

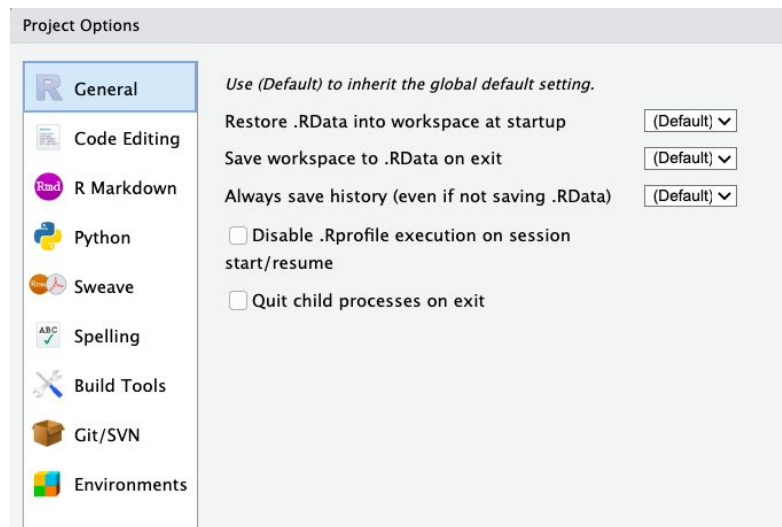
Data Visualization

Jen Cruz, Zichao Li, Isobel McEwen, Claire McLaughlin, Sanjana Srinivasan, Julien Chen



Set Up

- set up a new directory --> new project --> work in script --> save script regularly and especially before closing session
- use console for commands such as install
- use sections, section labels, and comments with command + shift + r and the #



Data Cleaning

1. Familiarize yourself with the data set
2. Check for structural errors
3. Check for data irregularities
4. Decide how to deal with missing values
5. Document data versions and changes made

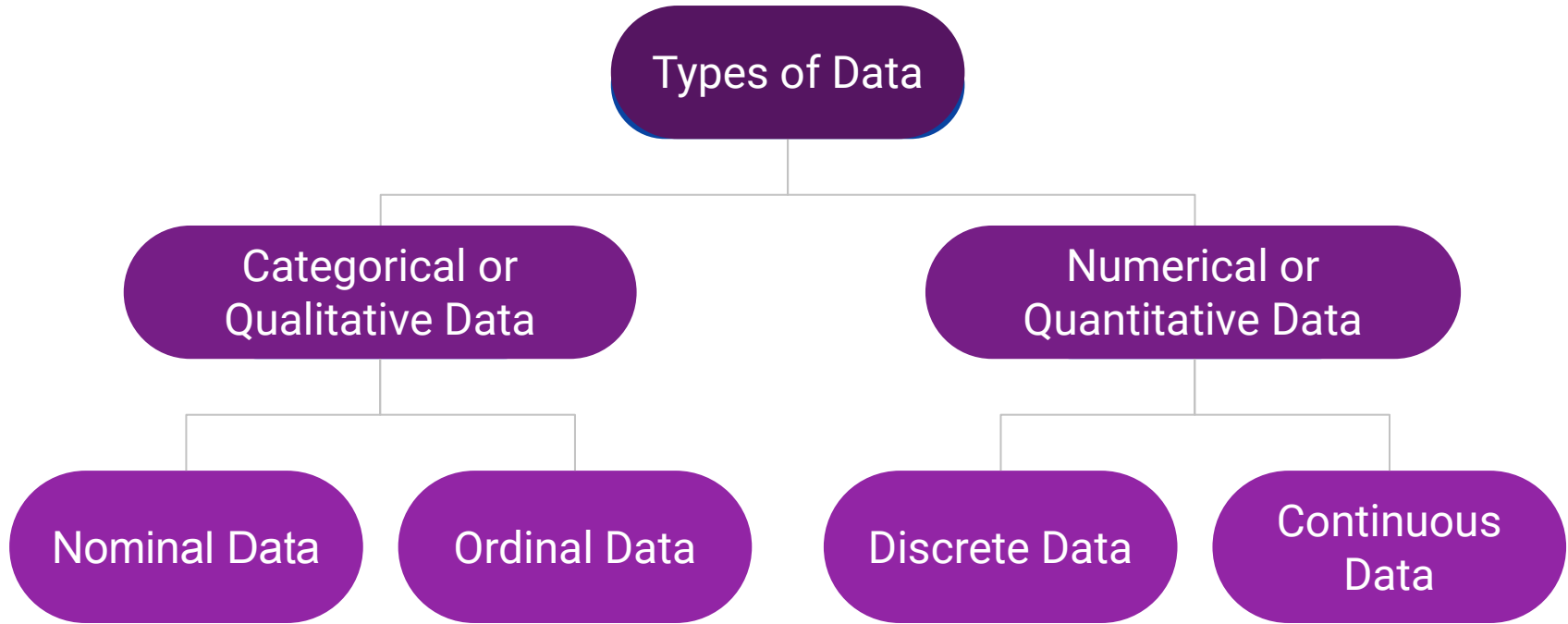
Dplyr

“Dplyr changed my life!” -Jarvis

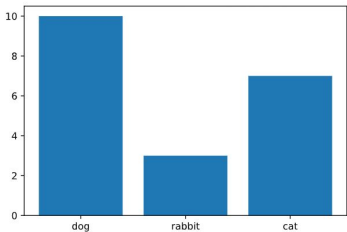
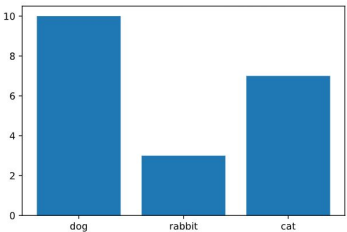
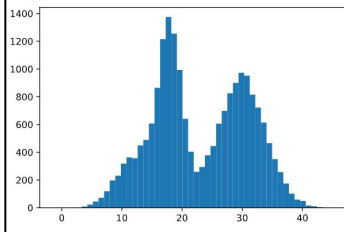
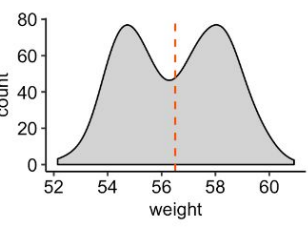
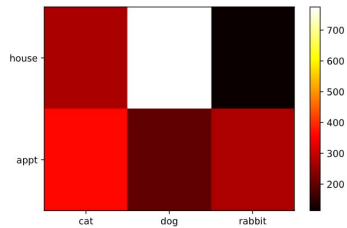
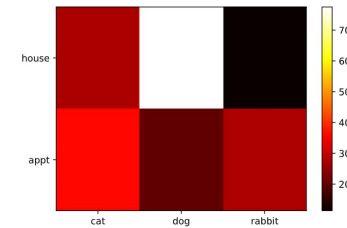
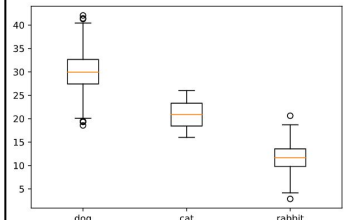
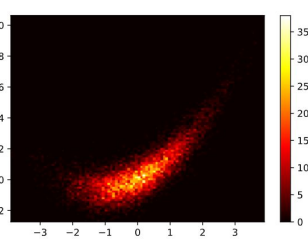
```
popsize_and_poverty <- popsize_and_poverty |>
  mutate(
    proportion_in_poverty = in_poverty / total_for_povert
  )

popsize_and_poverty <- popsize_and_poverty |>
  select(GEOID, popsize, proportion_in_poverty)
```

What type of data do you have?



Picking a way to visualize your data

| | Categorical | Ordinal | Discrete | Continuous | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|----------|--|----------|----------|------------|--------|--------|------|-----|----------|---|----------|-------|------|---|--------|-----|-----|--------|--|--|-----|----------|------|------|-----|------|--|----------|-----|----|--------|----|-----|-----|----|----|----|----|----|-----|----|----|----|----|----|--------|---|----|----|----|----|---|
| 1-D | <h3>Bar Chart</h3>  <table border="1"><thead><tr><th>Category</th><th>Count</th></tr></thead><tbody><tr><td>dog</td><td>10</td></tr><tr><td>rabbit</td><td>3</td></tr><tr><td>cat</td><td>7</td></tr></tbody></table> | Category | Count | dog | 10 | rabbit | 3 | cat | 7 | <h3>Bar Chart</h3>  <table border="1"><thead><tr><th>Category</th><th>Count</th></tr></thead><tbody><tr><td>dog</td><td>10</td></tr><tr><td>rabbit</td><td>3</td></tr><tr><td>cat</td><td>7</td></tr></tbody></table> | Category | Count | dog | 10 | rabbit | 3 | cat | 7 | <h3>Histogram</h3>  <p>A histogram showing the distribution of a discrete variable. The x-axis ranges from 0 to 40, and the y-axis (count) ranges from 0 to 1400. The distribution is bimodal, with peaks around 18 and 30.</p> | <h3>Density</h3>  <p>A density plot showing the distribution of a continuous variable. The x-axis is labeled 'weight' and ranges from 52 to 60. The y-axis is labeled 'count' and ranges from 0 to 80. The distribution is bimodal, with peaks around 54 and 58. A vertical dashed orange line is at weight 56.</p> | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Category | Count | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| dog | 10 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| rabbit | 3 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cat | 7 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Category | Count | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| dog | 10 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| rabbit | 3 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cat | 7 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2-D | <h3>Heat Maps</h3>  <table border="1"><thead><tr><th></th><th>cat</th><th>dog</th><th>rabbit</th></tr></thead><tbody><tr><th>house</th><td>High</td><td>Low</td><td>Very Low</td></tr><tr><th>appt</th><td>High</td><td>Low</td><td>High</td></tr></tbody></table> | | cat | dog | rabbit | house | High | Low | Very Low | appt | High | Low | High | <h3>Heat Maps</h3>  <table border="1"><thead><tr><th></th><th>cat</th><th>dog</th><th>rabbit</th></tr></thead><tbody><tr><th>house</th><td>High</td><td>Low</td><td>Very Low</td></tr><tr><th>appt</th><td>High</td><td>Low</td><td>High</td></tr></tbody></table> | | cat | dog | rabbit | house | High | Low | Very Low | appt | High | Low | High | <h3>Box and Whisker</h3>  <table border="1"><thead><tr><th>Category</th><th>Min</th><th>Q1</th><th>Median</th><th>Q3</th><th>Max</th></tr></thead><tbody><tr><td>dog</td><td>18</td><td>28</td><td>30</td><td>32</td><td>40</td></tr><tr><td>cat</td><td>16</td><td>18</td><td>20</td><td>22</td><td>24</td></tr><tr><td>rabbit</td><td>3</td><td>10</td><td>12</td><td>14</td><td>20</td></tr></tbody></table> | Category | Min | Q1 | Median | Q3 | Max | dog | 18 | 28 | 30 | 32 | 40 | cat | 16 | 18 | 20 | 22 | 24 | rabbit | 3 | 10 | 12 | 14 | 20 | <h3>Scatter Plot</h3>  <p>A scatter plot showing the relationship between two continuous variables. The x-axis ranges from -3 to 3, and the y-axis ranges from -2 to 10. The data points form a dense, curved, upward-sloping shape, indicating a positive correlation.</p> |
| | cat | dog | rabbit | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| house | High | Low | Very Low | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| appt | High | Low | High | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | cat | dog | rabbit | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| house | High | Low | Very Low | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| appt | High | Low | High | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Category | Min | Q1 | Median | Q3 | Max | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| dog | 18 | 28 | 30 | 32 | 40 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cat | 16 | 18 | 20 | 22 | 24 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| rabbit | 3 | 10 | 12 | 14 | 20 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

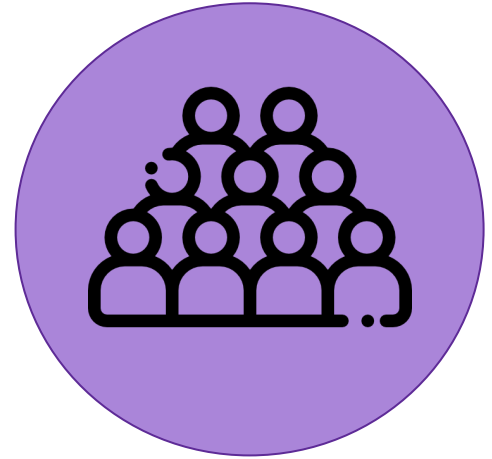
Goals of data visualization



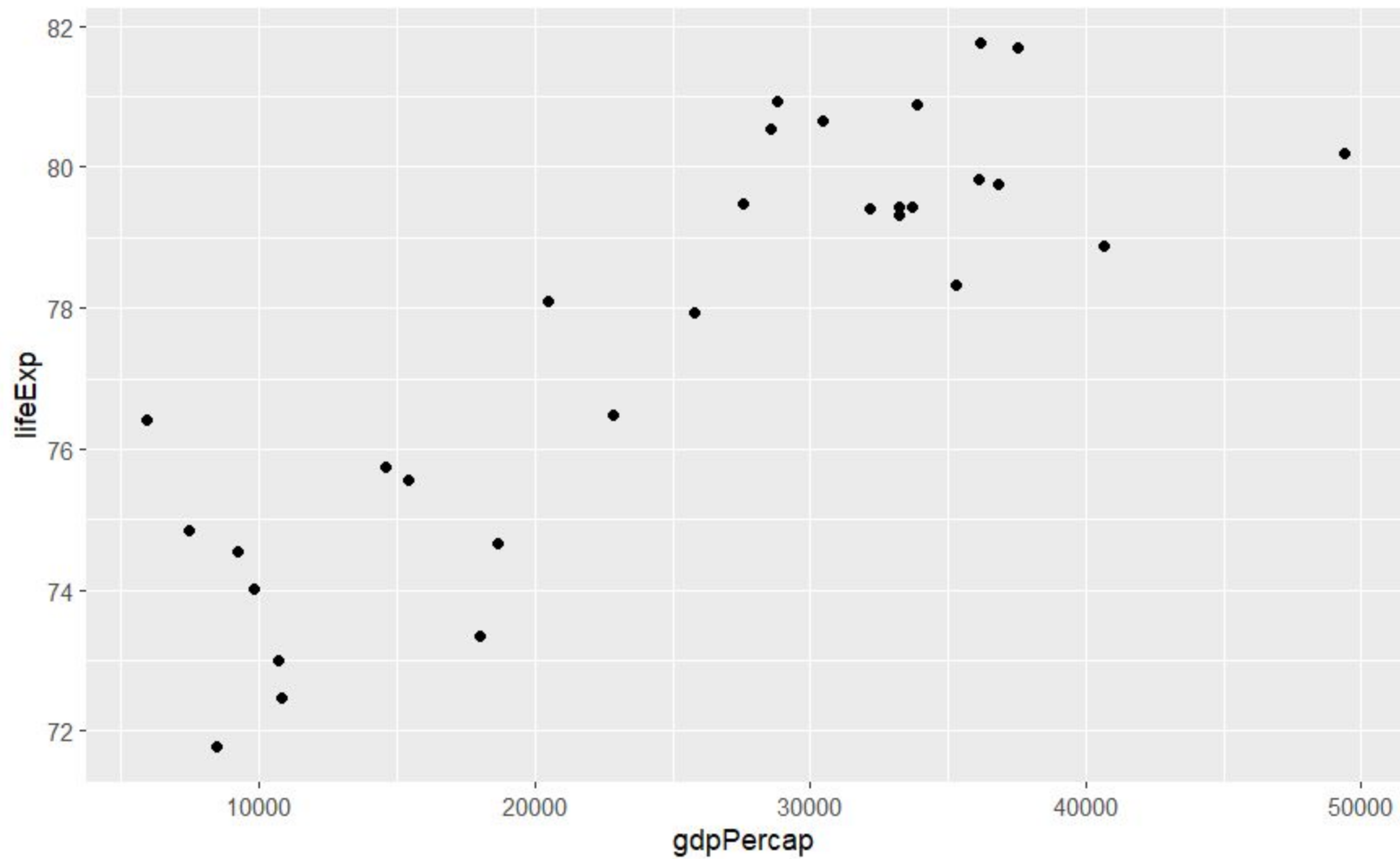
Understanding



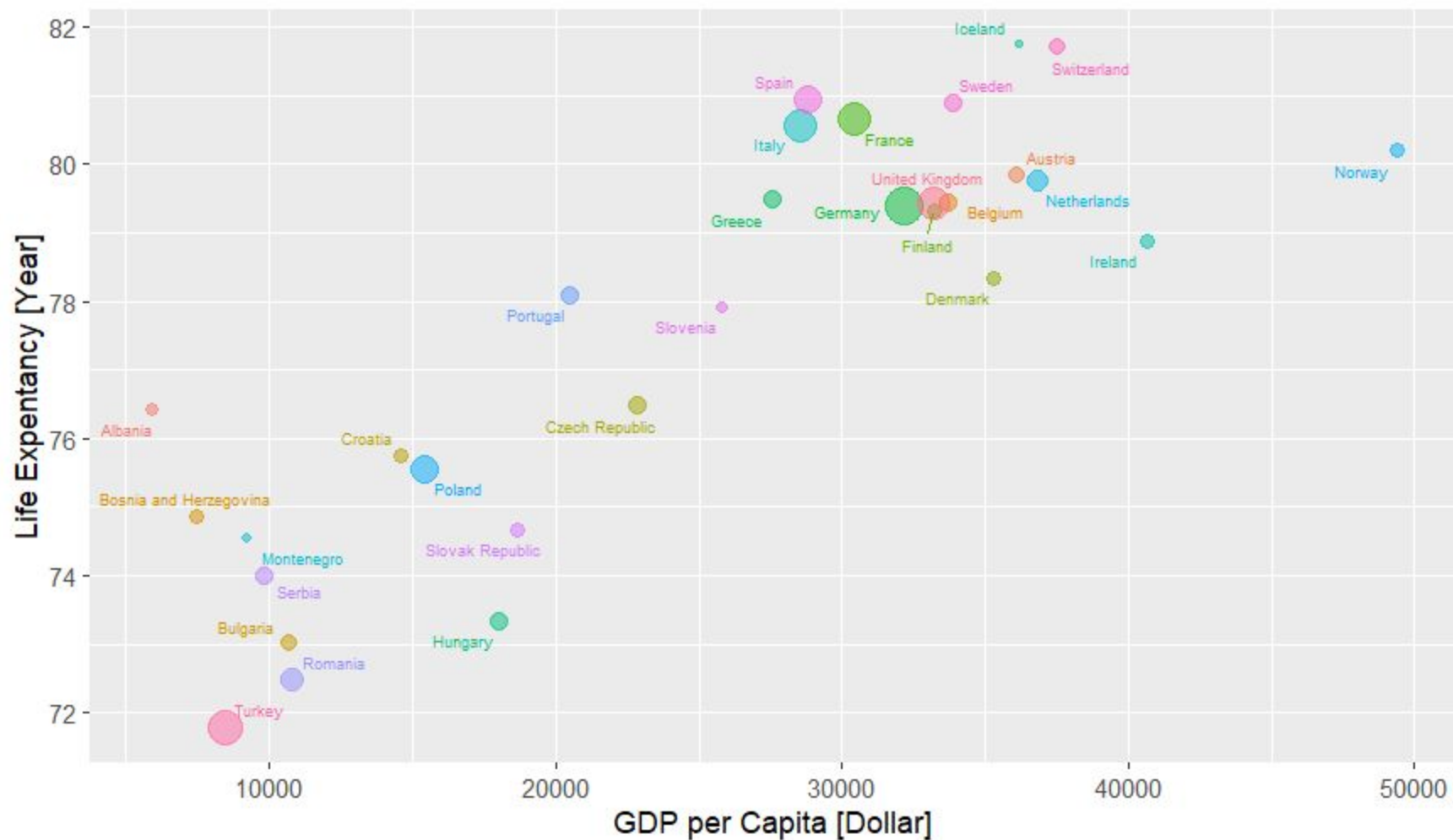
Storytelling



Accessibility



Relationship between life expectancy and GDP per Capita

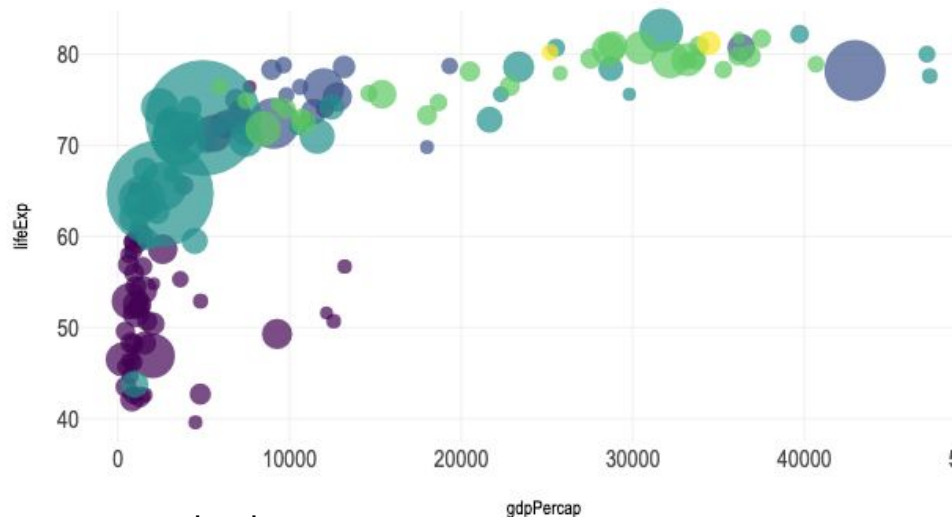
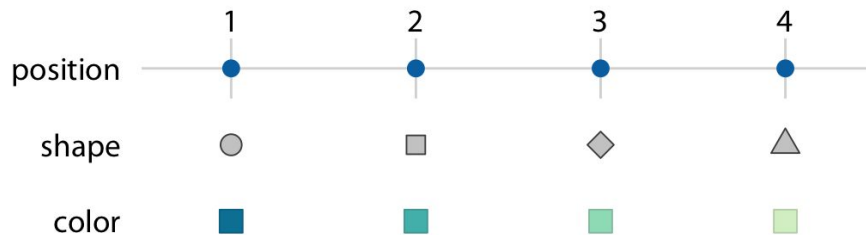


Scale

“A scale defines a unique mapping between data and aesthetics.”

Rule of thumb: Each data value must have a unique scale.

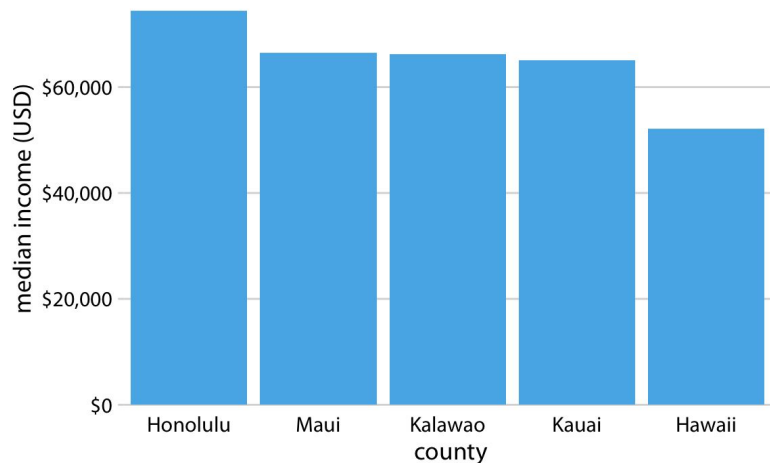
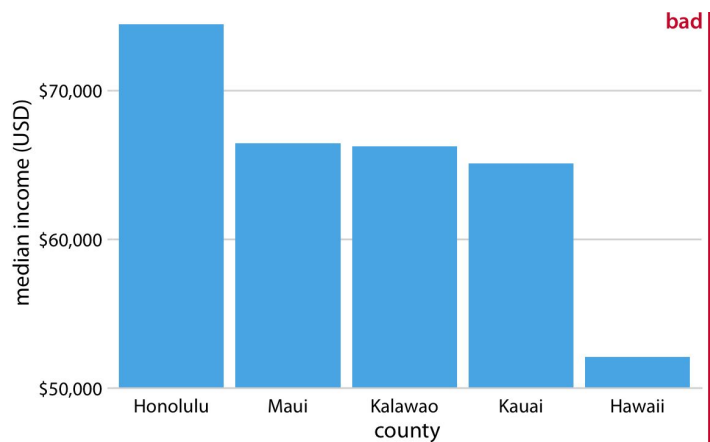
3 different scales:



Example plot:

Axis

Quantitative axes:

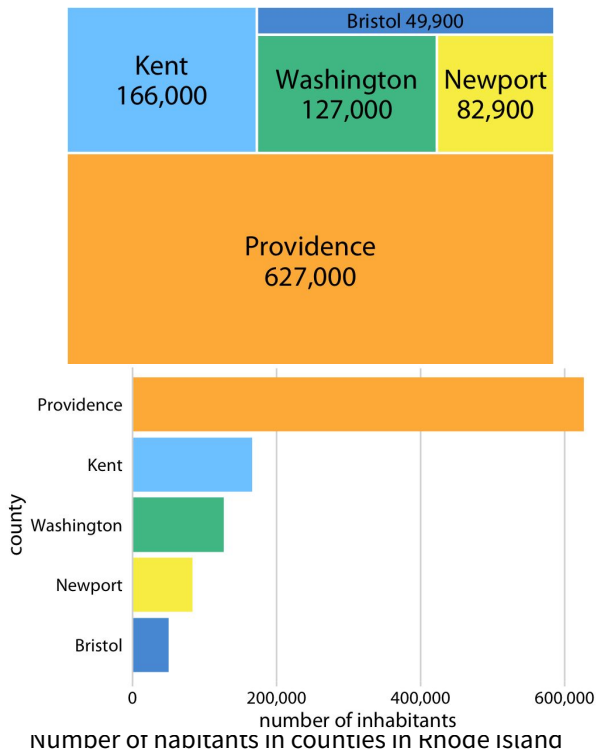


Example plot: 2 bar graphs of median income of 5 counties of Hawaii the state

Why do the income gaps look so different in the two graphs?

Axis

Categorical axes:

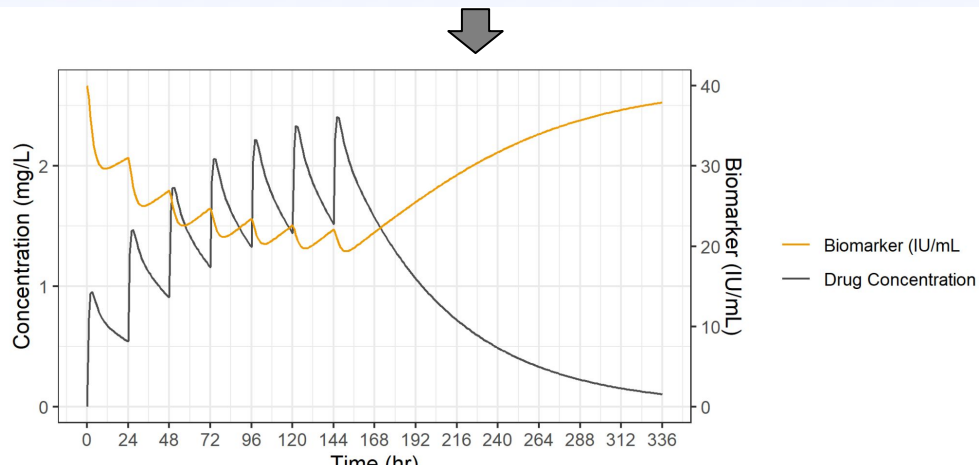


Source: <https://clauswilke.com/dataviz/proportional-ink.html>

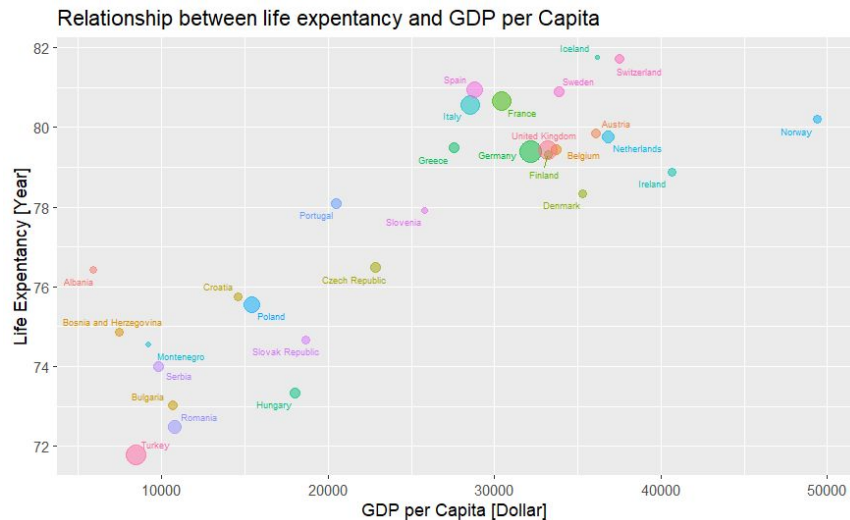
<https://finchstudio.io/blog/ggplot-dual-y-axes/>

Dual axes: two different y axes

```
scale = 15  
  
pkpd <- ggplot(res, aes(x = time, y = CP)) +  
  geom_line(aes(color = "Drug Concentration")) +  
  geom_line(aes(y = RESP/scale, color = "Biomarker (IU/mL)")) +  
  scale_x_continuous(breaks = seq(0, 336, 24)) +  
  scale_y_continuous(sec.axis = sec_axis(~.*scale, name="Biomarker (IU/mL)")) +  
  labs(x = "Time (hr)", y = "Concentration (mg/L)", color = "") +  
  scale_color_manual(values = c("orange2", "gray30"))  
  
print(pkpd)
```



Style for Accessibility



```
devtools::install_github("AndreaCirilloAC/paletter")
```

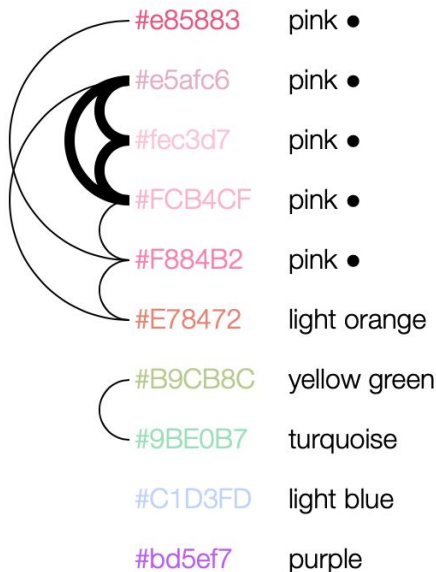
| | | | |
|---------|---------|---------|---------|
| #E8588E | #E5AFC6 | #FEC3D7 | #BD5EF7 |
| #FCB4CF | #B9CB8C | #E78472 | #F884B2 |
| #9BE0B7 | #C1D3FD | | |

Viz Palette by Elijah Meeks & Susie Lu

COLOR REPORT

Arcs link colors difficult
to tell apart as:

- Lines or small points
- Medium areas
- Large areas

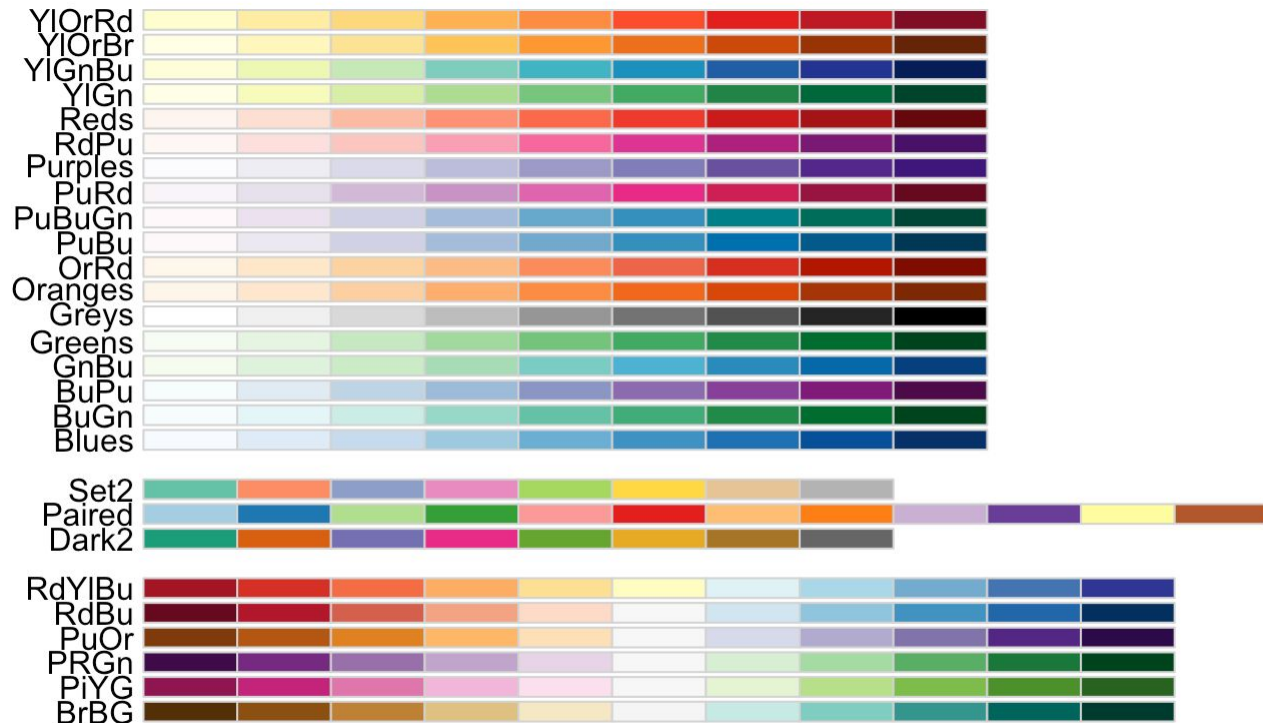


- Minimize name
conflicts for
categorical palettes

```
install.packages("RColorBrewer")
```

```
library(RColorBrewer)
```

```
display.brewer.all(colorblindFriendly = TRUE)
```



```
install.packages("wesanderson")  
library(wesanderson)
```

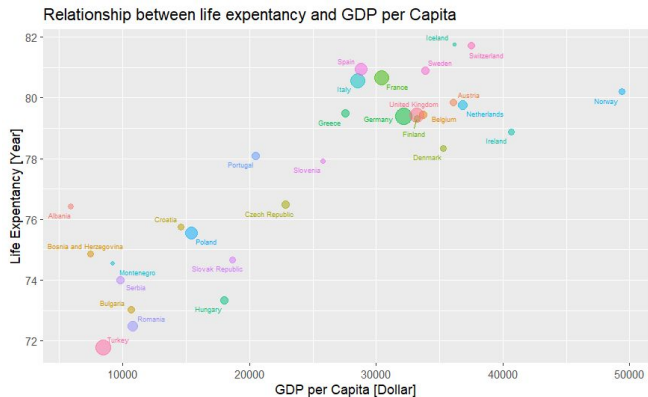
```
> names(wes_palettes)
```

```
[1] "BottleRocket1" "BottleRocket2" "Rushmore1"      "Rushmore"  
[5] "Royal1"         "Royal2"         "Zissou1"        "Darjeeling1"  
[9] "Darjeeling2"    "Chevalier1"     "FantasticFox1"   "Moonrise1"  
[13] "Moonrise2"     "Moonrise3"     "Cavalcanti1"    "GrandBudapest1"  
[17] "GrandBudapest2" "IsleofDogs1"   "IsleofDogs2"
```

```
wes_palette("FantasticFox1")
```



Style for Accessibility (cont.)

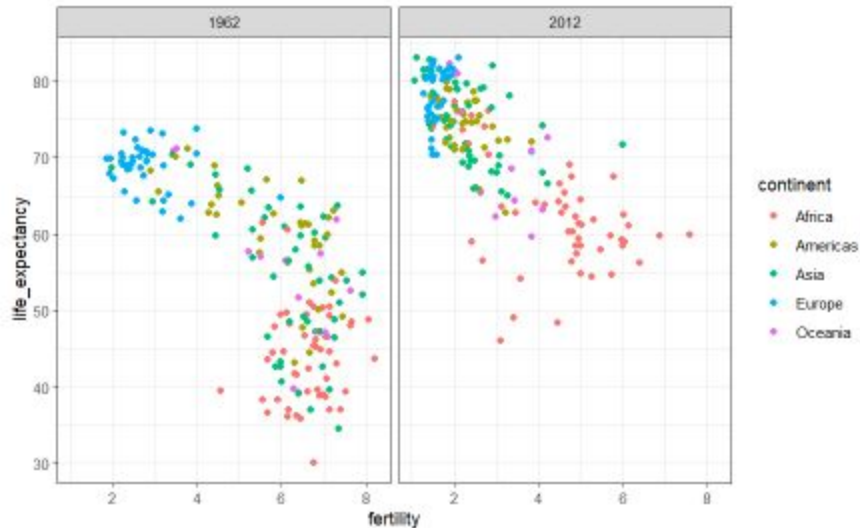


library(BrailleR)

```
## This chart has title 'Relationship between life expectancy and GDP per Capita'.  
## It has x-axis 'GDP per Capita [Dollar]' with labels 10000, 20000, 30000, 40000 and 50000.  
## It has y-axis 'Life Expectancy [Year]' with labels 72, 74, 76, 78, 80 and 82.  
## In this chart colour is used to show country. The legend that would normally indicate this has been hidden.  
## In this chart size is used to show pop. The legend that would normally indicate this has been hidden.  
## It has 2 layers.  
## Layer 1 is a set of 30 points.  
## Layer 1 has alpha set to 0.5.  
## Layer 2 is a textrepel graph that VI can not process.  
## Layer 2 has size set to 2.
```

facet_wrap

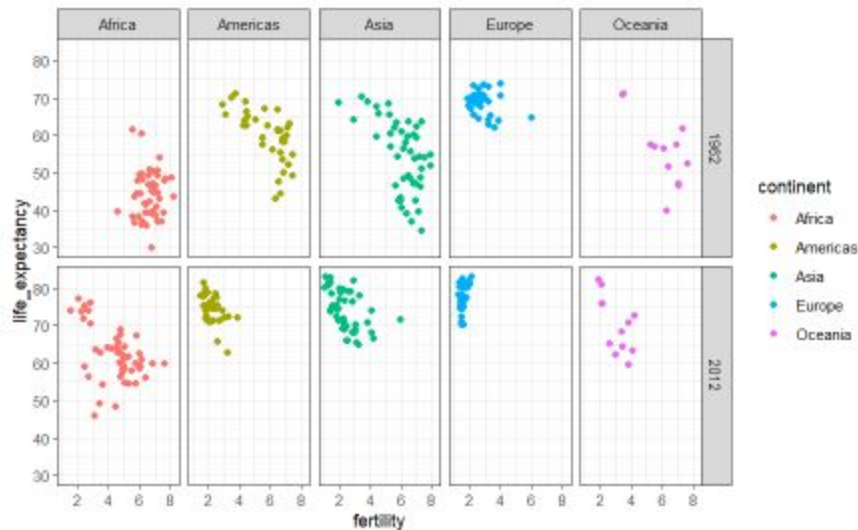
Faceting can group and visualize each subset. Facet_wrap can only visualize non-empty univariate.



```
```{r}
filter(gapminder, year%in%c(1962,
2012)) |> ggplot(aes(fertility,
life_expectancy, col = continent)) +
geom_point() +
facet_wrap(~year)
```
```

facet_grid

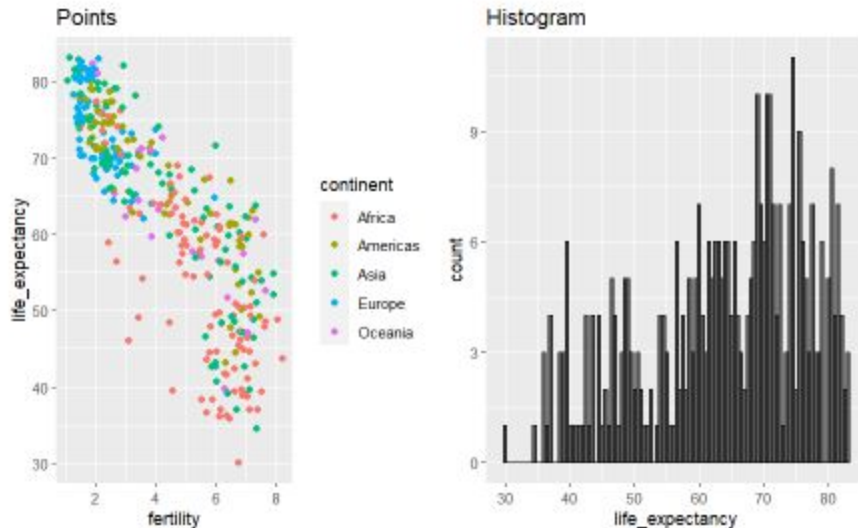
Facet_grid can visualize bivariate even some are empty.



```
```{r}
filter(gapminder, year%in%c(1962,
2012)) |> ggplot(aes(fertility,
life_expectancy, col = continent)) +
geom_point() +
facet_grid(year~continent)
```
```

grid.arrange

grid.arrange can layout multiple graphs on the same page



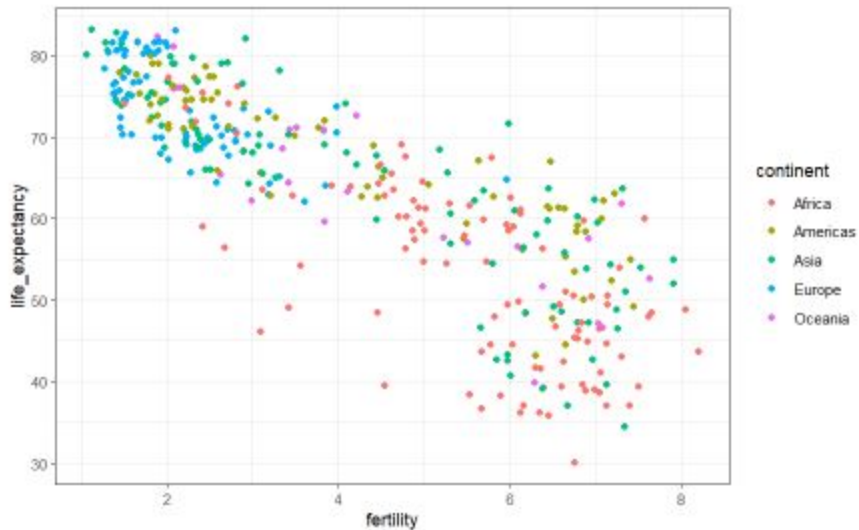
```
``{r}
```

```
p1=filter(gapminder, year%in%c(1962, 2012)) |> ggplot(aes(fertility, life_expectancy, col = continent)) +  
geom_point() + ggtitle("Points")
```

```
p2=filter(gapminder, year%in%c(1962, 2012)) |> ggplot(aes(life_expectancy)) +  
geom_histogram(binwidth = 0.5, color = "black") + ggtitle("Histogram")
```

```
grid.arrange(p1, p2, ncol = 2)``
```

ggplot2 can build a complex plot a layer at a time. Each layer can come from a different dataset and have a different aesthetic mapping. Layers are created using `geom_*` or `stat_*`



```
``{r}  
filter(gapminder, year%in%c(1962,  
2012)) |> ggplot(aes(fertility,  
life_expectancy, col = continent)) +  
geom_point()  
``
```

```
p + layer(  
  mapping = NULL,  
  data = NULL,  
  geom = "point",  
  stat = "identity",  
  position = "identity"  
)
```

Geom*

One variable:

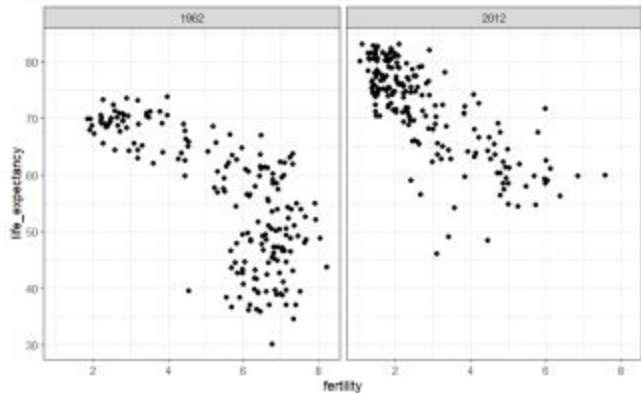
- `geom_bar()`: display distribution of discrete variable.
- `geom_histogram()`: bin and count continuous variable, display with bars.
- `geom_density()`: smoothed density estimate.
- `geom_dotplot()`: stack individual points into a dot plot.
- `geom_freqpoly()`: bin and count continuous variable, display with lines.

Two variables:

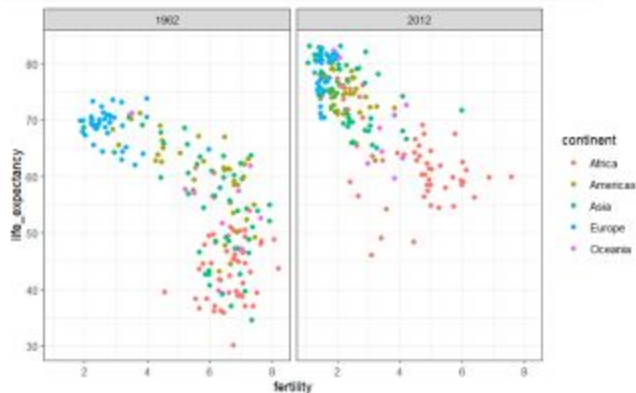
- `geom_point()`: scatterplot.
- `geom_quantile()`: smoothed quantile regression.
- `geom_rug()`: marginal rug plots.
- `geom_smooth()`: smoothed line of best fit.
- `geom_text()`: text labels.

Aesthetic Mapping*

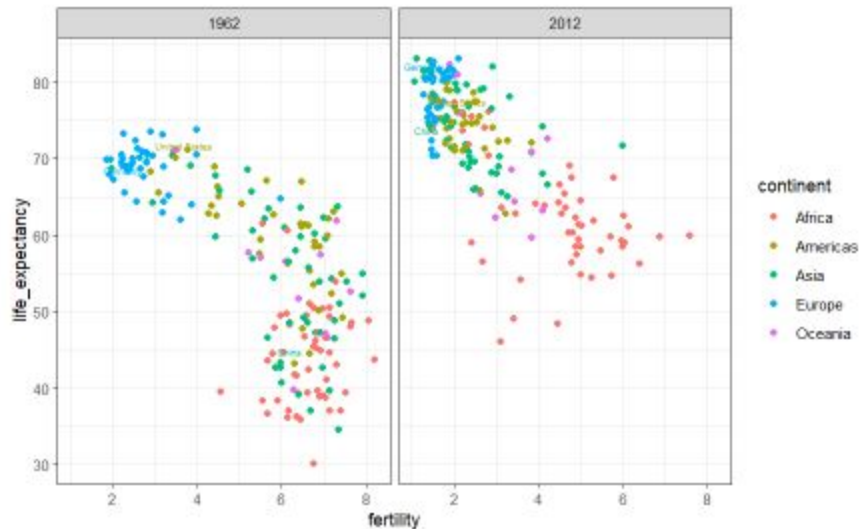
```
```{r}  
filter(gapminder, year==2012) |>
ggplot(aes(fertility, life_expectancy)) +
geom_point() +
Facet_wrap(. ~ year)
```
```



```
```{r}  
filter(gapminder, year==2012) |>
ggplot(aes(fertility, life_expectancy,
col=continent)) + geom_point() +
Facet_wrap(. ~ year)
```
```



Data



```
```{r}
highlight <- c("Germany", "China",
 "United States")
filter(gapminder, year %in% c(1962,
 2012)) %>% ggplot(aes(fertility,
 life_expectancy, color = continent)) +
 geom_point() +
 geom_text_repel(size = 2,
 show.legend = FALSE, aes(fertility,
 life_expectancy, label=country), data
 = filter(gapminder, year %in% c(1962,
 2012) & country %in% highlight)) +
 facet_grid(. ~ year)
```
```


Stats*

stat_bin(): geom_bar(), geom_freqpoly(), geom_histogram()

stat_bin2d(): geom_bin2d()

stat_bindot(): geom_dotplot()

stat_binhex(): geom_hex()

stat_boxplot(): geom_boxplot()

stat_contour(): geom_contour()

stat_quantile(): geom_quantile()

stat_smooth(): geom_smooth()

stat_sum(): geom_count()

Position*

Three position adjustments for points:

- `position_nudge()`: move points by a fixed offset.
- `position_jitter()`: add a little random noise to every position.
- `position_jitterdodge()`: dodge points within groups, then add a little random noise.

Three positions adjustments for bars:

- `position_stack()`: stack overlapping bars (or areas) on top of each other.
- `position_fill()`: stack overlapping bars, scaling so the top is always at 1.
- `position_dodge()`: place overlapping bars (or boxplots) side-by-side.

Saving and Exporting Visualizations

If using ggplot:

```
ggsave(  
  filename,           ← The name of the file locally  
  plot = last_plot(), ← Plot to export  
  device = NULL,      ← Type of export  
  path = NULL,        ← File location to be saved  
  scale = 1,          ← Multiplicative factor  
  width = NA,         ← Dimensions of graphic  
  height = NA,  
  units = c("in", "cm", "mm", "px"), ← unit of size  
  dpi = 300,          ← Set resolution  
  limitsize = TRUE,   ← max image size of 50 x 50 inches  
  bg = NULL,          ← Background color  
  ...  
)
```

ex. "mygraphic.png"
ex. plot1
ex. "Png", "jpeg"
ex. "mydocuments/folder1"

ex. 3,5,10 (inches, pixels...)

ex. "retina", 300, "screen"

ex. "gray"

Interactability

If looking for individuals to be able to change elements in your visualization, try creating an **RShiny App**

[Shiny App Example](#)

- Easy way to share ggplot figures, and others, in a way that allows users to change elements of the visualization easily
- See more information here: [R Shiny Tutorial](#)



Key Takeaways

- Don't take shortcuts in your code
- Be intentional with your visualizations
- Making good visualizations of data takes time
- Making accessible visualizations increases impact by broadening audience