

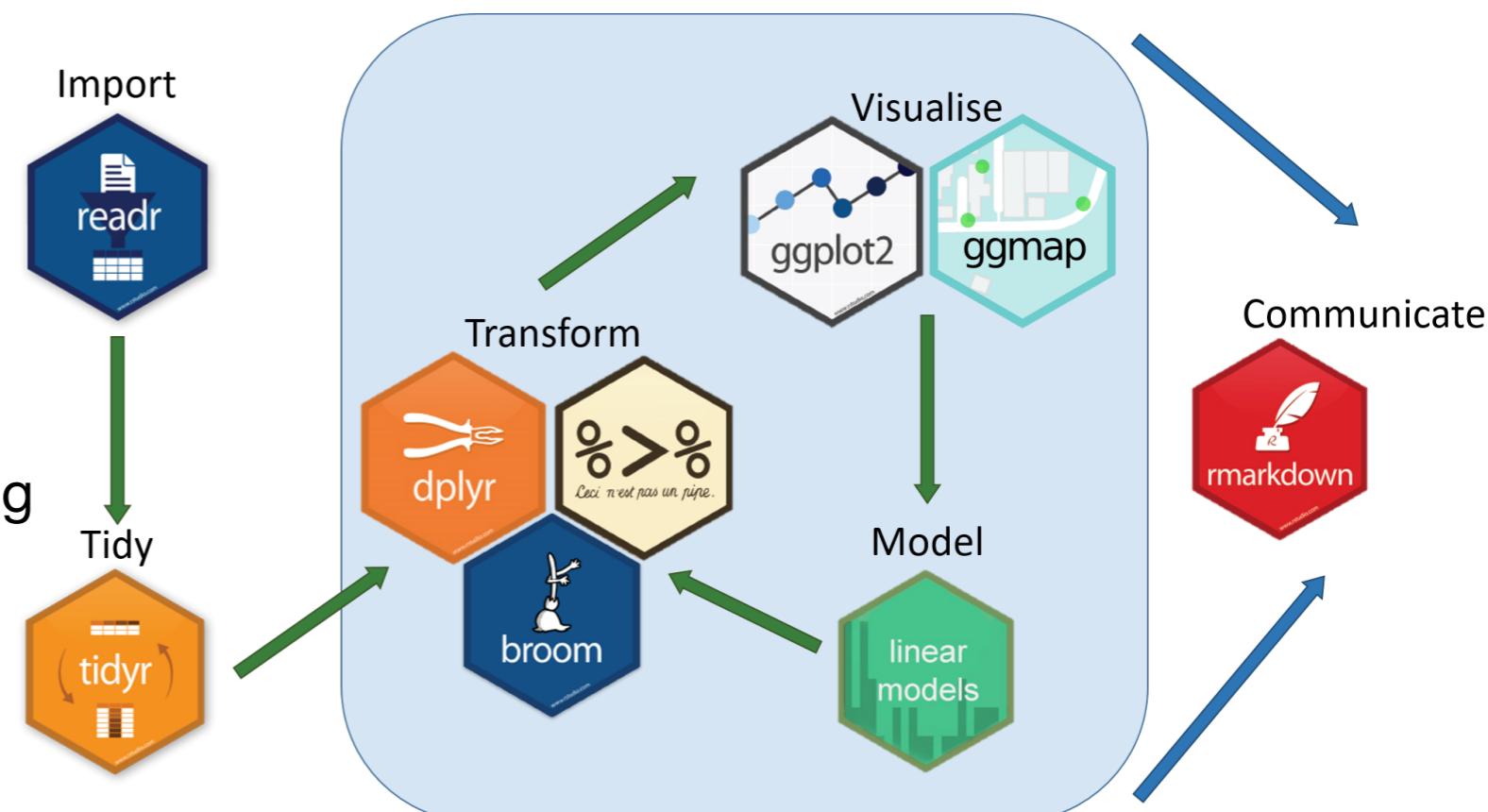
Data Science: GGplot2

Alex Di Genova

28/05/2024

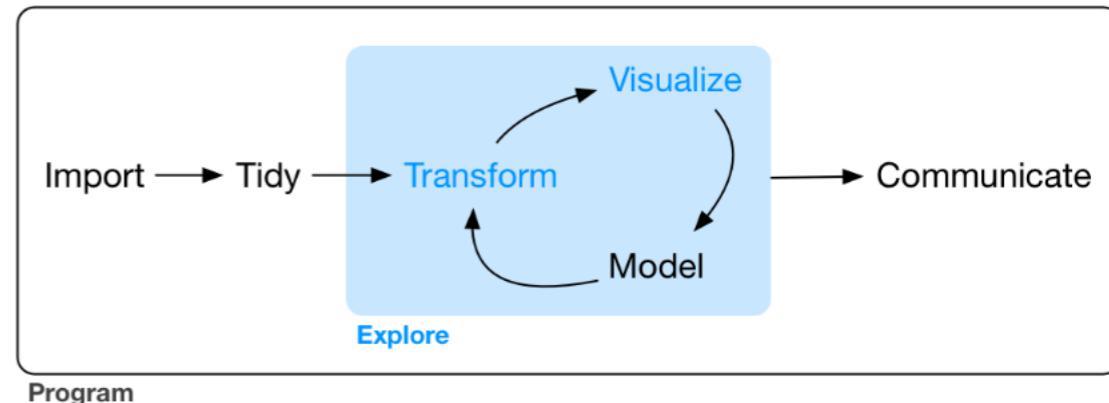
What is Tidyverse?

- A collection of R packages for data science
 - ggplot2 — data visualization
 - dplyr — data manipulation
 - tidyverse — data tidying
 - readr — data import
 - purrr — functional programming
 - tibble — modern dataframes
 - stringr — string manipulation
 - forcats — factor handling
- Data science workflow (import, clean, transform, visualize, model)



Tidyverse

Explore



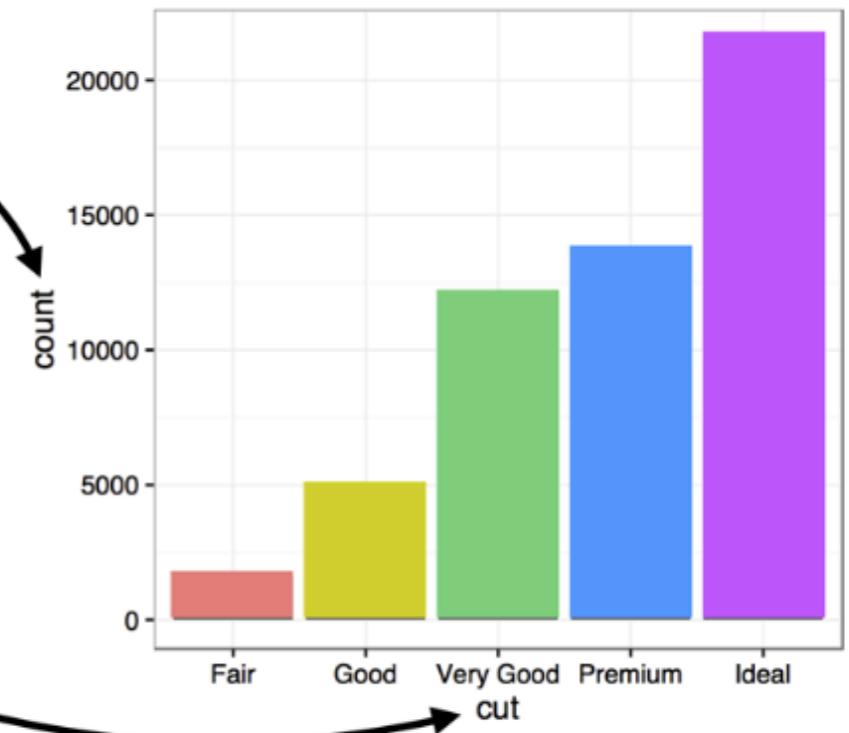
carat	cut	color	clarity	depth	table	price	x	y	z
0.23	Ideal	E	SI2	61.5	55	326	3.95	3.98	2.43
0.21	Premium	E	SI1	59.8	61	326	3.89	3.84	2.31
0.23	Good	E	VS1	56.9	65	327	4.05	4.07	2.31
0.29	Premium	I	VS2	62.4	58	334	4.20	4.23	2.63
0.31	Good	J	SI2	63.3	58	335	4.34	4.35	2.75
...

stat_count()

cut	count	prop
Fair	1610	1
Good	4906	1
Very Good	12082	1
Premium	13791	1
Ideal	21551	1

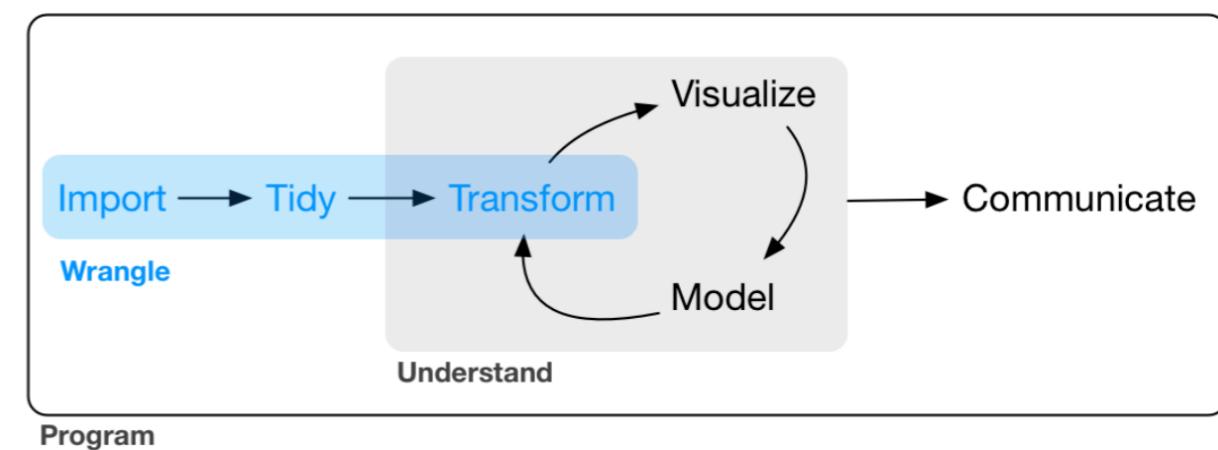
5. Place geoms in a cartesian coordinate system.

6. Map the y values to `..count..` and the x values to `cut`.



Tidyverse

Tidy



religion	<\$10k	\$10–20k	\$20–30k	\$30–40k	\$40–50k	\$50–75k
Agnostic	27	34	60	81	76	137
Atheist	12	27	37	52	35	70
Buddhist	27	21	30	34	33	58
Catholic	418	617	732	670	638	1116
Don't know/refused	15	14	15	11	10	35
Evangelical Prot	575	869	1064	982	881	1486
Hindu	1	9	7	9	11	34
Historically Black Prot	228	244	236	238	197	223
Jehovah's Witness	20	27	24	24	21	30
Jewish	19	19	25	25	30	95

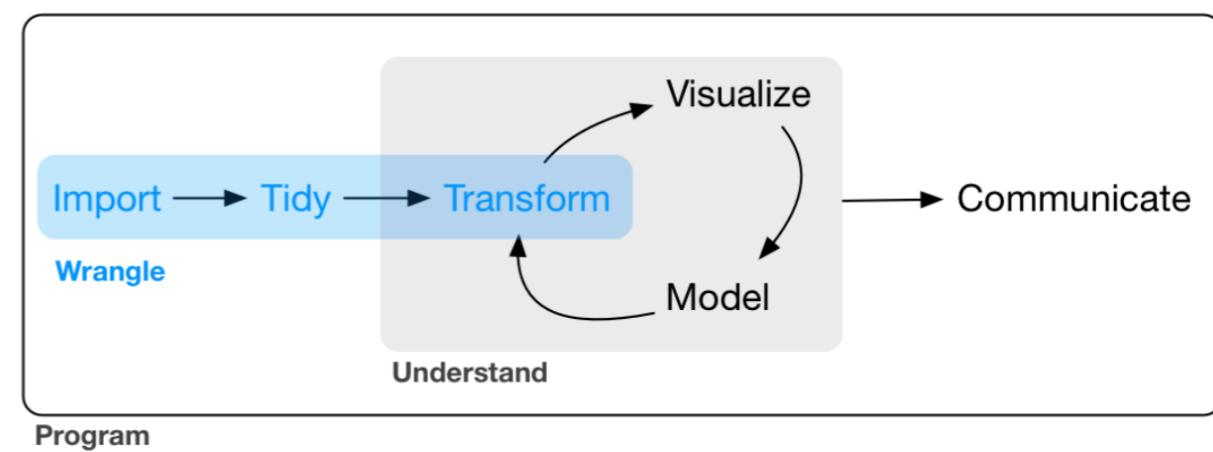
religion	income	freq
Agnostic	<\$10k	27
Agnostic	\$10–20k	34
Agnostic	\$20–30k	60
Agnostic	\$30–40k	81
Agnostic	\$40–50k	76
Agnostic	\$50–75k	137
Agnostic	\$75–100k	122
Agnostic	\$100–150k	109
Agnostic	>150k	84
Agnostic	Don't know/refused	96

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax	—	—	—	—	—	—	—	—
MX17004	2010	1	tmin	—	—	—	—	—	—	—	—
MX17004	2010	2	tmax	—	27.3	24.1	—	—	—	—	—
MX17004	2010	2	tmin	—	14.4	14.4	—	—	—	—	—
MX17004	2010	3	tmax	—	—	—	—	32.1	—	—	—
MX17004	2010	3	tmin	—	—	—	—	14.2	—	—	—
MX17004	2010	4	tmax	—	—	—	—	—	—	—	—
MX17004	2010	4	tmin	—	—	—	—	—	—	—	—
MX17004	2010	5	tmax	—	—	—	—	—	—	—	—
MX17004	2010	5	tmin	—	—	—	—	—	—	—	—

id	date	tmax	tmin
MX17004	2010-01-30	27.8	14.5
MX17004	2010-02-02	27.3	14.4
MX17004	2010-02-03	24.1	14.4
MX17004	2010-02-11	29.7	13.4
MX17004	2010-02-23	29.9	10.7
MX17004	2010-03-05	32.1	14.2
MX17004	2010-03-10	34.5	16.8
MX17004	2010-03-16	31.1	17.6
MX17004	2010-04-27	36.3	16.7
MX17004	2010-05-27	33.2	18.2

Tidyverse

Tidy



country	year	cases	population
Afghanistan	2019	745	1537071
Afghanistan	2000	1666	2059360
Brazil	1999	31737	172106362
Brazil	2000	81488	174104898
China	1999	213258	1272115272
China	2010	21166	1280128583

variables

country	year	cases	population
Afghanistan	2019	745	1537071
Afghanistan	2000	1666	2059360
Brazil	1999	31737	172106362
Brazil	2000	81488	174104898
China	1999	213258	1272115272
China	2010	21166	1280128583

observations

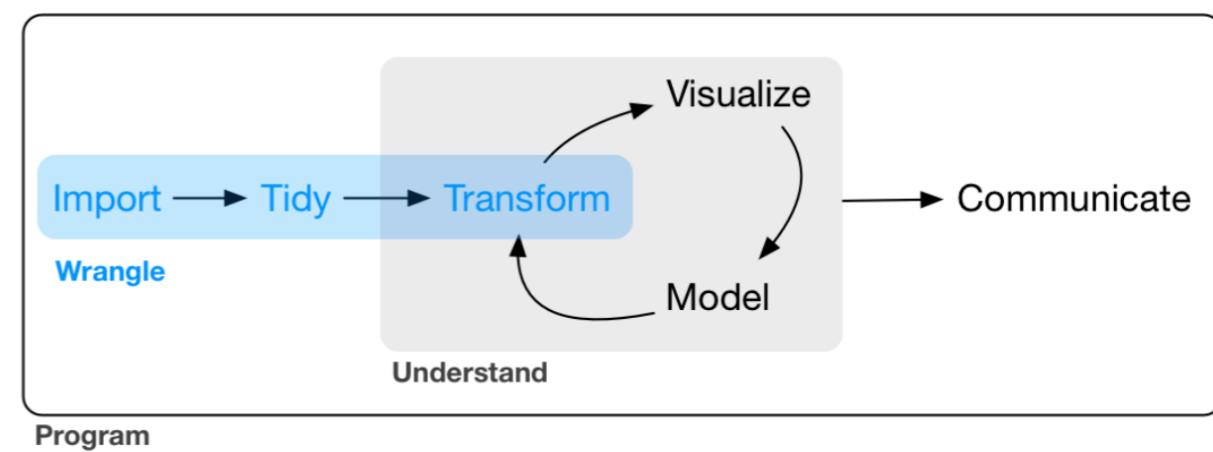
country	year	cases	population
Afghanistan	09	745	1537071
Afghanistan	00	1666	2059360
Brazil	99	31737	172106362
Brazil	00	81488	174104898
China	99	213258	1272115272
China	00	21166	1280128583

values

- Tidy dataset:
 - Variables are in columns
 - Observations in rows
 - Values in cells

Tidyverse

Tidy



- Spreading and Gathering
- Dateset where columns are values of a variable

```
```{r tidy1}
```

```
table4a
```

```
```
```

A tibble: 3 × 3

| country | 1999 | 2000 |
|-------------|--------|-------|
| <chr> | <dbl> | <dbl> |
| Afghanistan | 745 | 2 |
| Brazil | 37737 | 80 |
| China | 212258 | 213 |

3 rows

```
```{r tidy1}
```

```
table4a %>% gather(`1999`, `2000`, key="year", value="population")
```

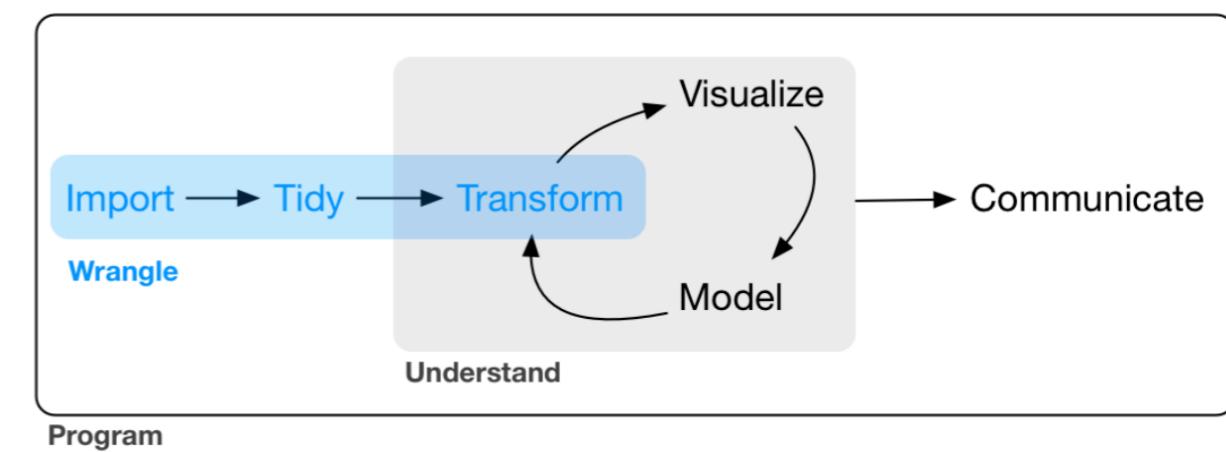
```
```
```

A tibble: 6 × 3

| country | year | population |
|-------------|-------|------------|
| <chr> | <chr> | <dbl> |
| Afghanistan | 1999 | 745 |
| Brazil | 1999 | 37737 |
| China | 1999 | 212258 |
| Afghanistan | 2000 | 2666 |
| Brazil | 2000 | 80488 |
| China | 2000 | 213766 |

Tidyverse

Tidy



- Spreading
 - Opposite of gathering

| country | year | key | value | country |
|-------------|------|------------|------------|-------------|
| Afghanistan | 1999 | cases | 745 | Afghanistan |
| Afghanistan | 1999 | population | 19987071 | Afghanistan |
| Afghanistan | 2000 | cases | 2666 | Brazil |
| Afghanistan | 2000 | population | 20595360 | Brazil |
| Brazil | 1999 | cases | 37737 | China |
| Brazil | 1999 | population | 172006362 | China |
| Brazil | 2000 | cases | 80488 | China |
| Brazil | 2000 | population | 174504898 | China |
| China | 1999 | cases | 212258 | China |
| China | 1999 | population | 1272915272 | China |
| China | 2000 | cases | 213766 | China |
| China | 2000 | population | 1280428583 | |

table2

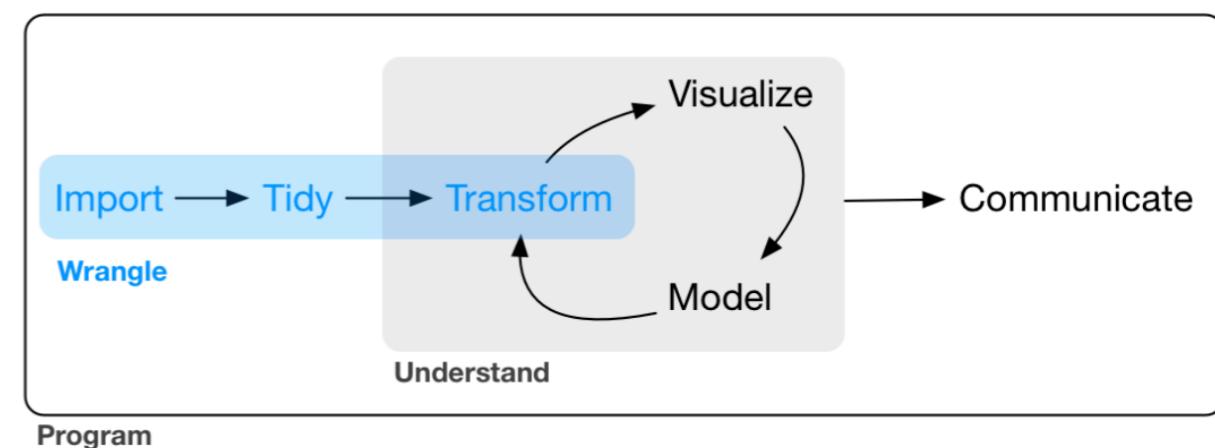
```
```{r tidy2}
table2 %>% spread(key=type, value=count)
````
```

A tibble: 6 × 4

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 745 | 19987071 |
| Afghanistan | 2000 | 2666 | 20595360 |
| Brazil | 1999 | 37737 | 172006362 |
| Brazil | 2000 | 80488 | 174504898 |
| China | 1999 | 212258 | 1272915272 |
| China | 2000 | 213766 | 1280428583 |

Tidyverse

Tidy



- Separate() and unite()

| country | year | rate |
|-------------|------|---------------------|
| Afghanistan | 1999 | 745 / 19987071 |
| Afghanistan | 2000 | 2666 / 20595360 |
| Brazil | 1999 | 37737 / 172006362 |
| Brazil | 2000 | 80488 / 174504898 |
| China | 1999 | 212258 / 1272915272 |
| China | 2000 | 213766 / 1280428583 |

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 745 | 19987071 |
| Afghanistan | 2000 | 2666 | 20595360 |
| Brazil | 1999 | 37737 | 172006362 |
| Brazil | 2000 | 80488 | 174504898 |
| China | 1999 | 212258 | 1272915272 |
| China | 2000 | 213766 | 1280428583 |

table3

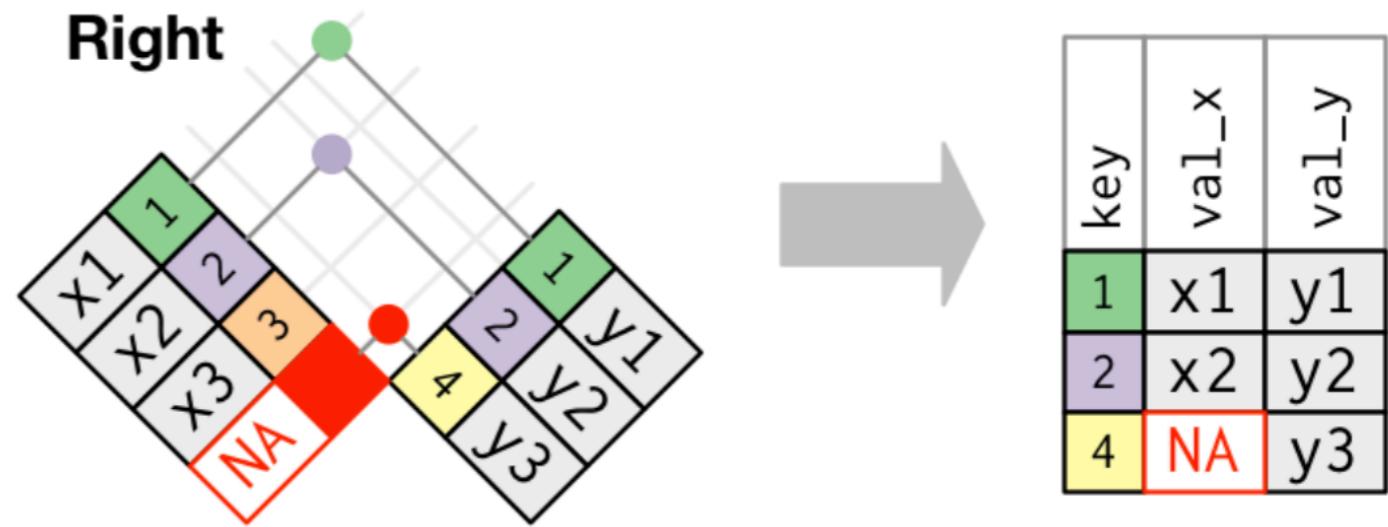
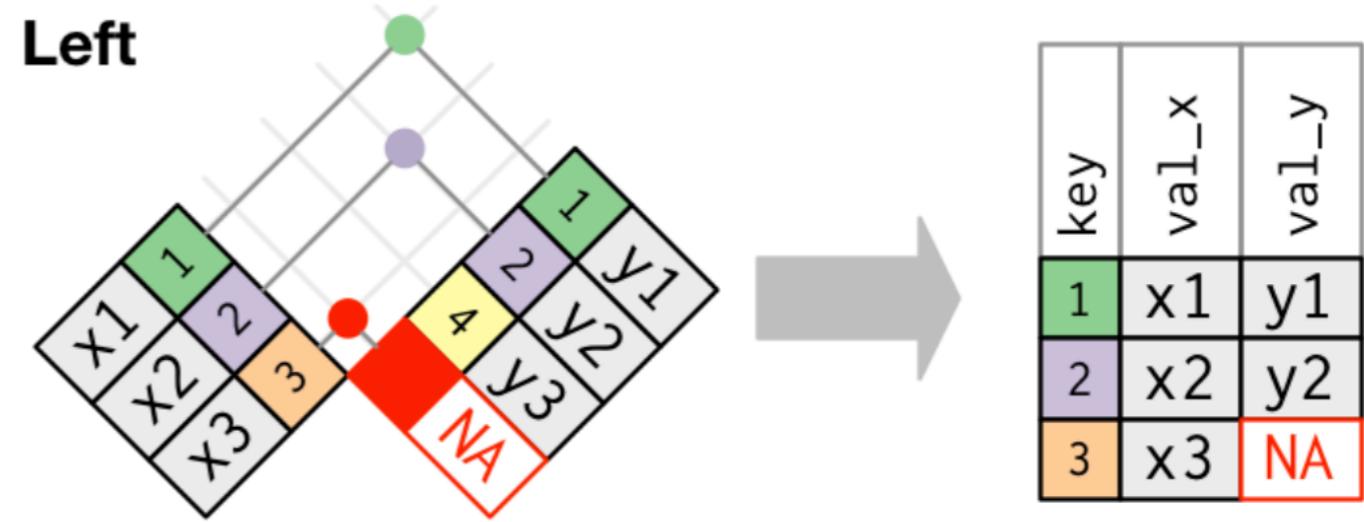
| country | year | rate |
|-------------|------|---------------------|
| Afghanistan | 1999 | 745 / 19987071 |
| Afghanistan | 2000 | 2666 / 20595360 |
| Brazil | 1999 | 37737 / 172006362 |
| Brazil | 2000 | 80488 / 174504898 |
| China | 1999 | 212258 / 1272915272 |
| China | 2000 | 213766 / 1280428583 |

| country | century | year | rate |
|-------------|---------|------|---------------------|
| Afghanistan | 19 | 99 | 745 / 19987071 |
| Afghanistan | 20 | 0 | 2666 / 20595360 |
| Brazil | 19 | 99 | 37737 / 172006362 |
| Brazil | 20 | 0 | 80488 / 174504898 |
| China | 19 | 99 | 212258 / 1272915272 |
| China | 20 | 0 | 213766 / 1280428583 |

table6

Tidyverse Joins

- Merging two data frames by a key.
 - `left_join(x,y, by= "key")`
 - `left_join(x,y,by=c("key1" = "key2"))`



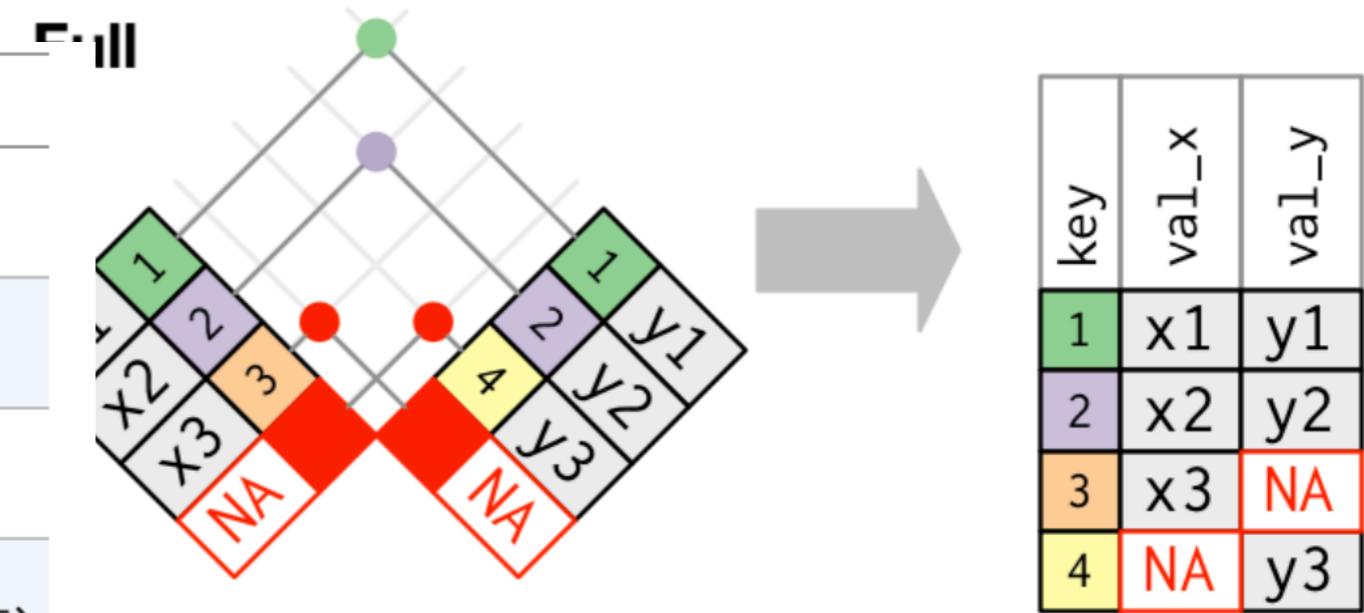
dplyr merge

```
inner_join(x, y)  merge(x, y)

left_join(x, y)   merge(x, y, all.x = TRUE)

right_join(x, y)  merge(x, y, all.y = TRUE),

full_join(x, y)   merge(x, y, all.x = TRUE, all.y = TRUE)
```

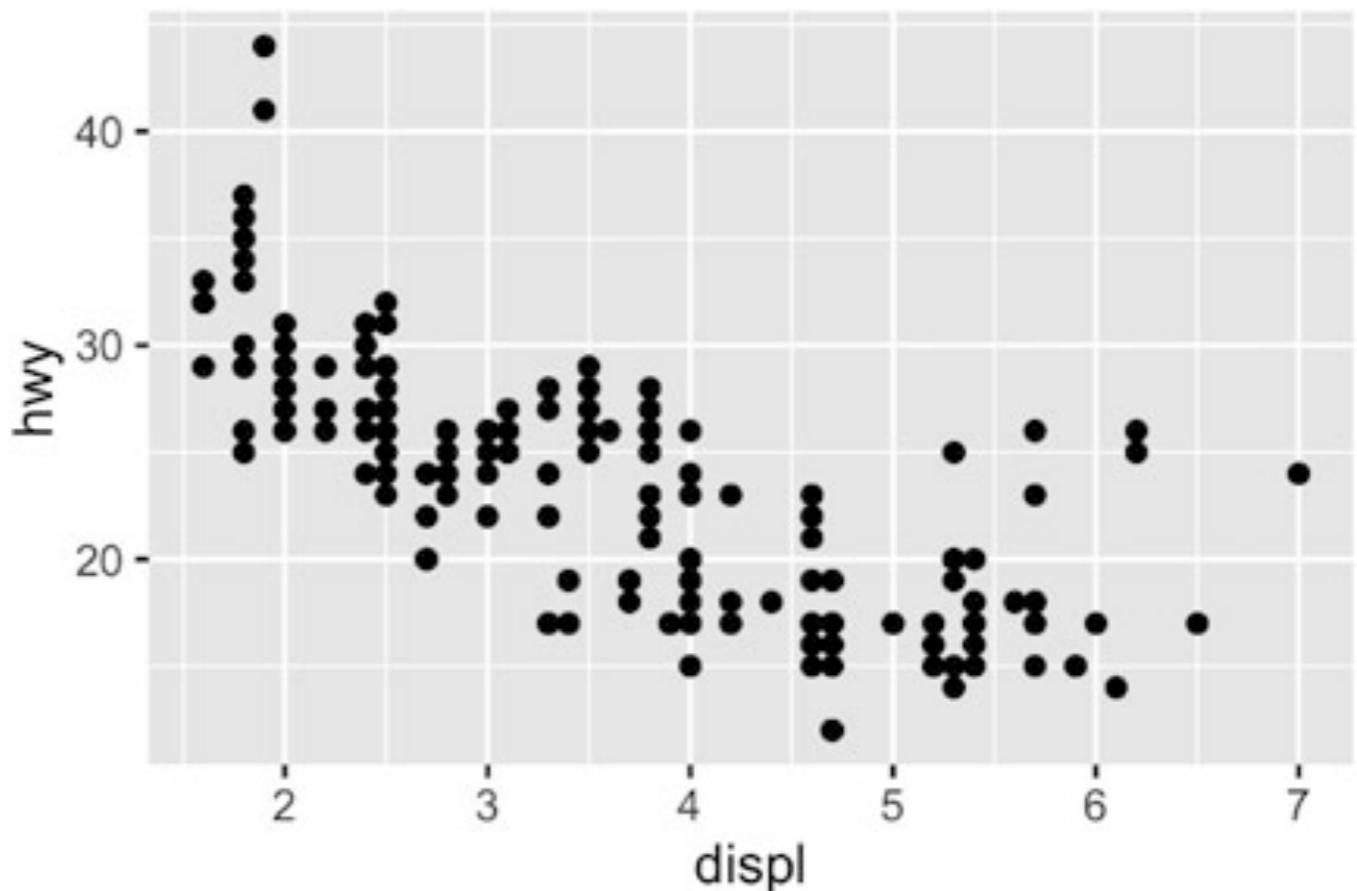


Tidyverse

Ggplot2

- Key Components
 - Data
 - aesthetic mappings -> variables in the data and visual properties
 - At least one layer that describes how to render each observation (geom_)

```
ggplot(mpg, aes(x = displ, y = hwy)) + geom_point()
```



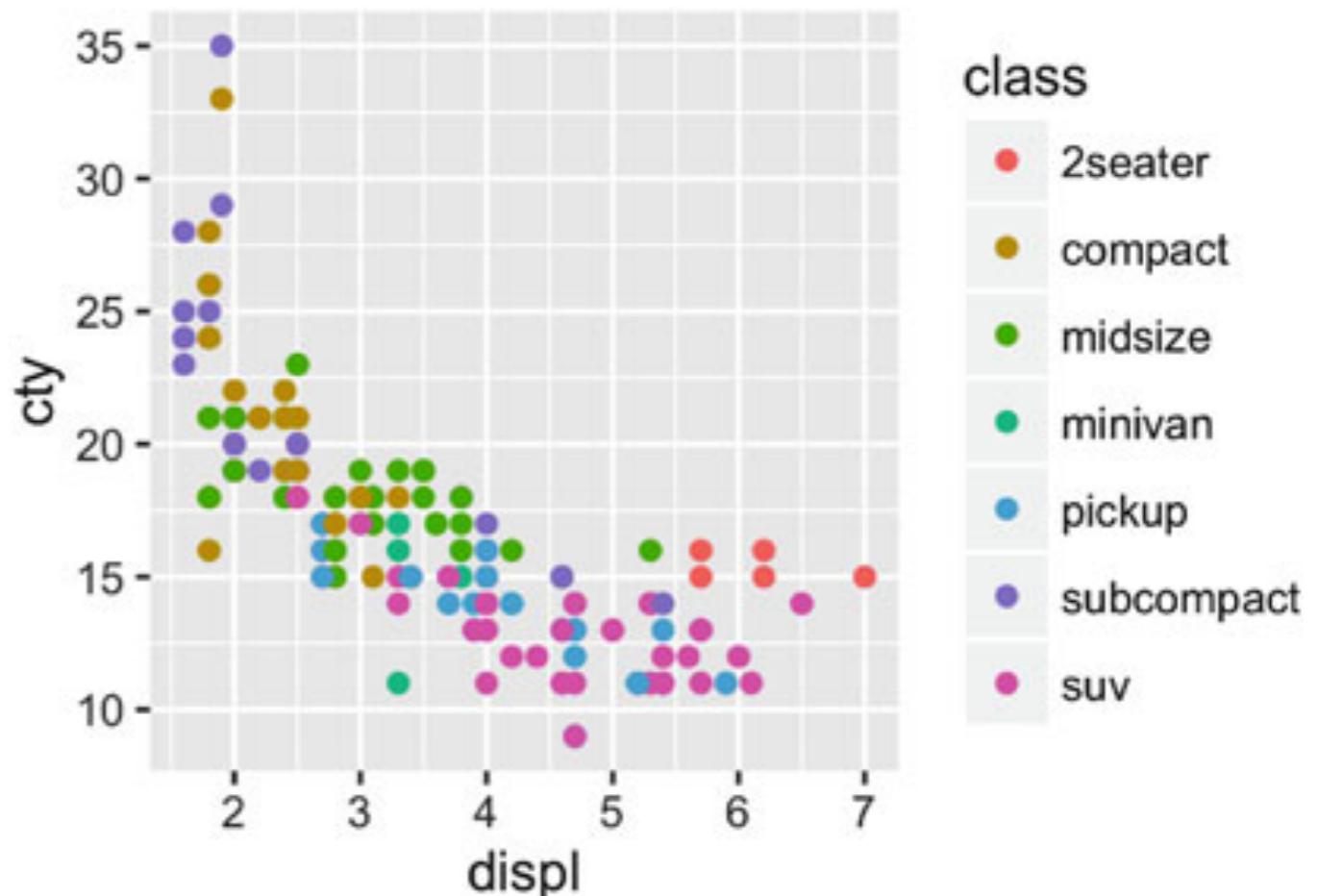
Tidyverse

Ggplot2

- Key Components

- Data
- aesthetic mappings -> variables in the data and visual properties
 - Colour, Size, Shape and Other Aesthetic Attributes
 - aes(displ, hwy, colour = class)
 - aes(displ, hwy, shape = drv)
 - aes(displ, hwy, size = cyl)

```
ggplot(mpg, aes(x = displ, y = hwy, colour=class)) +  
  geom_point()
```

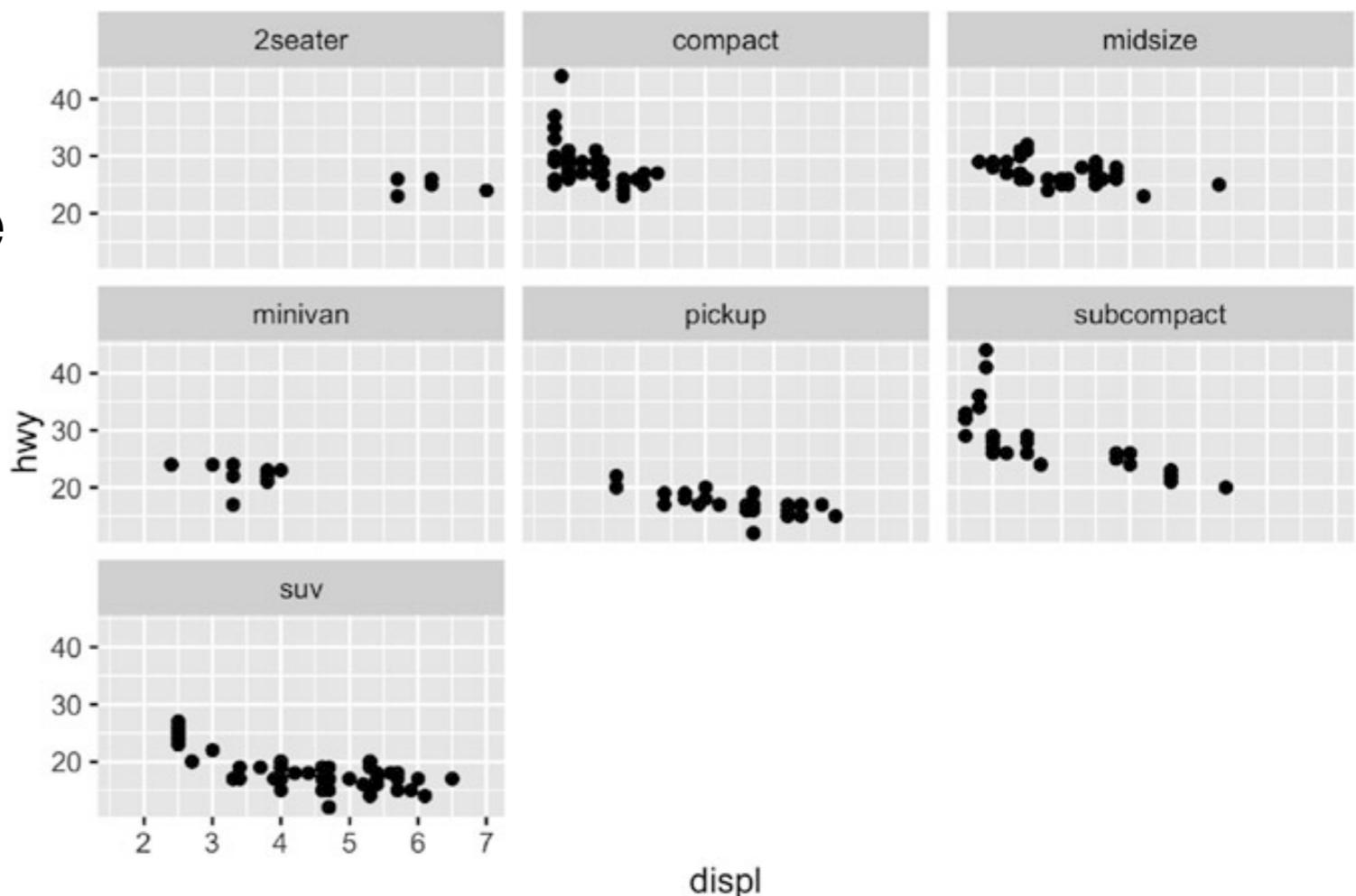


GGplot2

Facets

- Faceting
 - creates tables of graphics by splitting the data into subsets and displaying the same graph for each subset
 - Grid and wrapped
 - `facet_wrap()`

```
ggplot(mpg, aes(displ, hwy)) +  
  geom_point() +  
  facet_wrap(~class)
```

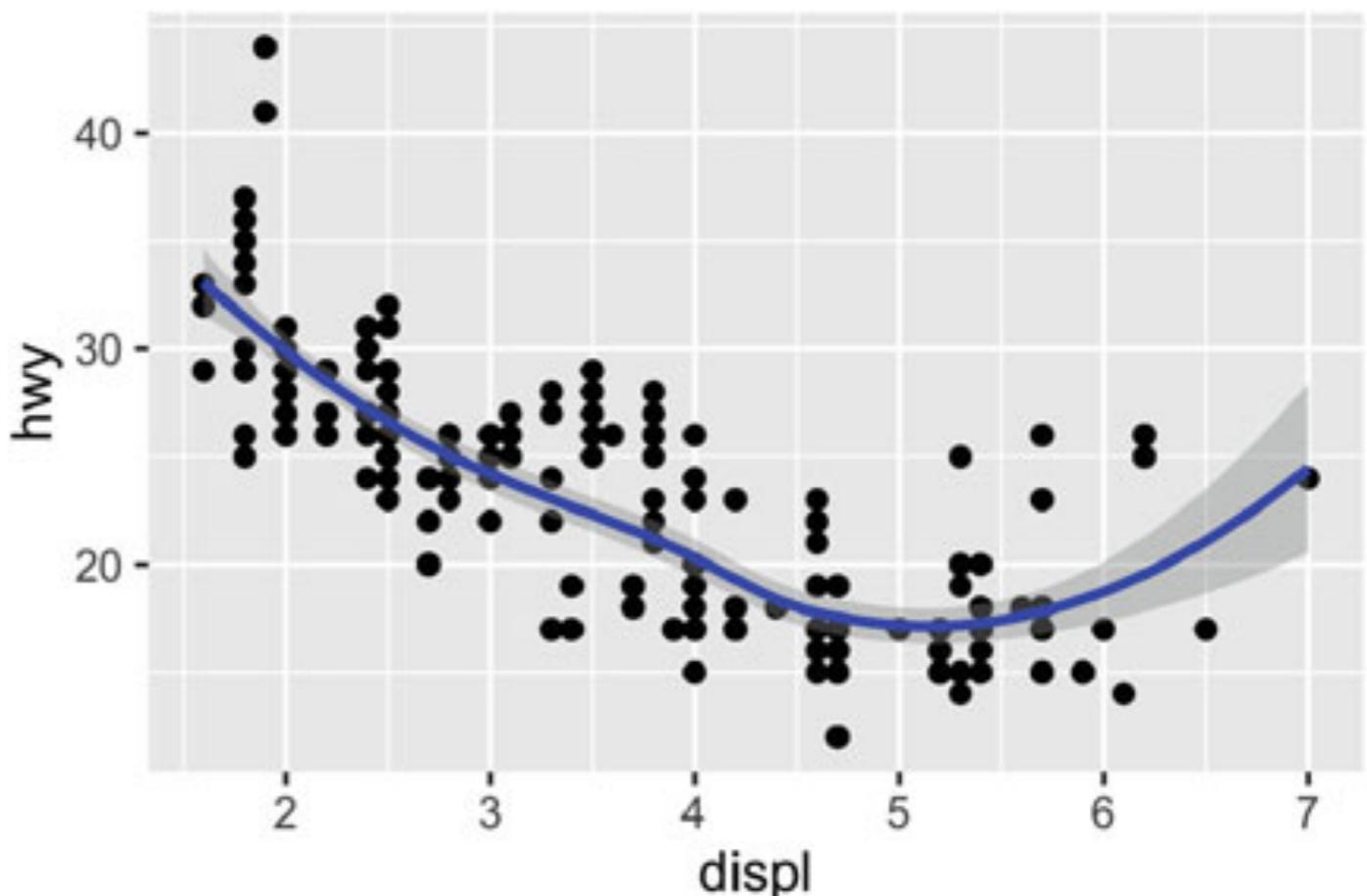


GGplot2

Geoms

- **geom_smooth()** fits a smoother to the data and displays the smooth and its standard error.
- **geom_boxplot()** produces a box-and-whisker plot to summarise the distribution of a set of points.
- **geom_histogram()** and **geom_freqpoly()** show the distribution of continuous variables.
- **geom_bar()** shows the distribution of categorical variables.
- **geom_path()** and **geom_line()** draw lines between the data points.

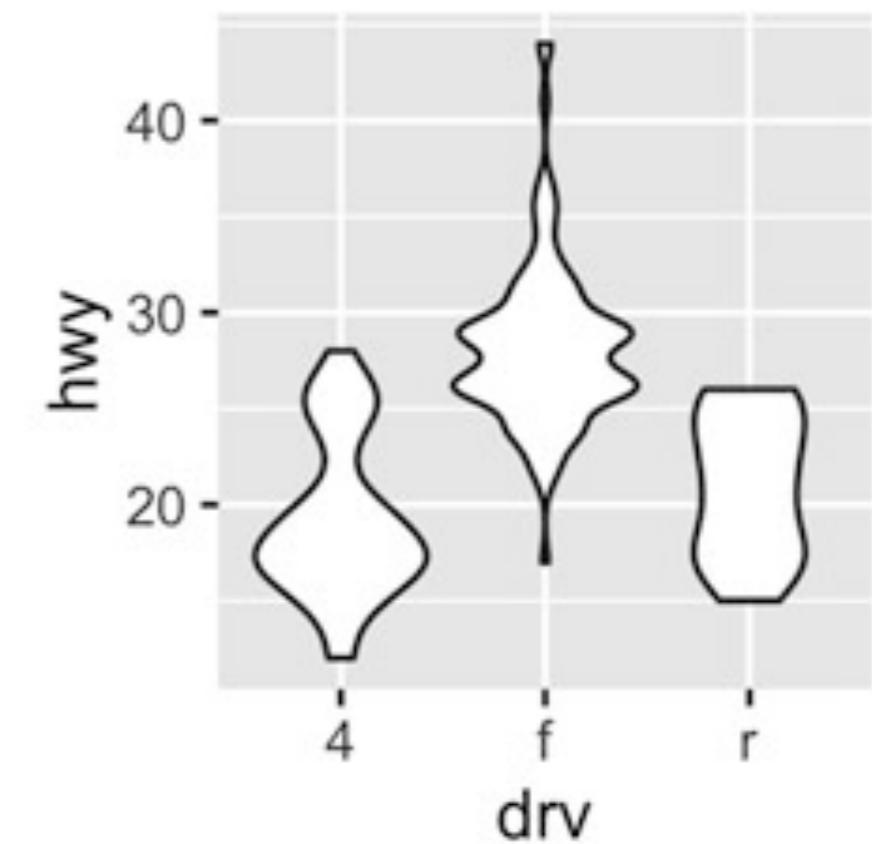
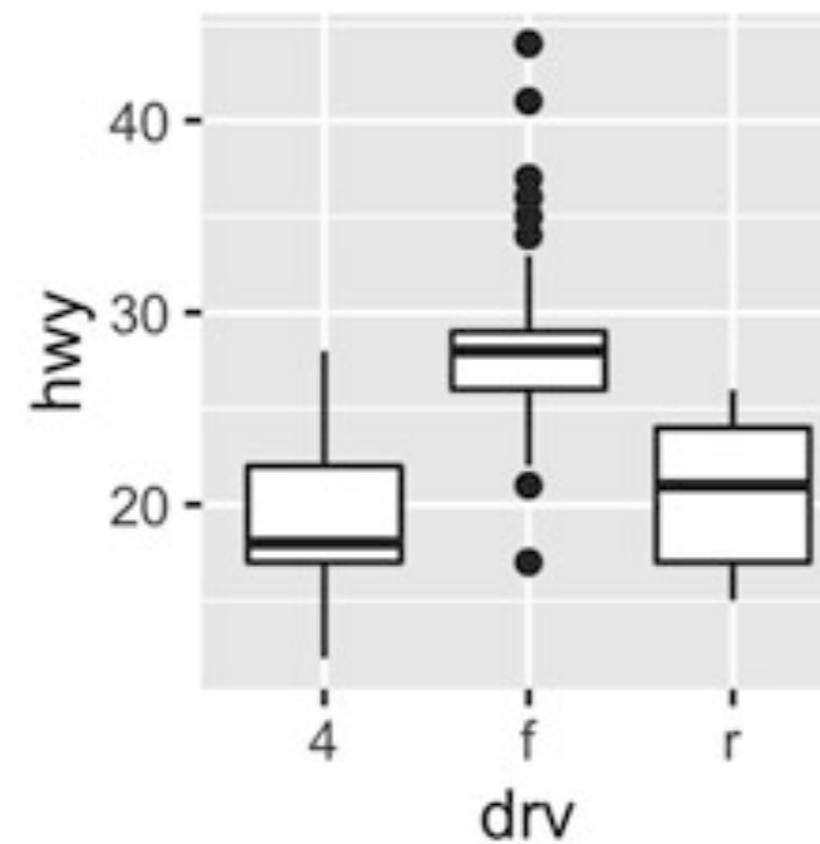
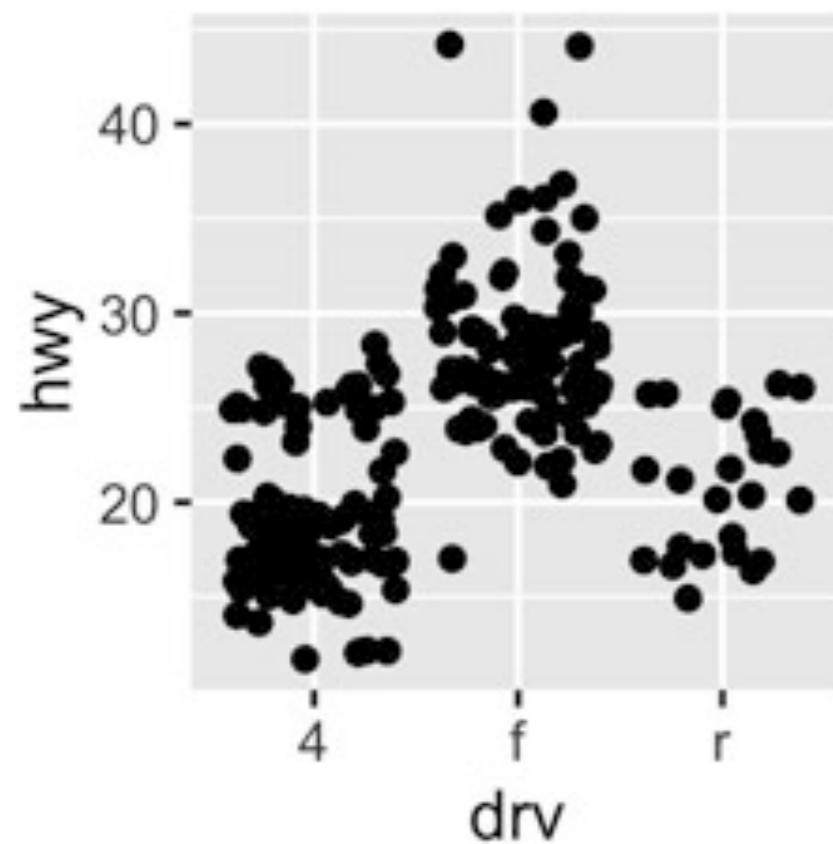
```
ggplot(mpg, aes(displ, hwy)) +  
  geom_point() +  
  geom_smooth()
```



GGplot2

Geoms

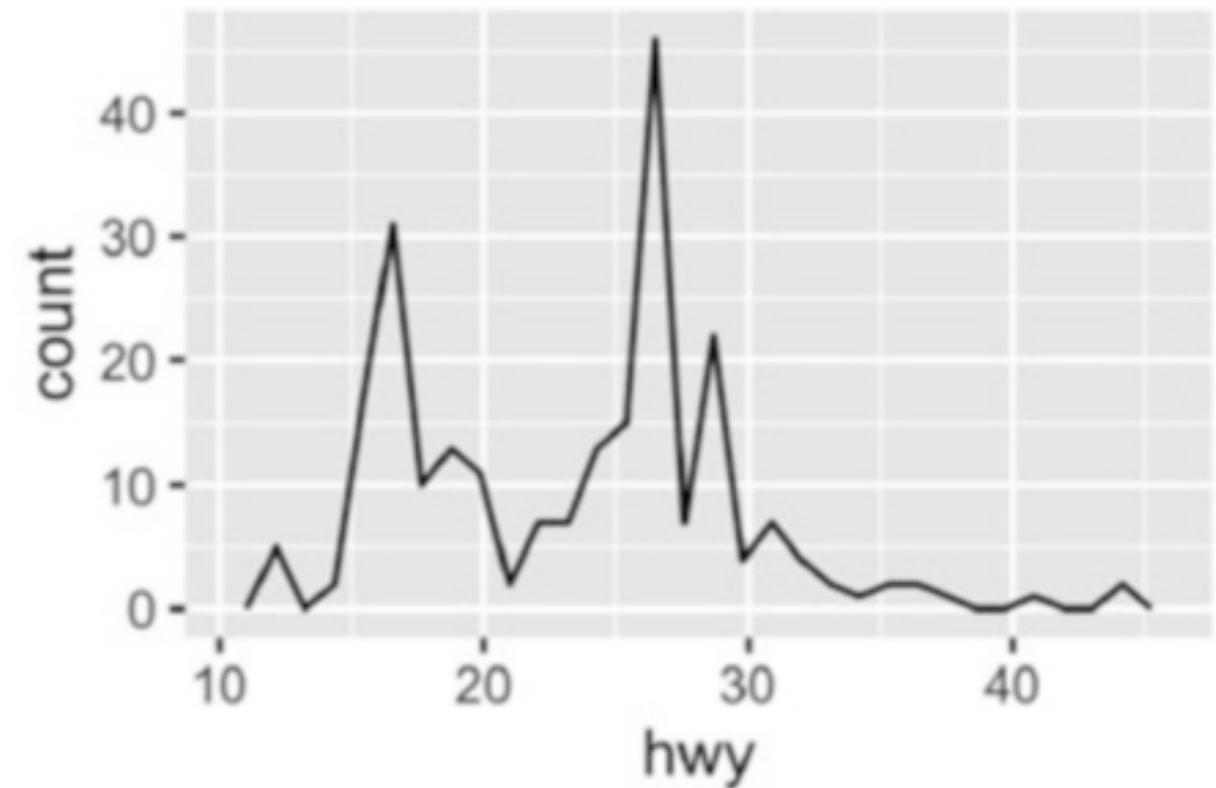
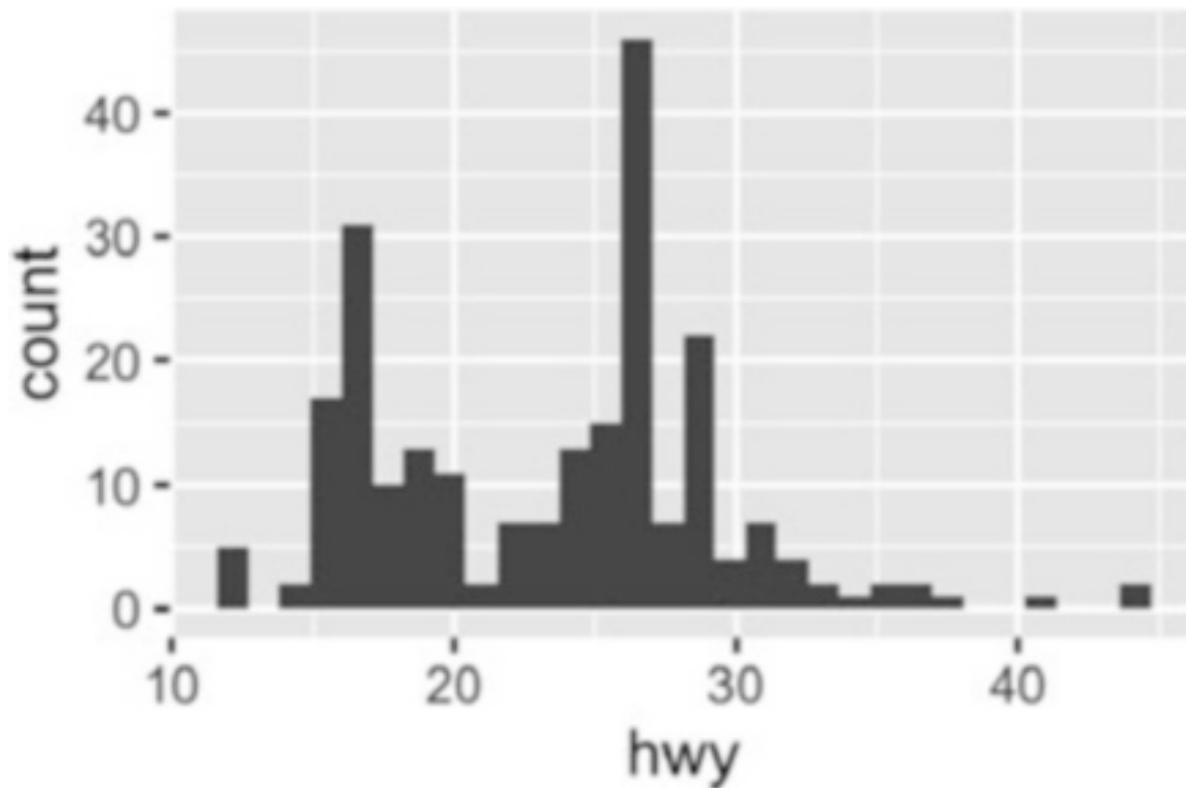
- `ggplot(mpg, aes(drv, hwy)) + geom_jitter()`
- `ggplot(mpg, aes(drv, hwy)) + geom_boxplot()`
- `ggplot(mpg, aes(drv, hwy)) + geom_violin()`



GGplot2

Geoms

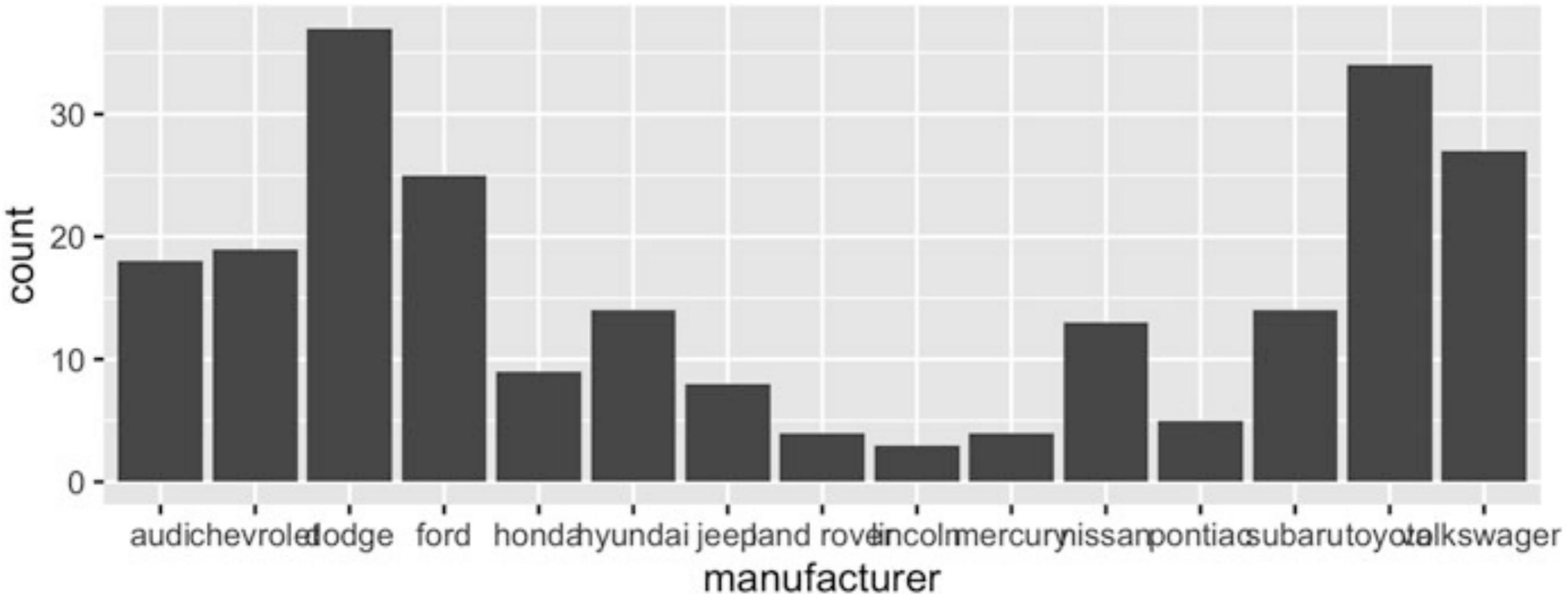
- `ggplot(mpg, aes(hwy)) + geom_histogram()`
- `aggplot(mpg, aes(hwy)) + geom_freqpoly()`



GGplot2

Geoms

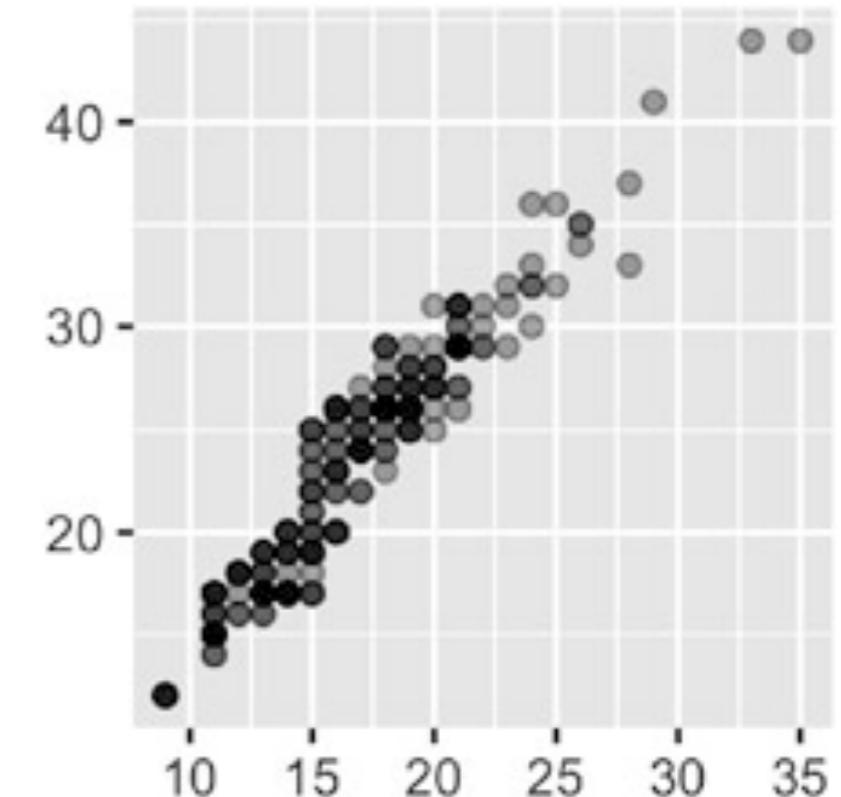
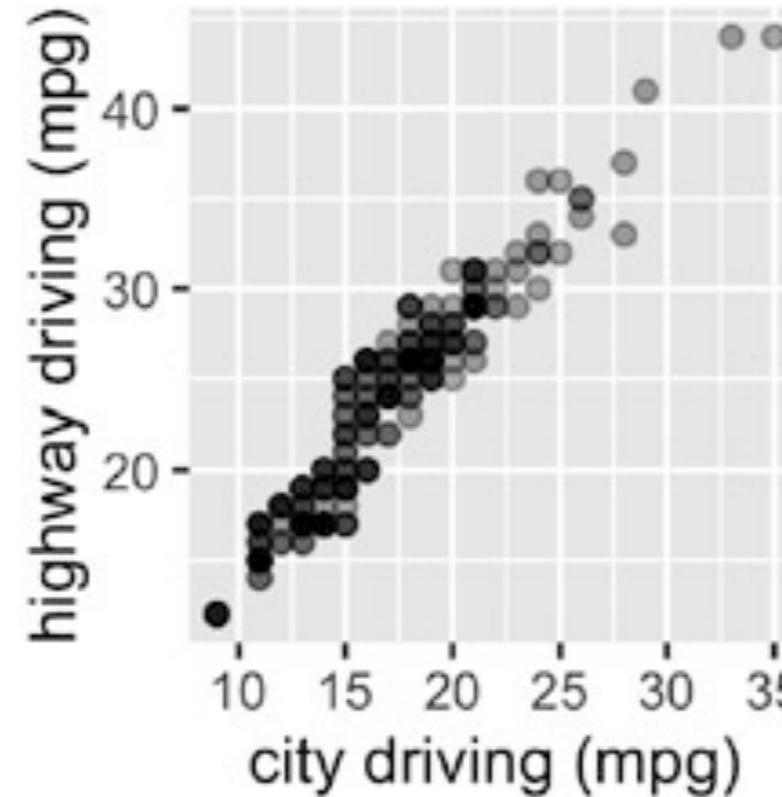
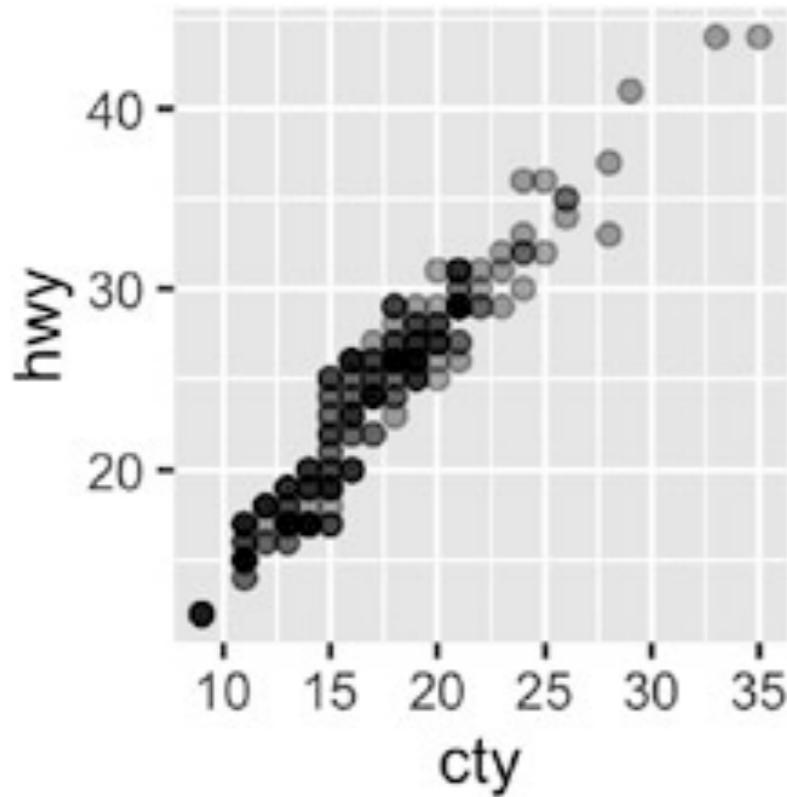
```
ggplot(mpg, aes(manufacturer)) +  
  geom_bar()
```



GGplot2

Geoms

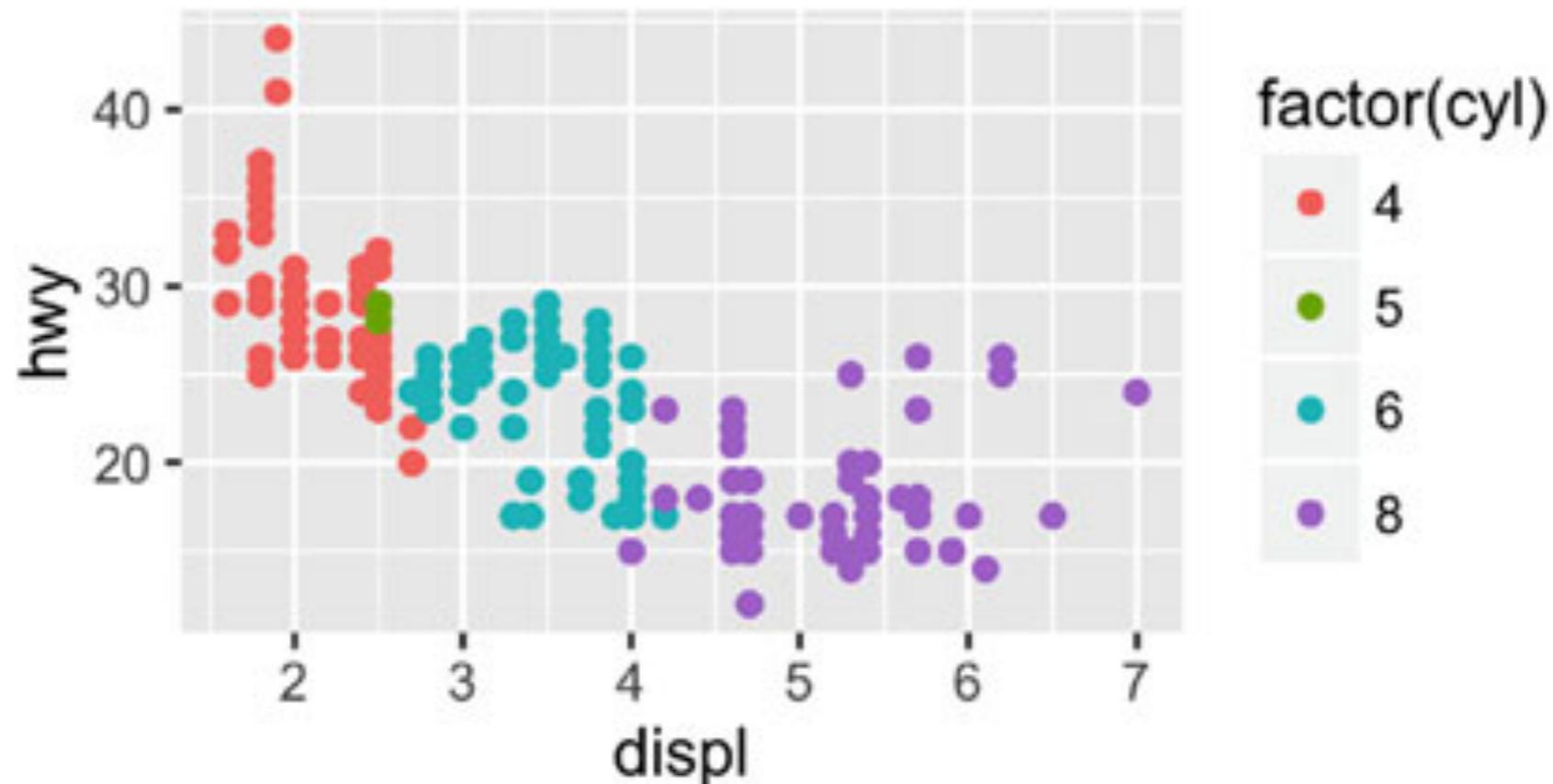
- `ggplot(mpg, aes(cty, hwy)) + geom_point(alpha = 1 / 3)`
- `ggplot(mpg, aes(cty, hwy)) + geom_point(alpha = 1 / 3) + xlab("city driving (mpg)") + ylab("highway driving (mpg)")`
- `ggplot(mpg, aes(cty, hwy)) + geom_point(alpha = 1 / 3) + xlab(NULL) + ylab(NULL)`



GGplot2

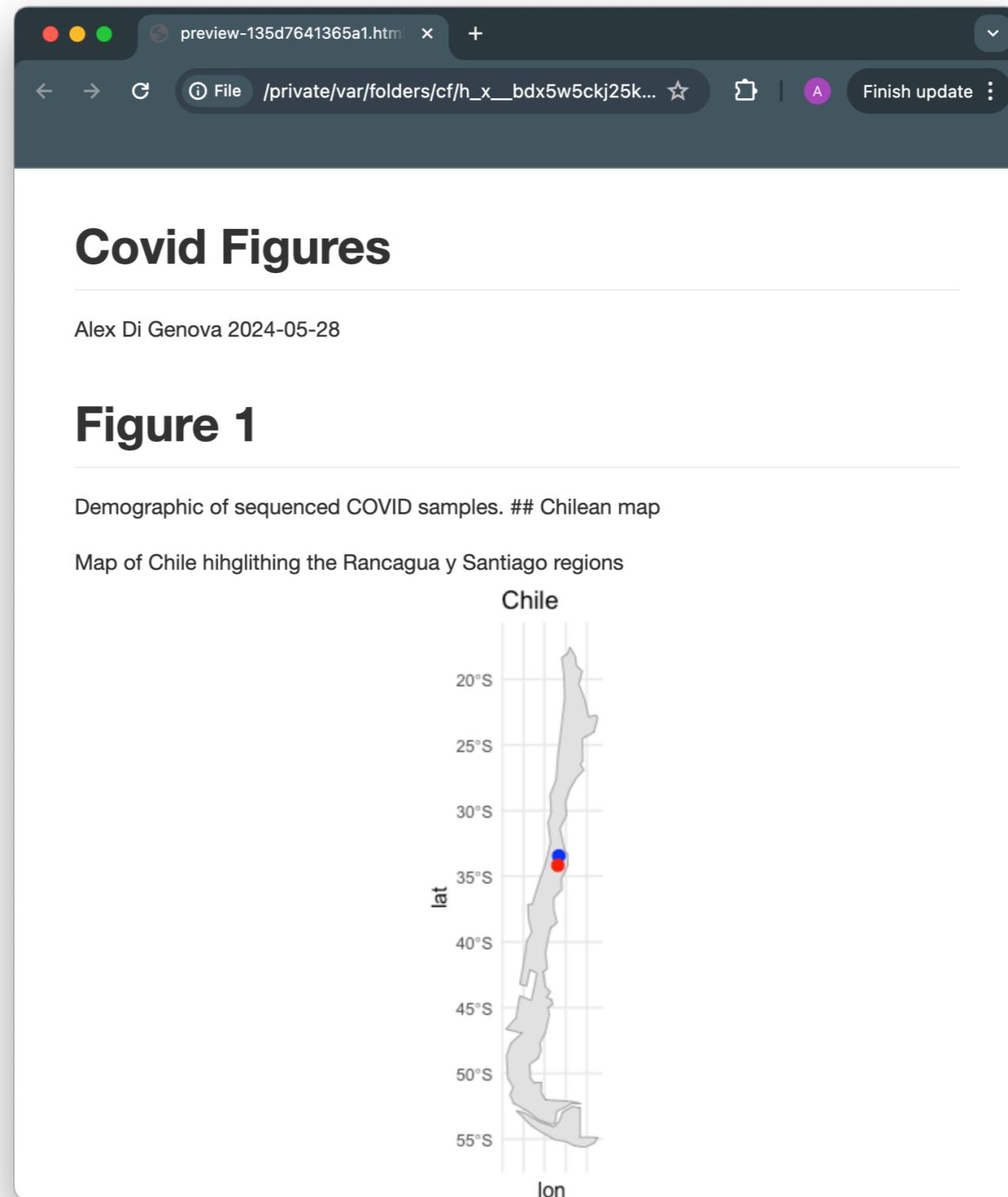
Output

```
p <- ggplot(mpg, aes(displ, hwy, colour = factor(cyl))) + geom_point()  
print(p)  
  
ggsave("plot.png", width = 5, height = 5)  
summary(p)
```



A complete example

COVID genomes



**Questions?
Practice!!!**