

Week 02

Tidy Data and Visualization

API209: Summer Math Camp

Rony Rodrigo Maximiliano Rodriguez-Ramirez

rrodriguezramirez@g.harvard.edu

Harvard University

August 20, 2024

Recap and Tidy Data

Wrangling your data {Recap}

- You are **highly encouraged** to read through **Hadley Wickham's chapter**. It's clear and concise.
- Also check out this great "cheatsheet" **here**.
- The package is organized around a set of **verbs**, i.e. *actions* to be taken.
- All *verbs* work as follows:

$$\text{verb}(\underbrace{\text{data.frame}}_{\text{1st argument}}, \underbrace{\text{what to do}}_{\text{2nd argument}})$$

- Alternatively you can (should) use the **pipe** operator **%>%**:

$$\underbrace{\text{data.frame}}_{\text{1st argument}} \underbrace{\%>\%}_{\text{"pipe" operator}} \text{verb}(\underbrace{\text{what to do}}_{\text{2nd argument}})$$

Tidy data

- In most cases, your datasets won't be **tidy**.

Tidy data: A dataset is said to be tidy if it satisfies the following conditions:

“**TIDY DATA** is a standard way of mapping the meaning of a dataset to its structure.”
—HADLEY WICKHAM

In tidy data:

- each **variable** forms a **column**
- each **observation** forms a **row**
- each **cell** is a **single measurement**

each column a variable

id	name	color
1	floof	gray
2	max	black
3	cat	orange
4	donut	gray
5	merlin	black
6	panda	calico

each row an observation

Wickham, H. (2014). Tidy Data. Journal of Statistical Software 59 (10). DOI: 10.18637/jss.v059.i10

Untidy data is pretty common

CITIZENSHIP	SOUTHWEST BORDER									
	BBT	DRT	ELC	EPT	LRT	RGV	SDC	TCA	YUM	SBO Total
AFGHANISTAN									1	1
ALBANIA				4		9		3		16
ALGERIA										0
ANGOLA		262	2							264
ANGUILLA				1						1
ARGENTINA	1	3				3		1	1	9
ARMENIA			4				1		1	6
AUSTRALIA										0
AZERBAIJAN										0
BAHAMAS										0
BANGLADESH		11	502		2	31	31		67	644
BELARUS		1								1
BELGIUM										0
BELIZE	1	3		5	1	22	1	3	2	38
BENIN		9	1				2		2	14
BOLIVIA		1		4	3	8				16
BRAZIL	9	347	392	5,185	47	143	337	13	473	6,946
BULGARIA				1						1
BURKINA FASO		3	1				7			11

However, storing data in wide form is easier to display in a printed table.

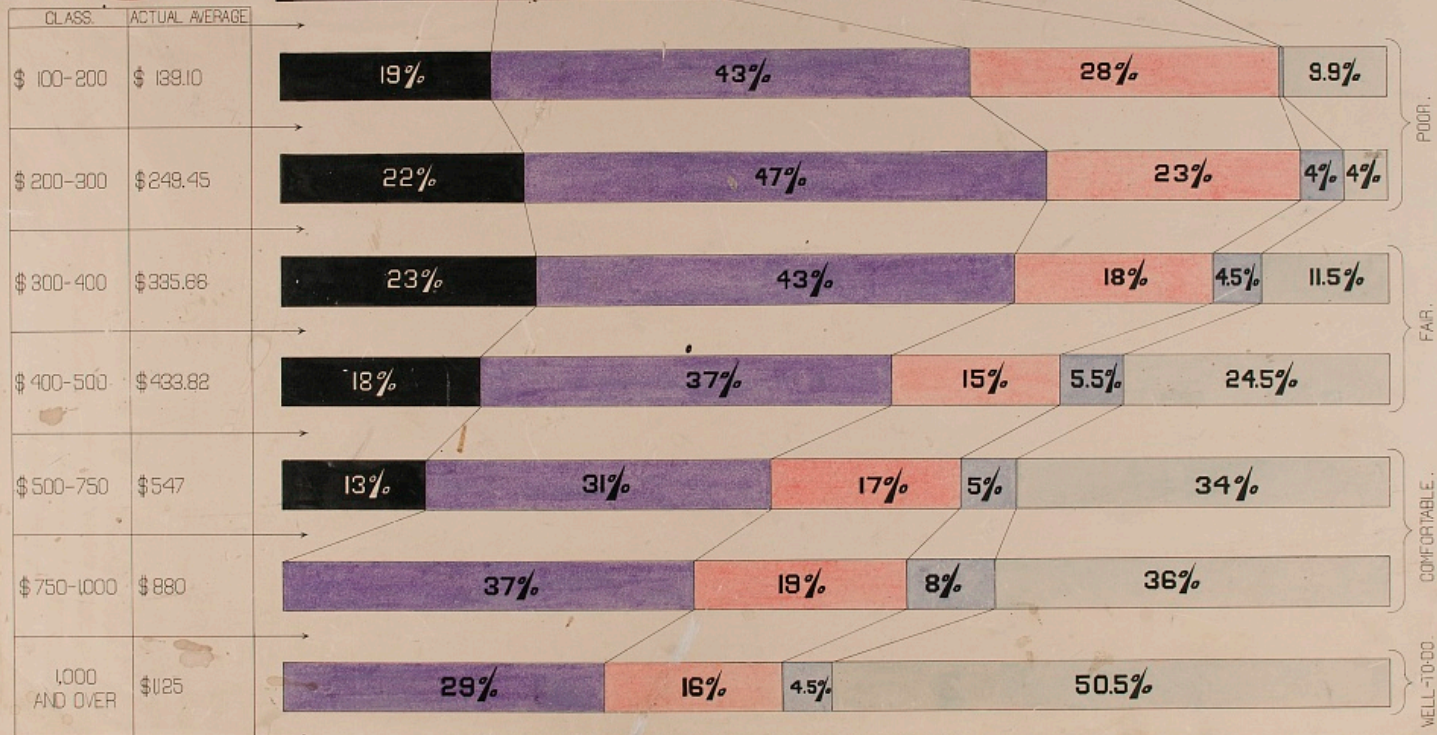
Tidy data
is data in
long format

Beautiful visualizations

INCOME AND EXPENDITURE OF 150 NEGRO FAMILIES IN ATLANTA, GA., U.S.A.



ANNUAL EXPENDITURE FOR				
RENT.	FOOD.	CLOTHES.	DIRECT TAXES.	OTHER EXPENSES AND SAVINGS.
	<p>DETAILED LIST OF WELL-TO-DO NEGRO FAMILIES FROM BULLETIN U.S. DEPARTMENT OF AGRICULTURE NO. 71.</p>		<p>THE STATE TAX RATE IS: 1880-\$3.50 PER \$1,000 1885-\$3.50 1890-\$3.98 1895-\$4.56 1899-\$5.38 STATE AND COUNTY TAXES RAISE THIS TO \$21 PER \$1,000 IN ATLANTA.</p>	<p>THE HIGHER LIFE. RELIGION. ART. EDUCATION. SICKNESS. SAVINGS. AMUSEMENTS. BOOKS AND PAPERS. TRAVEL.</p>



FOR FURTHER STATISTICS RAISE THIS FRAME.

What makes a great visualization?

Truthful

Functional

Beautiful

Insightful

Enlightening

Alberto Cairo, *The Truthful Art*

How do we express visuals in words?

- **Data** to be visualized
- **Geometric objects** that appear on the plot
- **Aesthetic mappings** from data to visual component
- **Statistics** transform data on the way to visualization
- **Coordinates** organize location of geometric objects
- **Scales** define the range of values for aesthetics
- **Facets** group into subplots

What makes a great visualization?

Good aesthetics

No substantive issues

No perceptual issues

Honesty + good judgment

Kieran Healy, *Data Visualization: A Practical Introduction*

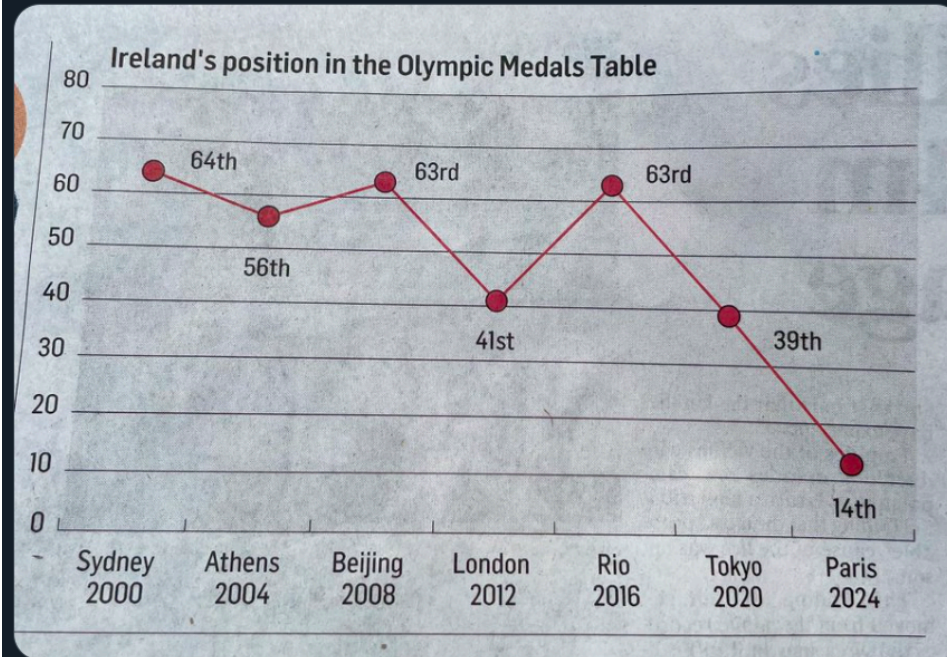
You see bad plots everywhere: What's wrong?



Dr Eemer Eivers

@EemerEivers

This is a masterclass (yet again) in @IrishTimes on how NOT to present information. Quick glance & you'd be sure 🇮🇪 has tumbled DOWN the medal tables since 2000. FFS. This is criminal level breaking of the rules of how to share data.

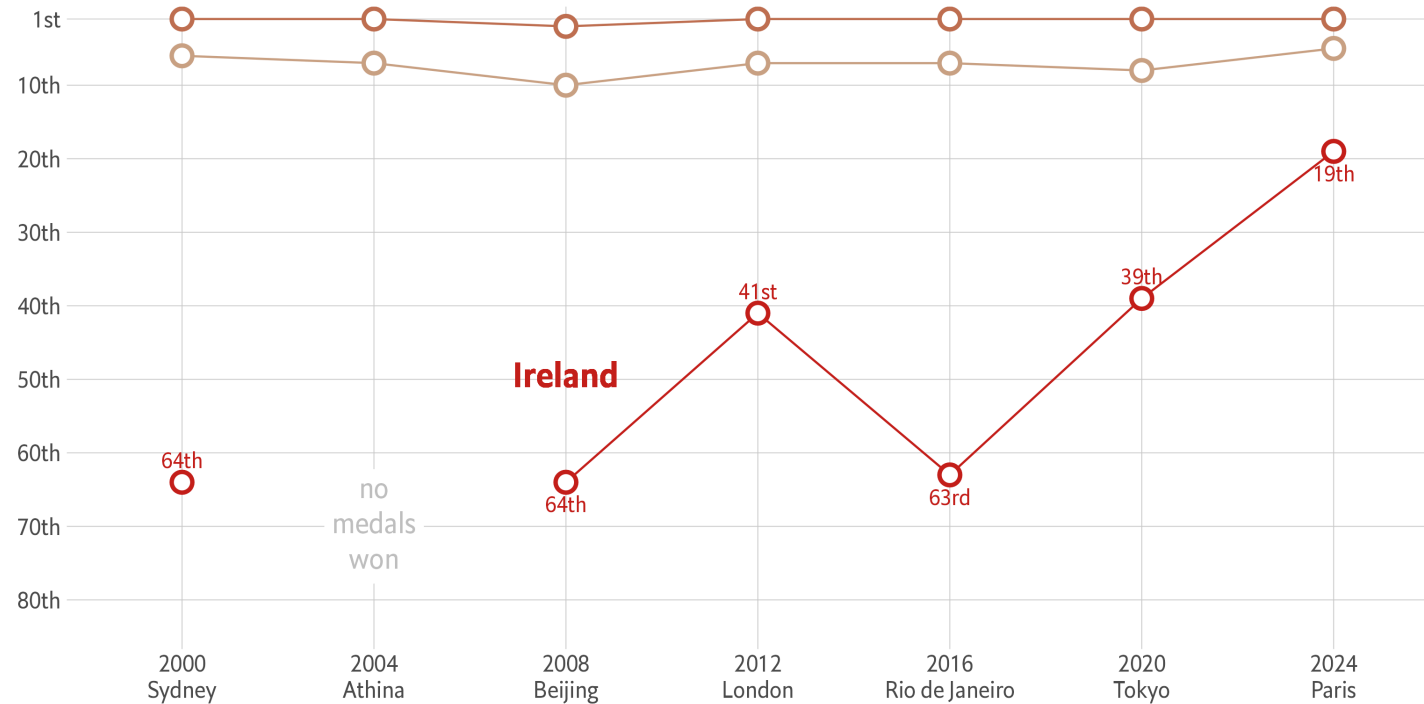


7:50 AM · Aug 10, 2024 · 868K Views

Is this right?

Ireland's position in the Olympics Medals Table

Compared to the position of the **United States** and **France**



Source: Wikipedia via Kaggle · Created with GGPlot · Original Chart: @lisacmuth · This Chart: @rrmaximiliano

Entering ggplot

ggplot

For this session, you'll use the ggplot2 package from the tidyverse meta-package.

- So, you can just load the `tidyverse` package when using ggplot.

1. Consistency with the **Grammar of Graphics**

- This book is the foundation of several data viz applications: `ggplot2`, `polaris-tableau`, `vega-lite`

2. Flexibility

3. Layering and theme customization

4. Community

It is a powerful and easy to use tool (once you understand its logic) that produces complex and multifaceted plots.

ggplot2: basic structure (template)

The basic ggplot structure is:

```
ggplot(data = DATA) +  
  GEOM_FUNCTION(mapping = aes(AESTHETIC MAPPINGS))
```


Mapping

Mappings do not directly specify the particular, e.g., colors, shapes, or line styles that will appear on the plot.

Rather, they establish which variables in the data will be represented by which visible elements on the plot.

ggplot2: full structure

```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(  
    mapping = aes(<MAPPINGS>),  
    stat = <STAT>,  
    position = <POSITION>  
  ) +  
  <COORDINATE_FUNCTION> +  
  <FACET_FUNCTION> +  
  <SCALE_FUNCTION> +  
  <THEME_FUNCTION>
```

1. **Data**: The data that you want to visualize
2. **Layers**: geom_ and stat_ → The geometric shapes and statistical summaries representing the data
3. **Aesthetics**: aes() → Aesthetic mappings of the geometric and statistical objects
4. **Scales**: scale_ → Maps between the data and the aesthetic dimensions
5. **Coordinate system**: coord_ → Maps data into the plane of the data rectangle
6. **Facets**: facet_ → The arrangement of the data into a grid of plots
7. **Visual themes**: theme() and theme_ → The overall visual defaults of a plot

ggplot2: decomposition

There are multiple ways to structure plots with ggplot

For this presentation, I will stick to Thomas Lin Pedersen's decomposition who is one of most prominent developers of the ggplot and gganimate package.

These components can be seen as layers, this is why we use the **+** sign in our ggplot syntax.



Exploratory Analysis

The most common `geoms` are:

- `geom_bar()`, `geom_col()`: bar charts.
- `geom_boxplot()`: box and whiskers plots.
- `geom_density()`: density estimates.
- `geom_jitter()`: jittered points.
- `geom_line()`: line plots.
- `geom_point()`: scatter plots.

If you want to know more about layers, you can refer to [this](#).

Step by step from Garrick Aden-Buie's gentle guide

Using the `gapminder` package, let's start with `lifeExp` vs `gdpPercap`

```
1 glimpse(gapminder)
```

```
Rows: 1,704
```

```
Columns: 6
```

```
$ country    <fct> "Afghanistan", "Afghanistan", "Afghanistan", "Afghanistan", ...
```

```
$ continent  <fct> Asia, Asia, Asia, Asia, Asia, Asia, Asia, Asia, Asia, Asia, ...
```

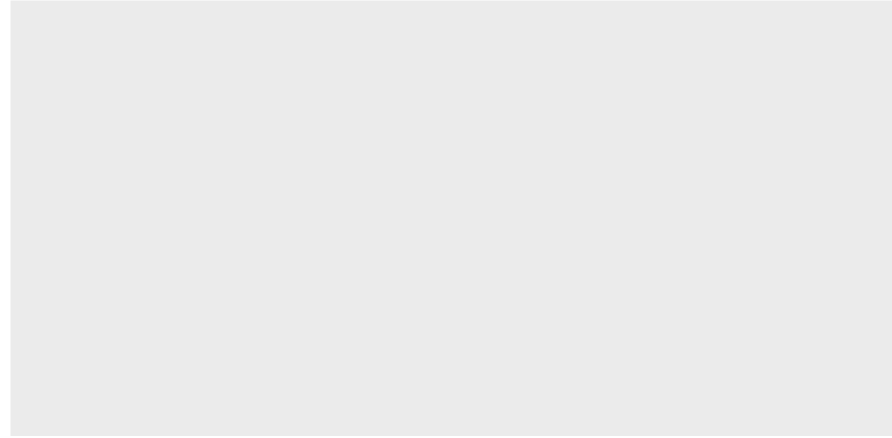
```
$ year       <int> 1952, 1957, 1962, 1967, 1972, 1977, 1982, 1987, 1992, 1997, ...
```

```
$ lifeExp    <dbl> 28.801, 30.332, 31.997, 34.020, 36.088, 38.438, 39.854, 40.8...
```

```
$ pop        <int> 8425333, 9240934, 10267083, 11537966, 13079460, 14880372, 12...
```

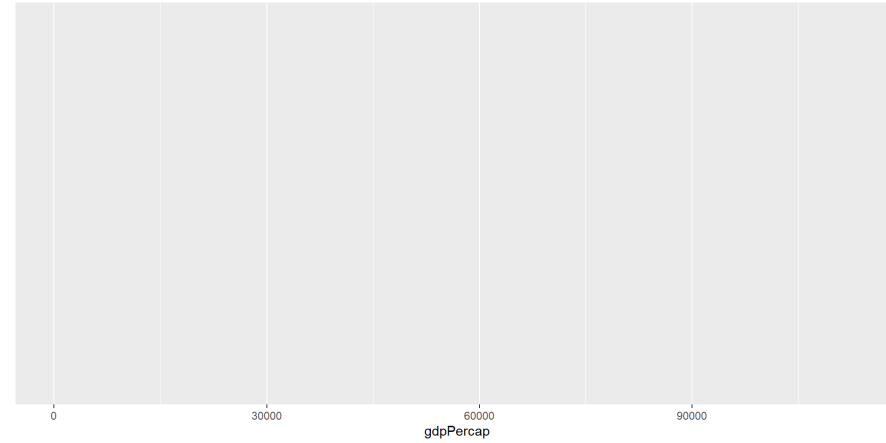
```
$ gdpPercap  <dbl> 779.4453, 820.8530, 853.1007, 836.1971, 739.9811, 786.1134, ...
```

```
1 ggplot(gapminder)
```



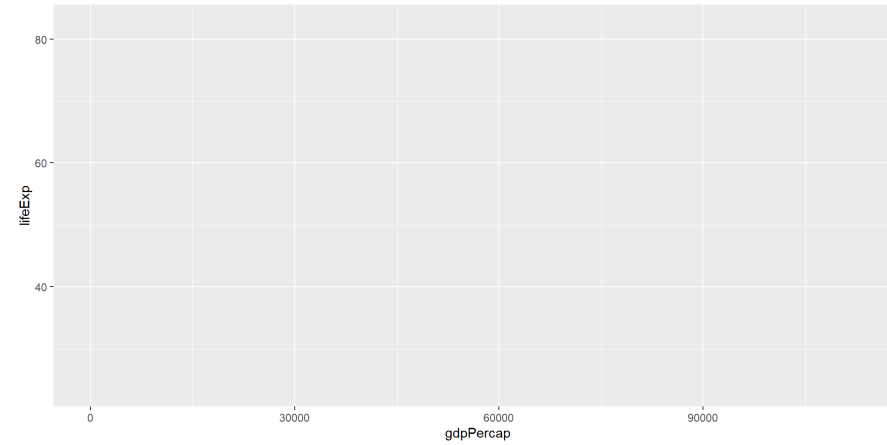
The Canvas

```
1 ggplot(gapminder) +  
2   aes(  
3     x = gdpPercap  
4   )
```



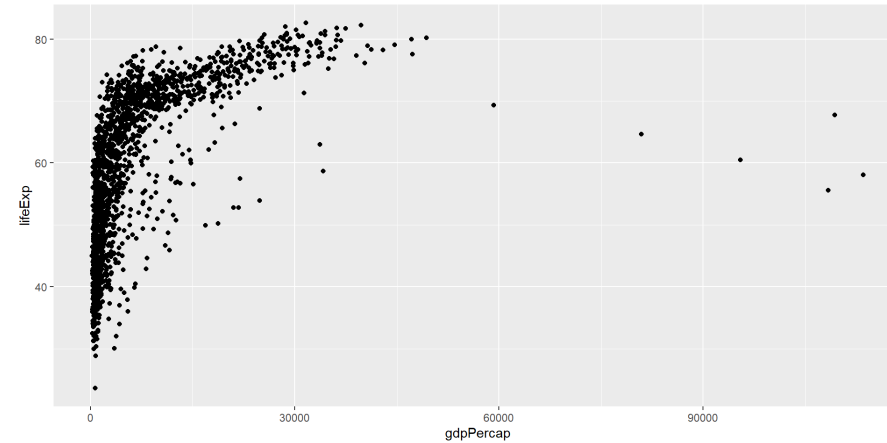
The Canvas

```
1 ggplot(gapminder) +  
2   aes(  
3     x = gdpPercap,  
4     y = lifeExp  
5   )
```



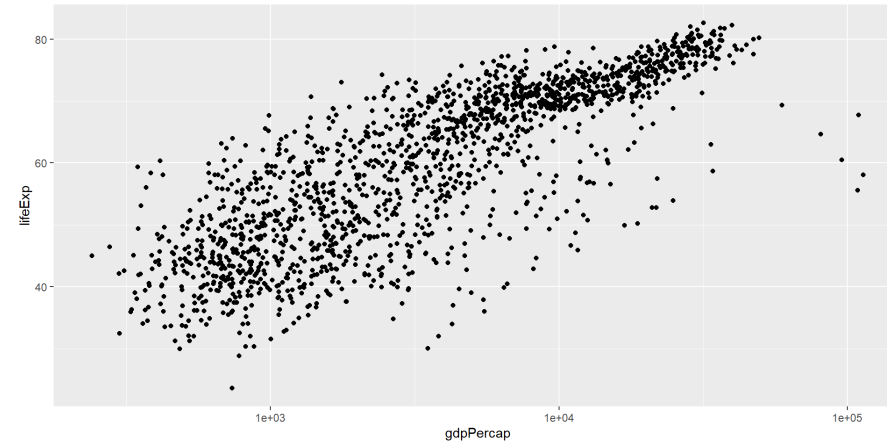
Let's add points...


```
1 ggplot(gapminder) +  
2   aes(  
3     x = gdpPercap,  
4     y = lifeExp  
5   ) +  
6   geom_point()
```



How can I tell countries apart? GDP is squished together on the left

```
1 ggplot(gapminder) +  
2   aes(  
3     x = gdpPercap,  
4     y = lifeExp  
5   ) +  
6   geom_point() +  
7   scale_x_log10()
```

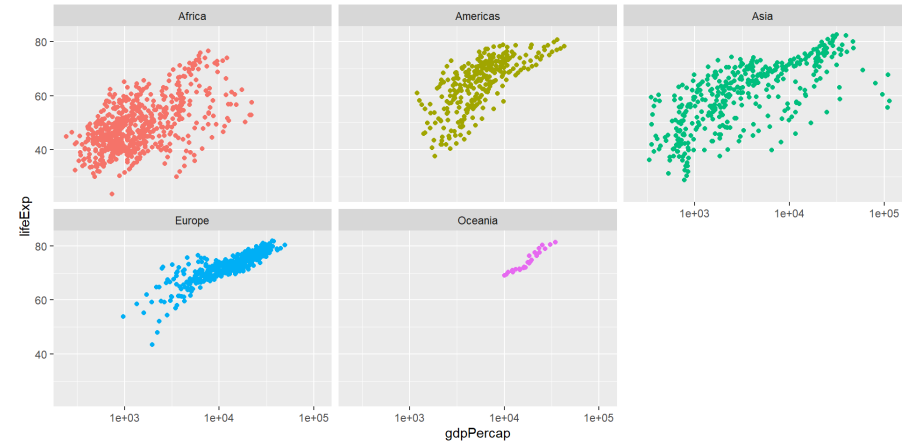


Still lots of overlap in the countries...

```

1 ggplot(gapminder) +
2   aes(
3     x = gdpPerCap,
4     y = lifeExp,
5     color = continent
6   ) +
7   geom_point() +
8   scale_x_log10() +
9   facet_wrap(~ continent) +
10  guides(color = FALSE)

```



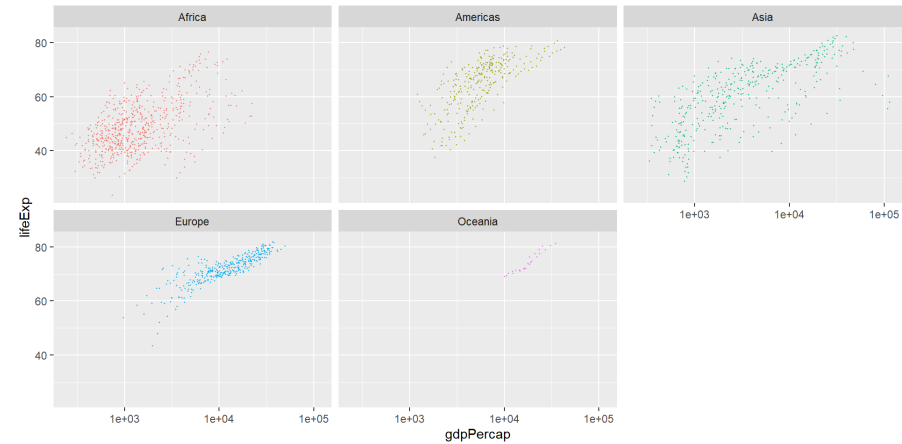
No need for color legend thanks to facet titles.

Lots of overplotting due to point size.

```

1 ggplot(gapminder) +
2   aes(
3     x = gdpPercap,
4     y = lifeExp,
5     color = continent
6   ) +
7   geom_point(size = 0.25) +
8   scale_x_log10() +
9   facet_wrap(~ continent) +
10  guides(color = FALSE)

```

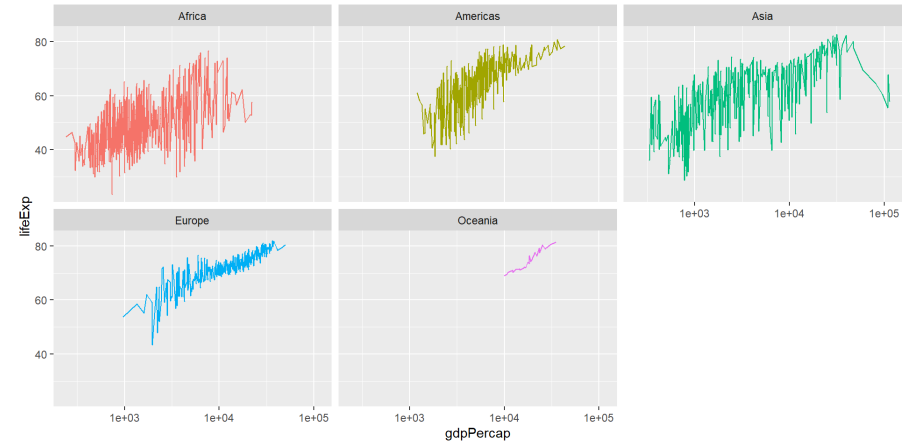


Is there a trend?

```

1 ggplot(gapminder) +
2   aes(
3     x = gdpPercap,
4     y = lifeExp,
5     color = continent
6   ) +
7   geom_line() + #<<
8   geom_point(size = 0.25) +
9   scale_x_log10() +
10  facet_wrap(~ continent) +
11  guides(color = FALSE)

```

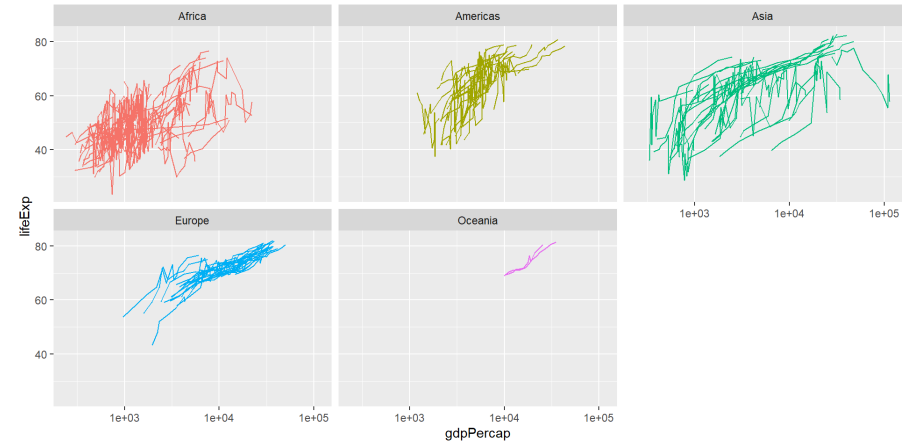


Okay, that line just connected all of the points sequentially...

```

1 ggplot(gapminder) +
2   aes(
3     x = gdpPercap,
4     y = lifeExp,
5     color = continent
6   ) +
7   geom_line(
8     aes(group = country)
9   ) +
10  geom_point(size = 0.25) +
11  scale_x_log10() +
12  facet_wrap(~ continent) +
13  guides(color = FALSE)

```

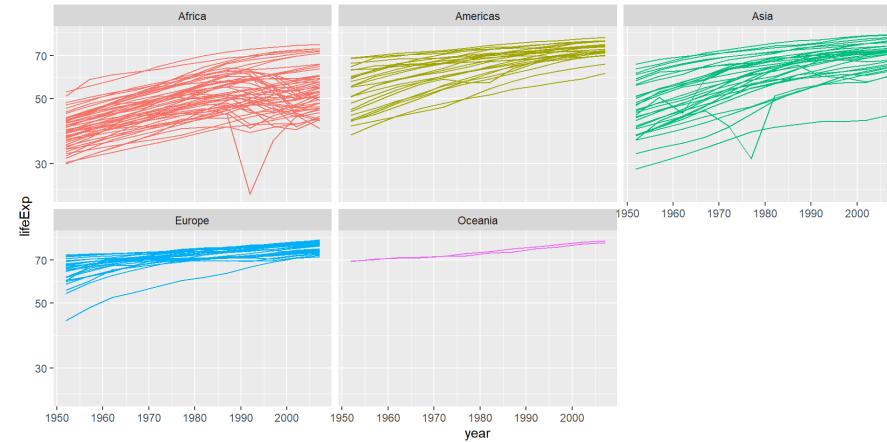


Oh no! Too confusing!

```

1 ggplot(gapminder) +
2   aes(
3     x = year,
4     y = lifeExp,
5     color = continent
6   ) +
7   geom_line(
8     aes(group = country)
9   ) +
10  geom_point(size = 0.25) +
11  scale_y_log10() +
12  facet_wrap(~ continent) +
13  guides(color = FALSE)

```

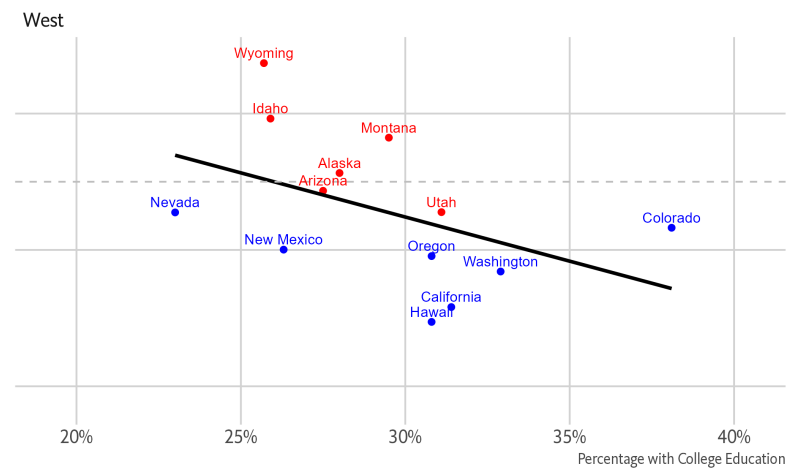
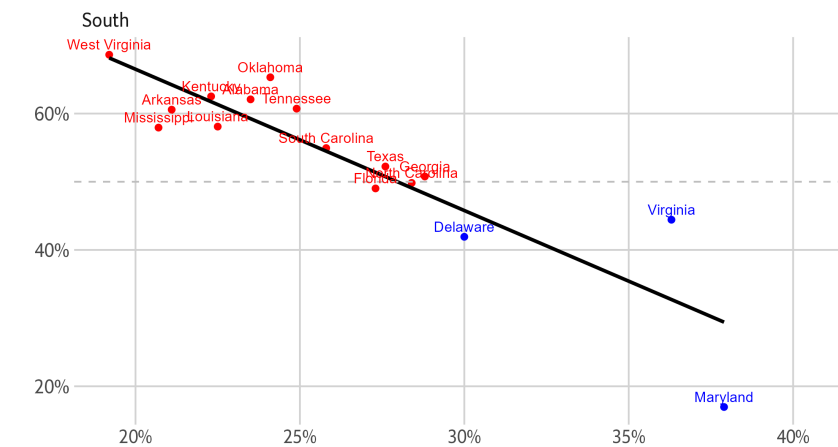
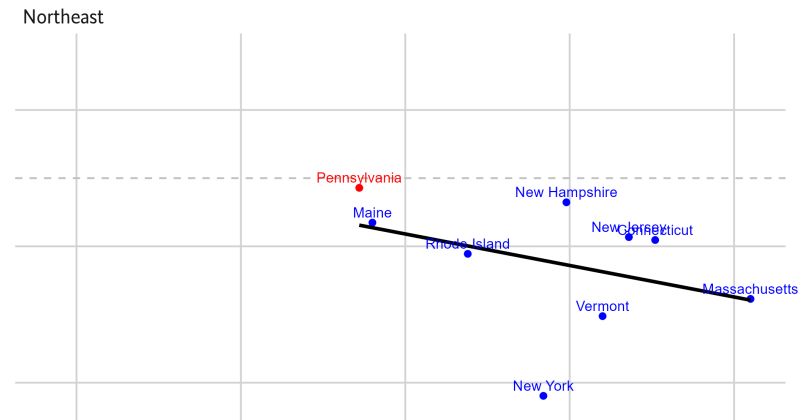
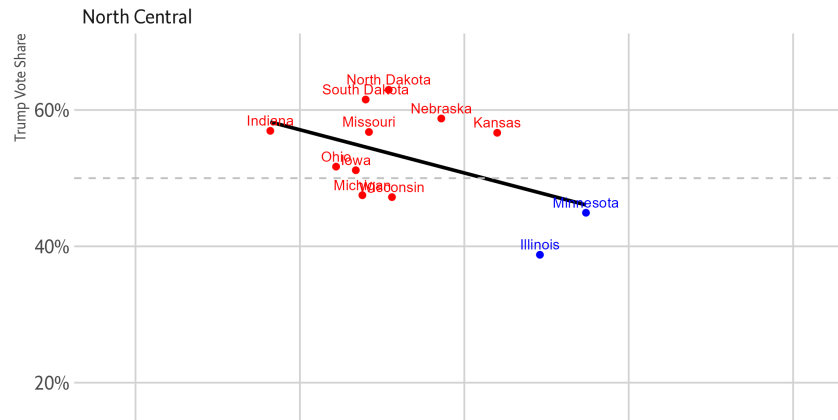


Let's add **year** in our x-axis instead of **gdp**!

Our goal

Trump Vote Share vs. College Education

Did Trump win the State? • No • Yes



Notes: Some notes here | Source: Source here | Plot by @