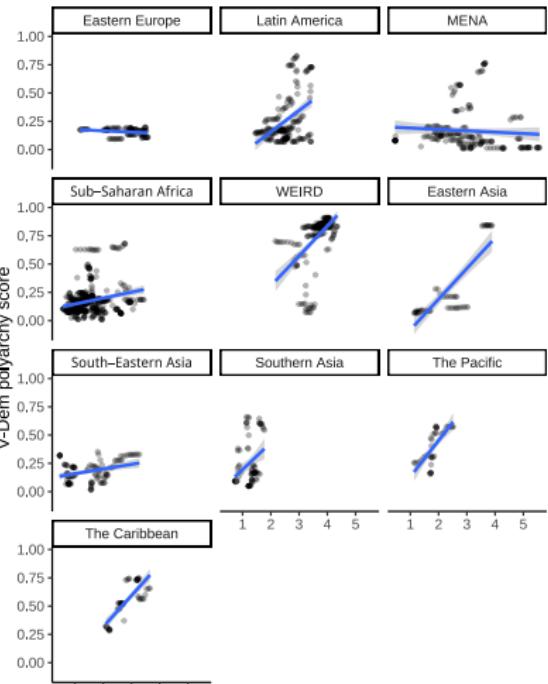
# Combining graphs with the patchwork package

Democracy and Development, 1970s

All countries



By region



Sant'Anna  
School of Advanced Studies – Pisa

# Exercises



**Sant'Anna**  
School of Advanced Studies – Pisa

# Exercises

Recreate the plots on the GitHub page.



Sant'Anna

School of Advanced Studies – Pisa