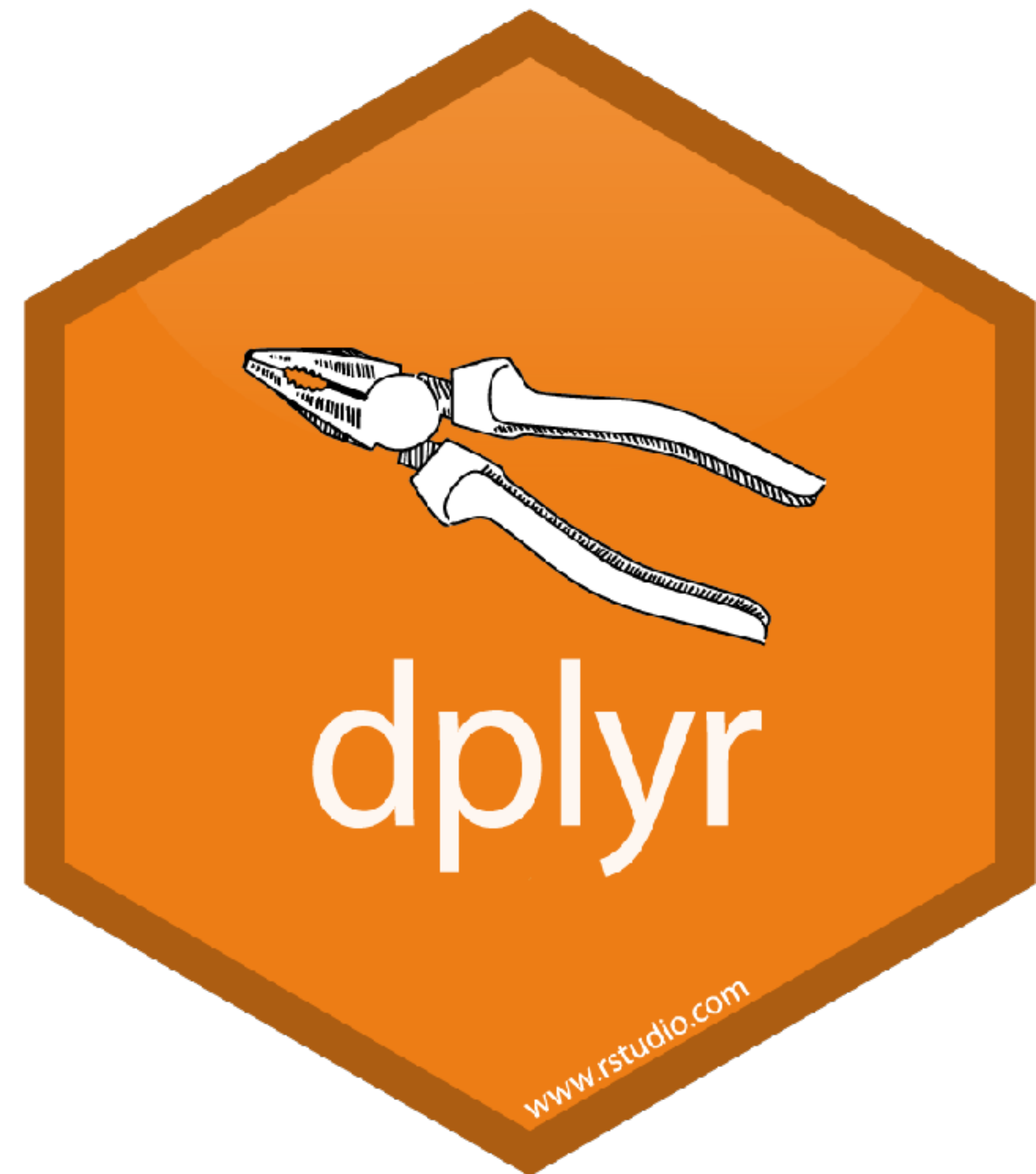


Lecture 8

Week 4, Nov 17th 2020
(4.3)

Lecture 8

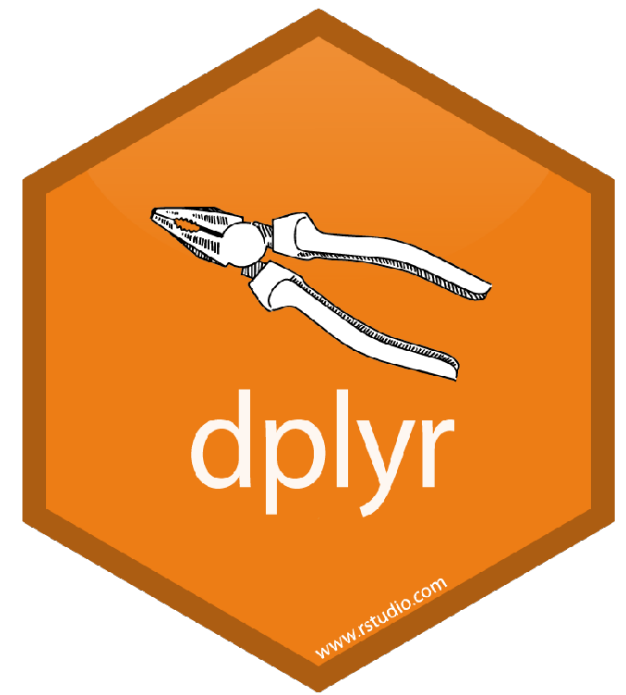
- Data transformation
 - `dplyr` package
- Tidy data
- Importing data
- Exporting images
- Labels
- Scales
- Choropleth map



Data Transformation

`dplyr + tidyr`

dplyr



- You are going to learn the 7 key dplyr functions that allow you to solve the vast majority of your data manipulation challenges:

Function	Description	Equivalent SQL
<code>select()</code>	selecting columns	SELECT
<code>filter()</code>	filtering rows / subsetting	WHERE
<code>group_by()</code>	grouping data	GROUP BY
<code>summarise()</code>	summarising / aggregating data	-
<code>arrange()</code>	sorting data	ORDER BY
<code>join()</code>	joining data tables	JOIN
<code>mutate()</code>	creating new columns	COLUMN ALIAS

Basic principle

- All **dp1yr** functions work similarly:
 - The first argument is a data frame
 - Subsequent argument describe what to do
 - Output is a new data frame

Prerequisites

- Download and load packages: `library(nycflights13)`
`library(tidyverse)`

```
flights
```

```
#> # A tibble: 336,776 x 19
```

```
#>   year month   day dep_time sched_dep_time dep_delay arr_time
```

```
#>   <int> <int> <int>   <int>         <int>         <dbl>   <int>
```

```
#> 1  2013     1     1     517           515           2     830
```

```
#> 2  2013     1     1     533           529           4     850
```

```
#> 3  2013     1     1     542           540           2     923
```

```
#> 4  2013     1     1     544           545          -1    1004
```

```
#> 5  2013     1     1     554           600          -6     812
```

```
#> 6  2013     1     1     554           558          -4     740
```

```
#> # ... with 3.368e+05 more rows, and 12 more variables: sched_arr_time <int>,
```

```
#> #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
```

```
#> #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
```

```
#> #   minute <dbl>, time_hour <dtm>
```

Your turn!

1. Are there vignettes for the **dplyr** package?
2. Can you find additional documentation explaining the **flights** dataset?
3. What variables are in the **flights** dataset? How many rows of data are in the **flights** dataset?

Solution

1. Are there vignettes for the dplyr package?

```
# Lists topics  
vignette(package = "dplyr")  
# Select a topic  
vignette(package = "dplyr", topic = "dplyr")
```

2. Can you find additional documentation explaining the flights dataset?

```
?flights
```

3. What variables are in the flights dataset? How many rows of data are in the flights dataset?

```
str(flights)
```

```
dim(flights)
```


filter()

pick observations by their values

Basic filtering

```
filter(flights, month == 1)
```

```
# A tibble: 27,004 × 19
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>
1	2013	1	1	517	515	2	830	819	11
2	2013	1	1	533	529	4	850	830	20
3	2013	1	1	542	540	2	923	850	33
4	2013	1	1	544	545	-1	1004	1022	-18
5	2013	1	1	554	600	-6	812	837	-25
6	2013	1	1	554	558	-4	740	728	12
7	2013	1	1	555	600	-5	913	854	19
8	2013	1	1	557	600	-3	709	723	-14
9	2013	1	1	557	600	-3	838	846	-8
10	2013	1	1	558	600	-2	753	745	8

```
# ... with 26,994 more rows, and 10 more variables: carrier <chr>, flight <int>,
```

Basic filtering

```
filter(flights, month == 1, day == 1)
```

```
# A tibble: 842 × 19
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>
1	2013	1	1	517	515	2	830	819	11
2	2013	1	1	533	529	4	850	830	20
3	2013	1	1	542	540	2	923	850	33
4	2013	1	1	544	545	-1	1004	1022	-18
5	2013	1	1	554	600	-6	812	837	-25
6	2013	1	1	554	558	-4	740	728	12
7	2013	1	1	555	600	-5	913	854	19
8	2013	1	1	557	600	-3	709	723	-14
9	2013	1	1	557	600	-3	838	846	-8
10	2013	1	1	558	600	-2	753	745	8

```
# ... with 832 more rows, and 10 more variables: carrier <chr>, flight <int>, tailnum <chr>,
```

Basic filtering

```
filter(flights, month == 1, day == 1, dep_delay > 0)
```

```
# A tibble: 352 × 19
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>
1	2013	1	1	517	515	2	830	819	11
2	2013	1	1	533	529	4	850	830	20
3	2013	1	1	542	540	2	923	850	33
4	2013	1	1	601	600	1	844	850	-6
5	2013	1	1	608	600	8	807	735	32
6	2013	1	1	611	600	11	945	931	14
7	2013	1	1	613	610	3	925	921	4
8	2013	1	1	623	610	13	920	915	5
9	2013	1	1	632	608	24	740	728	12
10	2013	1	1	644	636	8	931	940	-9

```
# ... with 342 more rows, and 10 more variables: carrier <chr>, flight <int>, tailnum <chr>,
```

Save a new data frame

Save data frame using assignment operator "<-"

```
jan1 <- filter(flights, month == 1, day == 1)
```

Save and print filtered data frame by wrapping entire code with ()

```
(dec25 <- filter(flights, month == 12, day == 25))
```

Logical tests

```
12 == 12
```

```
12 <= c(12, 11)
```

```
12 %in% c(12, 11, 8)
```

```
x <- c(12, NA, 11, NA, 8)  
is.na(x)
```

<

Less than

>

Greater than

==

Equal to

<=

Less than or equal to

>=

Greater than or equal to

!=

Not equal to

%in%

Group membership

is.na

Is NA

!is.na

Is not NA

Comparison

- What will these operations produce?

- 1) `filter(flights, month == 12)`
- 2) `filter(flights, month != 12)`
- 3) `filter(flights, month %in% c(11, 12))`
- 4) `filter(flights, arr_delay <= 120)`
- 5) `filter(flights, !(arr_delay <= 120))`
- 6) `filter(flights, is.na(tailnum))`

Menti Quiz

Multiple logical tests

?base::Logical

```
12 == 12 & 12 < 14
```

```
12 == 12 & 12 < 10
```

```
12 == 12 | 12 < 10
```

```
any(12 == 12, 12 < 10)
```

```
all(12 == 12, 12 < 10)
```

&	boolean and
	boolean or
xor	exclusively x or y
!	not
any	any true
all	all true

Multiple comparisons

Using comma is same as using &

```
filter(flights, month == 12, day == 25)
```

```
filter(flights, month == 12 & day == 25)
```

Use %in% as a shortcut for |

```
filter(flights, month == 11 | month == 12)
```

```
filter(flights, month %in% c(11, 12))
```

Are the outputs the same?

```
filter(flights, !(arr_delay > 120 | dep_delay > 120))
```

```
filter(flights, arr_delay <= 120, dep_delay <= 120)
```

arrange()

Reorder the rows

Ordering your data

```
arrange(flights, dep_delay)
```

```
# A tibble: 336,776 × 19
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay	carrier	flight
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>	<chr>	<int>
1	2013	12	7	2040	2123	-43	40	2352	48	B6	97
2	2013	2	3	2022	2055	-33	2240	2338	-58	DL	1715
3	2013	11	10	1408	1440	-32	1549	1559	-10	EV	5713
4	2013	1	11	1900	1930	-30	2233	2243	-10	DL	1435
5	2013	1	29	1703	1730	-27	1947	1957	-10	F9	837
6	2013	8	9	729	755	-26	1002	955	7	MQ	3478
7	2013	10	23	1907	1932	-25	2143	2143	0	EV	4361
8	2013	3	30	2030	2055	-25	2213	2250	-37	MQ	4573
9	2013	3	2	1431	1455	-24	1601	1631	-30	9E	3318
10	2013	5	5	934	958	-24	1225	1309	-44	B6	375

```
# ... with 336,766 more rows, and 8 more variables: tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
```

```
# distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Ordering your data

```
arrange(flights, dep_delay, arr_delay)
# A tibble: 336,776 × 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier flight
  <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>      <dbl>   <chr>   <int>
1  2013    12     7    2040           2123      -43     40           2352        48 ↓     B6      97
2  2013     2     3    2022           2055      -33    2240           2338       -58 ↓     DL    1715
3  2013    11    10    1408           1440      -32    1549           1559       -10 ↓     EV    5713
4  2013     1    11    1900           1930      -30    2233           2243       -10 ↓     DL    1435
5  2013     1    29    1703           1730      -27    1947           1957       -10 ↓     F9     837
6  2013     8     9     729           755       -26    1002           955         7 ↓     MQ    3478
7  2013     3    30    2030           2055      -25    2213           2250       -37 ↓     MQ    4573
8  2013    10    23    1907           1932      -25    2143           2143         0 ↓     EV    4361
9  2013     5     5     934           958       -24    1225           1309      -44 ↓     B6     375
10 2013     9    18    1631           1655      -24 ↓    1812           1845      -33 ↓     AA    2223
# ... with 336,766 more rows, and 8 more variables: tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
# distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Ordering your data

```
arrange(flights, desc(dep_delay))
```

```
# A tibble: 336,776 × 19
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay	carrier	flight
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>	<chr>	<int>
1	2013	1	9	641	900	1301	1242	1530	1272	HA	51
2	2013	6	15	1432	1935	1137	1607	2120	1127	MQ	3535
3	2013	1	10	1121	1635	1126	1239	1810	1109	MQ	3695
4	2013	9	20	1139	1845	1014	1457	2210	1007	AA	177
5	2013	7	22	845	1600	1005	1044	1815	989	MQ	3075
6	2013	4	10	1100	1900	960	1342	2211	931	DL	2391
7	2013	3	17	2321	810	911	135	1020	915	DL	2119
8	2013	6	27	959	1900	899	1236	2226	850	DL	2007
9	2013	7	22	2257	759	898	121	1026	895	DL	2047
10	2013	12	5	756	1700	896	1058	2020	878	AA	172

```
# ... with 336,766 more rows, and 8 more variables: tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
```

```
# distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Missing data order

- Note: missing values are always sorted at the end:

```
df <- tibble(x = c(5, 2, 5, NA))  
# A tibble: 4 × 1  
      x  
  <dbl>  
1     5  
2     2  
3     5  
4    NA
```

```
arrange(df, x)  
# A tibble: 4 × 1  
      x  
  <dbl>  
1     2  
2     5  
3     5  
4    NA
```

```
arrange(df, desc(x))  
# A tibble: 4 × 1  
      x  
  <dbl>  
1     5  
2     5  
3     2  
4    NA
```

`select()`

Pick variables by their names

Selecting variables

```
select(flights, year, month, day)
# A tibble: 336,776 × 3
   year month   day
  <int> <int> <int>
1  2013     1     1
2  2013     1     1
3  2013     1     1
4  2013     1     1
5  2013     1     1
6  2013     1     1
7  2013     1     1
8  2013     1     1
9  2013     1     1
10 2013     1     1
# ... with 336,766 more rows
```

```
select(flights, year:day)
# A tibble: 336,776 × 3
   year month   day
  <int> <int> <int>
1  2013     1     1
2  2013     1     1
3  2013     1     1
4  2013     1     1
5  2013     1     1
6  2013     1     1
7  2013     1     1
8  2013     1     1
9  2013     1     1
10 2013     1     1
# ... with 336,766 more rows
```


Deselecting variable

- Deselect with “-” sign

```
select(flights, -(year:day))
```

```
# A tibble: 336,776 × 16
```

	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time	arr_delay	carrier	flight
	<int>	<int>	<dbl>	<int>	<int>	<dbl>	<chr>	<int>
1	517	515	2	830	819	11	UA	1545
2	533	529	4	850	830	20	UA	1714
3	542	540	2	923	850	33	AA	1141
4	544	545	-1	1004	1022	-18	B6	725
5	554	600	-6	812	837	-25	DL	461
6	554	558	-4	740	728	12	UA	1696
7	555	600	-5	913	854	19	B6	507
8	557	600	-3	709	723	-14	EV	5708
9	557	600	-3	838	846	-8	B6	79
10	558	600	-2	753	745	8	AA	301

Useful select functions

-	Select everything but
:	Select range
<code>contains()</code>	Select columns whose name contains a character string
<code>ends_with()</code>	Select columns whose name ends with a string
<code>everything()</code>	Select every column
<code>matches()</code>	Select columns whose name matches a regular expression
<code>num_range()</code>	Select columns named x1, x2, x3, x4, x5
<code>one_of()</code>	Select columns whose names are in a group of names
<code>starts_with()</code>	Select columns whose name starts with a character string

select by name patterns

```
select(flights, ends_with("time"))
```

```
# A tibble: 336,776 × 5
```

	dep_time	sched_dep_time	arr_time	sched_arr_time	air_time
	<int>	<int>	<int>	<int>	<dbl>
1	517	515	830	819	227
2	533	529	850	830	227
3	542	540	923	850	160
4	544	545	1004	1022	183
5	554	600	812	837	116
6	554	558	740	728	150
7	555	600	913	854	158
8	557	600	709	723	53
9	557	600	838	846	140
10	558	600	753	745	138

```
# ... with 336,766 more rows
```

Multiple name patterns

```
select(flights, c(carrier, ends_with("time"), contains("delay")))
```

```
# A tibble: 336,776 × 8
```

	carrier	dep_time	sched_dep_time	arr_time	sched_arr_time	air_time	dep_delay	arr_delay
	<chr>	<int>	<int>	<int>	<int>	<dbl>	<dbl>	<dbl>
1	UA	517	515	830	819	227	2	11
2	UA	533	529	850	830	227	4	20
3	AA	542	540	923	850	160	2	33
4	B6	544	545	1004	1022	183	-1	-18
5	DL	554	600	812	837	116	-6	-25
6	UA	554	558	740	728	150	-4	12
7	B6	555	600	913	854	158	-5	19
8	EV	557	600	709	723	53	-3	-14
9	B6	557	600	838	846	140	-3	-8
10	AA	558	600	753	745	138	-2	8

```
# ... with 336,766 more rows
```

Reorder variables

```
select(flights, time_hour, air_time, everything())
```

```
# A tibble: 336,776 × 19
```

	time_hour	air_time	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time
	<dtm>	<dbl>	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>
1	2013-01-01 05:00:00	227	2013	1	1	517	515	2	830
2	2013-01-01 05:00:00	227	2013	1	1	533	529	4	850
3	2013-01-01 05:00:00	160	2013	1	1	542	540	2	923
4	2013-01-01 05:00:00	183	2013	1	1	544	545	-1	1004
5	2013-01-01 06:00:00	116	2013	1	1	554	600	-6	812
6	2013-01-01 05:00:00	150	2013	1	1	554	558	-4	740
7	2013-01-01 06:00:00	158	2013	1	1	555	600	-5	913
8	2013-01-01 06:00:00	53	2013	1	1	557	600	-3	709
9	2013-01-01 06:00:00	140	2013	1	1	557	600	-3	838
10	2013-01-01 06:00:00	138	2013	1	1	558	600	-2	753

```
# ... with 336,766 more rows, and 10 more variables: sched_arr_time <int>, arr_delay <dbl>,
```

```
#   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>, distance <dbl>,
```

Renaming variables

```
rename(flights, newName = dep_delay)
```

```
# A tibble: 336,776 × 19
```

	year	month	day	dep_time	sched_dep_time	newName	arr_time	sched_arr_time	arr_delay
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>	<dbl>
1	2013	1	1	517	515	2	830	819	11
2	2013	1	1	533	529	4	850	830	20
3	2013	1	1	542	540	2	923	850	33
4	2013	1	1	544	545	-1	1004	1022	-18
5	2013	1	1	554	600	-6	812	837	-25
6	2013	1	1	554	558	-4	740	728	12
7	2013	1	1	555	600	-5	913	854	19
8	2013	1	1	557	600	-3	709	723	-14
9	2013	1	1	557	600	-3	838	846	-8
10	2013	1	1	558	600	-2	753	745	8

```
# ... with 336,766 more rows, and 10 more variables: carrier <chr>, flight <int>,
```

```
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
```

Your turn!

- 1) What happens if you include the name of a variable multiple times in a **select()** call?
- 2) What does the **any_of()** function do? Why might it be helpful in conjunction with this vector?

```
vars <- c("year", "month", "day", "dep_delay", "arr_delay")
```
- 3) Does the result of running the following code surprise you? How do the select helpers deal with case by default? How can you change that default?

```
select(flights, contains("TIME"))
```

Solution

```
# 1. what happens if you call a variable multiple times in select()
```

```
select(flights, month, month)
```

```
# A tibble: 336,776 × 1
```

```
  month
```

```
  <int>
```

```
1      1
```

```
2      1
```

```
3      1
```

```
4      1
```

```
5      1
```

```
6      1
```

```
7      1
```


Solution

```
# 2. what does any_of() do?
```

```
# Matches variable names in a character vector.
```

```
vars <- c("year", "month", "day", "dep_delay", "arr_delay")  
select(flights, any_of(vars))
```

```
## # A tibble: 336,776 x 5
```

```
##   year month   day dep_delay arr_delay
```

```
##   <int> <int> <int>      <dbl>      <dbl>
```

```
## 1  2013     1     1         2         11
```

```
## 2  2013     1     1         4         20
```

```
## 3  2013     1     1         2         33
```

```
## 4  2013     1     1        -1        -18
```

```
## 5  2013     1     1        -6        -25
```

Solution

#3 Does the result of running the following code surprise you? How do the select helpers deal with case by default? How can you change that default?

The default helper functions are insensitive to case. This can be changes by setting `ignore.case=FALSE`.

```
select(flights, contains("TIME"))  
select(flights, contains("TIME", ignore.case = FALSE))
```

mutate()

create new variables

create new variables

```
# Smaller dataset for demo
flights_sml <- select(flights,
  year:day,
  ends_with("delay"),
  distance,
  air_time
)

mutate(flights_sml,
  gain = dep_delay - arr_delay,
  speed = distance / air_time * 60
)
```

Useful tips

```
# You can refer to columns that you've just created  
mutate(flights_sml,  
  gain = dep_delay - arr_delay,  
  hours = air_time / 60,  
  gain_per_hour = gain / hours  
)
```

```
# If you only want to keep the new variables  
transmute(flights,  
  gain = dep_delay - arr_delay,  
  hours = air_time / 60,  
  gain_per_hour = gain / hours  
)
```

Useful creation functions

- There are a wide variety of functions you can use with **mutate()**

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

arithmetic

```
transmute(flights,  
  normalized_delay =  
    dep_delay/(mean(dep_delay, na.rm = TRUE)))
```

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

Modular arithmetic

```
transmute(flights,
  dep_time,
  hour = dep_time %/% 100,
  minute = dep_time %% 100
)
```

%/% (Integer division)
%% (Remainder)

```
#> # A tibble: 336,776 x 3
#>   dep_time hour minute
#>   <int> <dbl> <dbl>
#> 1     517     5     17
#> 2     533     5     33
#> 3     542     5     42
#> 4     544     5     44
#> 5     554     5     54
```

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

non-linear transformations

```
transmute(flights,  
  log_air_time = log2(air_time),  
  exp_delay = exp(dep_delay))
```

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

offsets & cumulative aggregates

```
transmute(flights,  
  dep_delay = dep_delay,  
  lag_delay = lag(dep_delay),  
  sum_delay = cumsum(dep_delay))
```

```
## # A tibble: 336,776 x 3  
##   dep_delay lag_delay sum_delay  
##   <dbl>      <dbl>    <dbl>  
## 1         2      NA         2  
## 2         4        2         6  
## 3         2        4         8  
## 4        -1        2         7  
## 5        -6       -1         1
```

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

Ranking

```
y <- c(1, 2, 2, NA, 3, 4)

min_rank(y)
#> [1] 1 2 2 NA 4 5

row_number(y)
#> [1] 1 2 3 NA 4 5

dense_rank(y)
#> [1] 1 2 2 NA 3 4

percent_rank(y)
#> [1] 0.00 0.25 0.25 NA 0.75 1.00

cume_dist(y)
#> [1] 0.2 0.6 0.6 NA 0.8 1.0
```

Functions	Description
<code>+, -, *, /, ^</code>	arithmetic
<code>x / sum(x)</code>	arithmetic w/aggregate functions
<code>%/%, %%</code>	modular arithmetic
<code>log, exp, sqrt</code>	transformations
<code>lag, lead</code>	offsets
<code>cumsum, cumprod, cum...</code>	cum/rolling aggregates
<code>>, >=, <, <=, !=, ==</code>	logical comparisons
<code>min_rank, dense_rank, etc</code>	ranking
<code>between</code>	are values between a and b?
<code>ntile</code>	bin values into buckets

ntile

```
transmute(flights,  
  arr_delay = arr_delay,  
  bucket = ntile(arr_delay, 10))
```

```
## # A tibble: 336,776 x 2  
##   arr_delay bucket  
##   <dbl>   <int>  
## 1      11      8  
## 2     20      8  
## 3     33      9  
## 4    -18      3  
## 5    -25      2  
## 6     12      8  
## 7     19      8  
## 8    -14      3  
## 9     -8      5  
## 10      8      7  
## # ... with 336,766 more rows
```

Functions	Description
+, -, *, /, ^	arithmetic
x / sum(x)	arithmetic w/aggregate functions
%/%, %%	modular arithmetic
log, exp, sqrt	transformations
lag, lead	offsets
cumsum, cumprod, cum...	cum/rolling aggregates
>, >=, <, <=, !=, ==	logical comparisons
min_rank, dense_rank, etc	ranking
between	are values between a and b?
ntile	bin values into buckets

Your turn!

- 1) Create a new variable **distance_km** that converts distance in miles to kilometres
- 2) Create a **time_per_km** variable based on **air_time** and **distance_km**.

Solution

```
transmute(flights,  
  distance_km = distance * 1.60934,  
  time_per_km = air_time / distance_km)
```

summarise() & group_by()

Create grouped summaries

summarise()

- **summarise()** collapses a data frame into a single row

```
summarise(flights, delay = mean(dep_delay, na.rm = TRUE))
```

```
#> # A tibble: 1 x 1  
#>   delay  
#>   <dbl>  
#> 1  12.6
```

Why is this important?
Try without this argument.

- **summarise()** is not really useful without pairing it with **group_by()**

group_by()

country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3

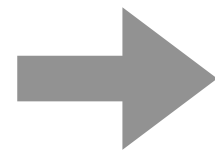


country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3

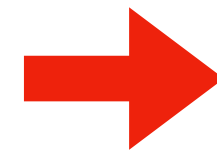
```
group_by(data, country)
```

group_by()

country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



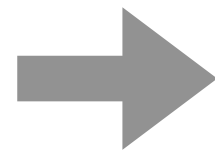
country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3

`group_by(data, country)`

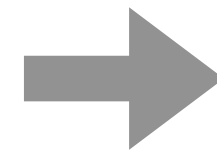
`group_by(data, country, year)`

group_by()

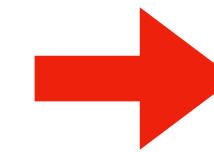
country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3



country	year	sex	cases
Afghanistan	1999	female	1
Afghanistan	1999	male	1
Afghanistan	2000	female	1
Afghanistan	2000	male	1
Brazil	1999	female	2
Brazil	1999	male	2
Brazil	2000	female	2
Brazil	2000	male	2
China	1999	female	3
China	1999	male	3
China	2000	female	3
China	2000	male	3

`group_by(data, country)`

`group_by(data, country, year)`

`ungroup(data)`

Pipe operator

Chaining functions together

Example

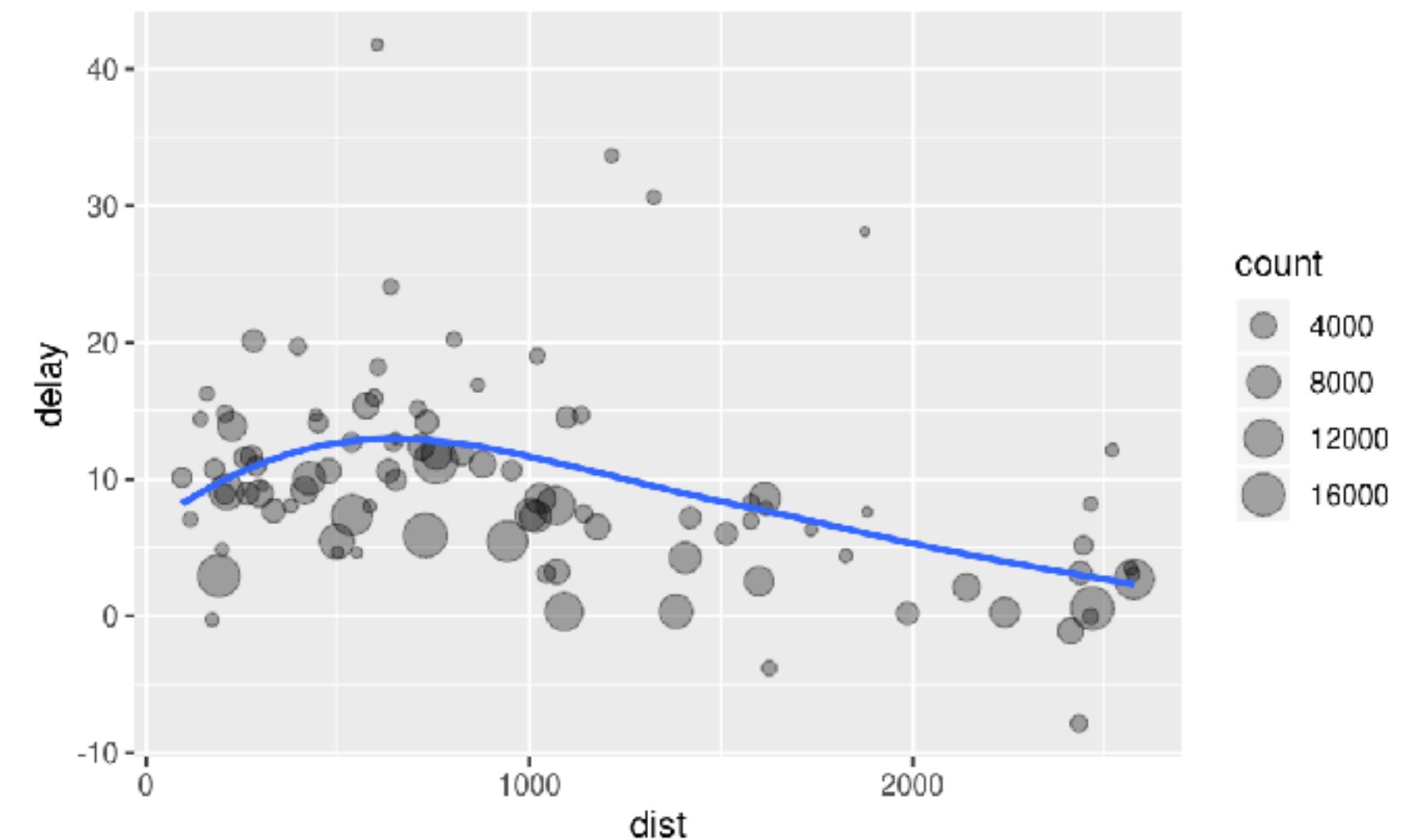
- Imagine that we want to explore the relationship between the distance and average delay for each location. Using what you know about dplyr, you might write code like this:

```
by_dest <- group_by(flights, dest)

delay <- summarise(by_dest,
  count = n(),
  dist = mean(distance, na.rm = TRUE),
  delay = mean(arr_delay, na.rm = TRUE)
)

delay <- filter(delay, count > 20, dest != "HNL")

ggplot(data = delay, mapping = aes(x = dist, y = delay)) +
  geom_point(aes(size = count), alpha = 1/3) +
  geom_smooth(se = FALSE)
```



Example

- Imagine that we want to explore the relationship between the distance and average delay for each location. Using what you know about dplyr, you might write code like this:

```
by_dest <- group_by(flights, dest)
```

grouping by destination

```
delay <- summarise(by_dest,  
  count = n(),  
  dist = mean(distance, na.rm = TRUE),  
  delay = mean(arr_delay, na.rm = TRUE)  
)
```

summarising count, distance and arrival delay

```
delay <- filter(delay, count > 20, dest != "HNL")
```

filtering out low counts and Honolulu

```
ggplot(data = delay, mapping = aes(x = dist, y = delay)) +  
  geom_point(aes(size = count), alpha = 1/3) +  
  geom_smooth(se = FALSE)
```

creating a plot

Example

```
by_dest <- group_by(flights, dest)

delay <- summarise(by_dest,
  count = n(),
  dist = mean(distance, na.rm = TRUE),
  delay = mean(arr_delay, na.rm = TRUE)
)

delay <- filter(delay, count > 20, dest != "HNL")
```

```
delays <- flights %>%
  group_by(dest) %>%
  summarise(
    count = n(),
    dist = mean(distance, na.rm = TRUE),
    delay = mean(arr_delay, na.rm = TRUE)
  ) %>%
  filter(count > 20, dest != "HNL")
```

Command/CTRL + Shift + M

Pipe operator

```
# without pipe to subset "suv" cars  
filter(mpg, class == "suv")  
  
# equivalent form with pipe  
mpg %>% filter(class == "suv")
```

Command/CTRL + Shift + M

*“Tidy datasets are all alike,
but every messy dataset is messy in its own way.”
– Hadley Wickham*

Tidy data

data structure

Data Import Cheatsheet

Data Import: : CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.

The front side of this sheet shows how to read text files into R with **readr**.

The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

Save Data

Save **x**, an R object, to **path**, a file path, as:

Comma delimited file
`write_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

File with arbitrary delimiter
`write_delim(x, path, delim = ";", na = "NA", append = FALSE, col_names = !append)`

CSV for excel
`write_excel_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

String to file
`write_file(x, path, append = FALSE)`

String vector to file, one element per line
`write_lines(x, path, na = "NA", append = FALSE)`

Object to RDS file
`write_rds(x, path, compress = c("none", "gz", "bz2", "xz", ...))`

Tab delimited files
`write_tsv(x, path, na = "NA", append = FALSE, col_names = !append)`

Read Tabular Data

These functions share the common arguments:

```
read_* (file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),  
quoted_na = TRUE, comment = "#", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,  
n_max), progress = interactive())
```

Comma Delimited Files
`read_csv("file.csv")`
To make file.csv run:
`write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")`

Semi-colon Delimited Files
`read_csv2("file2.csv")`
`write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")`

Files with Any Delimiter
`read_delim("file.txt", delim = ";")`
`write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file.txt")`

Fixed Width Files
`read_fwf("file.fwf", col_positions = c(1, 3, 5))`
`write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")`

Tab Delimited Files
`read_tsv("file.tsv")` Also `read_table()`
`write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")`

USEFUL ARGUMENTS

Example file <code>write_file("a,b,c\n1,2,3\n4,5,NA", "file.csv")</code> <code>f <- "file.csv"</code>	Skip lines <code>read_csv(f, skip = 1)</code>
No header <code>read_csv(f, col_names = FALSE)</code>	Read in a subset <code>read_csv(f, n_max = 1)</code>
Provide header <code>read_csv(f, col_names = c("x", "y", "z"))</code>	Missing Values <code>read_csv(f, na = c("1", " "))</code>

Read Non-Tabular Data

Read a file into a single string <code>read_file(file, locale = default_locale())</code>	Read a file into a raw vector <code>read_file_raw(file)</code>
Read each line into its own string <code>read_lines(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())</code>	Read each line into a raw vector <code>read_lines_raw(file, skip = 0, n_max = -1L, progress = interactive())</code>
Read Apache style log files <code>read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())</code>	

Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specifications:  
##      col1  
##      age = col_integer(),    age is an integer  
##      sex = col_character(),  sex is a character  
##      eern = col_double()     eern is a double (numeric)
```

1. Use `problems()` to diagnose problems.
`x <- read_csv("file.csv"); problems(x)`

2. Use a `col_` function to guide parsing.
- `col_guess()` - the default
 - `col_character()`
 - `col_double()`, `col_euro_double()`
 - `col_datetime()` (format = "%") Also `col_date()` (format = "%")
 - `col_factor()` (levels, ordered = FALSE)
 - `col_integer()`
 - `col_logical()`
 - `col_number()`, `col_numeric()`
 - `col_skip()`
- `x <- read_csv("file.csv", col_types = cols(
A = col_double(),
B = col_logical(),
C = col_factor()))`

3. Use `readr` as character vectors then parse with a `parse_` function.

- `parse_guess()`
 - `parse_character()`
 - `parse_datetime()` Also `parse_date()` and `parse_time()`
 - `parse_double()`
 - `parse_factor()`
 - `parse_integer()`
 - `parse_logical()`
 - `parse_number()`
- `x$A <- parse_number(x$A)`

Tibbles - an enhanced data frame

The **tibble** package provides a new S3 class for storing tabular data, the **tibble**. Tibbles inherit the data frame class, but improve three behaviors:

- **Subsetting** - `[]` always returns a new tibble, `[[` and `$` always return a vector.
- **No partial matching** - You must use full column names when subsetting
- **Display** - When you print a tibble, R provides a concise view of the data that fits on one screen

A large table to display

tibble display

data frame display

- Control the default appearance with options:
`options(tibble.print_max = n,
tibble.print_min = n, tibble.width = Inf)`
- View full data set with `View()` or `glimpse()`
- Revert to data frame with `as.data.frame()`

CONSTRUCT A TIBBLE IN TWO WAYS

`tibble(...)`
Construct by columns.
`tibble(x = 1:3, y = c("a", "b", "c"))`

`tribble(...)`
Construct by rows.
`tribble(~x, ~y,
1, "a",
2, "b",
3, "c")`

Both make this tibble

`as_tibble(x, ...)` Convert data frame to tibble

`enframe(x, name = "name", value = "value")`
Convert named vector to a tibble

`is_tibble(x)` Test whether x is a tibble.

Tidy Data with tidyr

Tidy data is a way to organize tabular data. It provides a consistent data structure across packages.

A table is tidy if:

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

Reshape Data - change the layout of values in a table

Use `gather()` and `spread()` to reorganize the values of a table into a new layout.

`gather(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor_key = FALSE)`

`gather()` moves column names into a **key** column, gathering the column values into a single **value** column.

`spread(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)`

`spread()` moves the unique values of a **key** column into the column names, spreading the values of a **value** column across the new columns.

`gather(table4, "1999", "2000",
key = "year", value = "cases")`

Handle Missing Values

`drop_na(data, ...)`

Drop rows containing NA's in ... columns

`fill(data, ..., direction = c("down", "up"))`

Fill in NA's in ... columns with most recent non-NA values.

`replace_na(data, replace = list(), ...)`

Replace NA's by column.

`complete(data, ..., fill = list())`

Adds to the data missing combinations of the values of the variables listed in ...

`expand(data, ...)`

Create new tibble with all possible combinations of the values of the variables listed in ...

Split Cells

Use these functions to split or combine cells into individual, isolated values.

`separate(data, col, into, sep = "[^:alnum:]")`
+ `remove = TRUE`, `convert = FALSE`, `extra = "warn"`, `fill = "warn"`, ...)

Separate each cell in a column to make several columns.

`separate_rows(data, ..., sep = "[^:alnum:]")`
+ `convert = FALSE`

Separate each cell in a column to make several rows.

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite_rows(table3, col, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`

`unite(table3, col, into, sep = " ", na.rm = TRUE)`



3 rules of tidy data

- Each variable must have its own column
- Each observation must have its own row
- Each value must have its own cell

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	1280425583

variables

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	1280425583

observations

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	1280425583

values

Your turn!

Is this a tidy table?

Table 1:

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

Your turn!

Is this a tidy table?

Table 2:

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

Your turn!

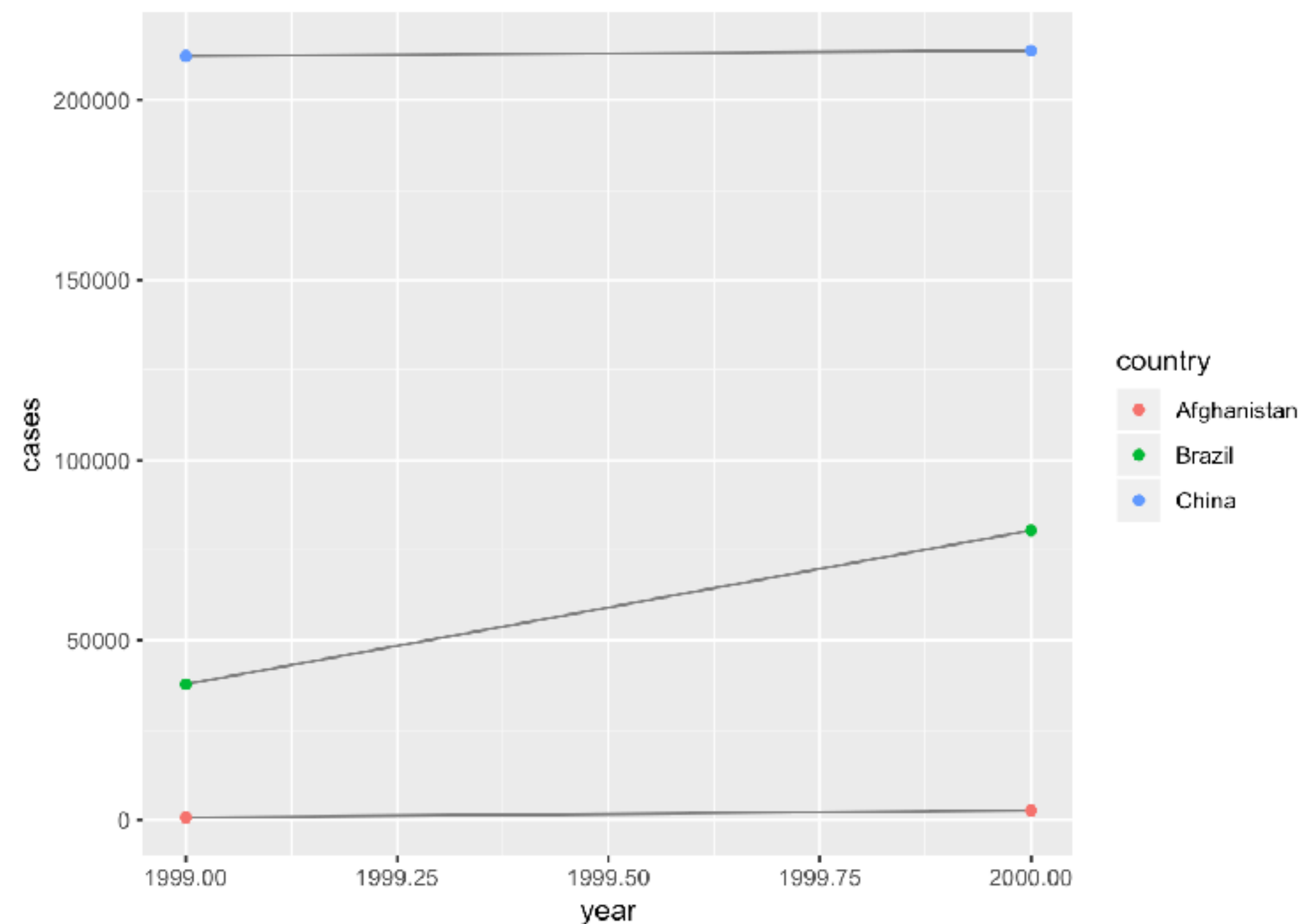
Is this a tidy table?

Table 3:

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

Tidy data example

```
# using a sample dataset from tidyr (tidyr::table1)
# check the table View(tidyr::table1)
ggplot(data = table1, mapping = aes(year, cases)) +
  geom_line(mapping = aes(group = country), colour = "grey50") +
  geom_point(mapping = aes(colour = country))
```



4 key functions

Function	Old name	Description
<code>pivot_longer()</code>	<code>gather()</code>	transforms data from wide to long
<code>pivot_wider()</code>	<code>spread()</code>	transforms data from long to wide
<code>separate()</code>		splits a single column into multiple columns
<code>unite()</code>		combines multiple columns into a single column

country	year	cases	population
Afghanistan	1999	37745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

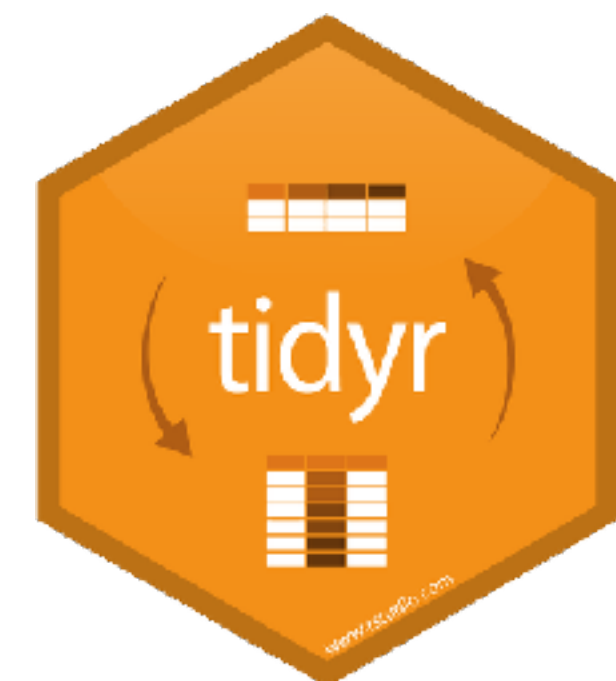
variables

country	year	cases	population
Afghanistan	1999	37745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

observations

country	year	cases	population
Afghanistan	1999	37745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	128042583

values



`pivot_longer()`

Transform data from wide to long

pivot_longer() example

```
table4a %>%  
  pivot_longer(cols = c('1999', '2000'), names_to = "year", values_to = "cases")
```

country	1999	2000
Afghanistan	745	2666
Brazil	37737	80488
China	212258	213766

pivot_longer()

country	year	cases
Afghanistan	1999	745
Afghanistan	2000	2666
Brazil	1999	37737
Brazil	2000	80488
China	1999	212258
China	2000	213766

Your turn!

1. Download the data (**bomber_wide.rds**) in the data folder
2. Import the .rds with **read_rds()** function
3. Reshape this data from wide to long

Solution

```
# One example
df <- read_rds("Lecture_8_data/bomber_wide.rds")
df %>%
  pivot_longer(cols = `1996`:`2014`,
               names_to = "year",
               values_to = "value")
```

`pivot_wider()`

Transform data from long to wide

pivot_wider() example

```
table2 %>%  
  pivot_wider(names_from = type, values_from = count)
```

		names_from	values_from
country	year	type	count
Afghanistan	1999	cases	745
Afghanistan	1999	population	19987071
Afghanistan	2000	cases	2666
Afghanistan	2000	population	20595360
Brazil	1999	cases	37737
Brazil	1999	population	172006362
Brazil	2000	cases	80488
Brazil	2000	population	174504898
China	1999	cases	212258
China	1999	population	1272915272
China	2000	cases	213766
China	2000	population	1280428583

pivot_wider()

		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

Your turn!

1. Download the data (**bomber_long.rds**) in the data folder
2. Import the .rds with **read_rds()** function
3. Reshape this data from long to wide

Solution

```
df <- read_rds("Lecture_8_data/bomber_long.rds")
df %>%
  pivot_wider( names_from = Output,
               values_from = Value)
```


separate()

split a single column into multiple columns

separate() example

```
table3 %>%  
  separate(rate, into = c("cases", "population"), sep="/", convert = TRUE)
```

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

table3

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

separate() example

```
table3 %>%  
  separate(year, into = c("century", "year"), sep = 2)  
  
#> # A tibble: 6 x 4  
#>   country      century year      rate  
#>   <chr>      <chr>    <chr> <chr>  
#> 1 Afghanistan 19      99    745/19987071  
#> 2 Afghanistan 20      00    2666/20595360  
#> 3 Brazil      19      99    37737/172006362  
#> 4 Brazil      20      00    80488/174504898  
#> 5 China       19      99    212258/1272915272  
#> 6 China       20      00    213766/1280428583
```

Your turn!

1. Download the data (**bomber_combined.rds**) in the data folder
2. Import the .rds with **read_rds()** function
3. Separate the AC variable into “Type” and “MD”

Solution

```
read_rds("data/bomber_combined.rds") %>%  
  separate(AC, into = c("Type", "MD"), sep = " ")
```


	Type	MD	FY	Cost	FH	Gallons
1	Bomber	B-1	1996	72753781	26914	88594449
2	Bomber	B-1	1997	71297263	25219	85484074
3	Bomber	B-1	1998	84026805	24205	85259038
4	Bomber	B-1	1999	71848336	23306	79323816
5	Bomber	B-1	2000	58439777	25013	86230284
6	Bomber	B-1	2001	94946077	25059	86892432

`unite()`

combine multiple columns into a single column

unite() example

```
table5 %>%  
  unite(new, century, year, sep = "")
```



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

Your turn!

1. Download the data (bomber_prefix.rds) in the data folder
2. Import the .rds with **read_rds()** function
3. Unite the prefix and number columns into a “MD” variable with “-” as separator

Solution

```
read_rds("data/bomber_combined.rds") %>%  
  unite(MD, prefix, number, sep = "_")
```

	Type	MD	FY	Output	Value
1	Bomber	B-1	1996	FH	26914
2	Bomber	B-1	1997	FH	25219
3	Bomber	B-1	1998	FH	24205
4	Bomber	B-1	1999	FH	23306
5	Bomber	B-1	2000	FH	25013

Your turn!

1. Download the data (**bomber_mess.rds**) in the data folder
2. Import the .rds with `read_rds()` function
3. Clean this data up so it looks like:

```
# A tibble: 57 × 6
  Type      MD    FY      Cost    FH  Gallons
*   <chr> <chr> <chr>    <int> <int>    <int>
1 Bomber   B-1  1996  72753781 26914  88594449
2 Bomber   B-1  1997  71297263 25219  85484074
3 Bomber   B-1  1998  84026805 24205  85259038
4 Bomber   B-1  1999  71848336 23306  79323816
5 Bomber   B-1  2000  58439777 25013  86230284
6 Bomber   B-1  2001  94946077 25059  86892432
7 Bomber   B-1  2002  96458536 26581  89198262
8 Bomber   B-1  2003  68650070 21491  74485788
9 Bomber   B-1  2004 101895634 28118 101397707
10 Bomber  B-1  2005 124816690 21859  78410415
# ... with 47 more rows
```

Solution

```
read_rds("data/bomber_mess.rds") %>%  
  unite(col = MD, prefix:number, sep = "-") %>%  
  separate(Metric, into = c("FY", "Output")) %>%  
  spread(Output, Value) %>%  
  as_tibble()
```

```
# A tibble: 57 × 6
```

	Type	MD	FY	Cost	FH	Gallons
*	<chr>	<chr>	<chr>	<int>	<int>	<int>
1	Bomber	B-1	1996	72753781	26914	88594449
2	Bomber	B-1	1997	71297263	25219	85484074
3	Bomber	B-1	1998	84026805	24205	85259038
4	Bomber	B-1	1999	71848336	23306	79323816
5	Bomber	B-1	2000	58439777	25013	86230284
6	Bomber	B-1	2001	94946077	25059	86892432
7	Bomber	B-1	2002	96458536	26581	89198262
8	Bomber	B-1	2003	68650070	21491	74485788
9	Bomber	B-1	2004	101895634	28118	101397707
10	Bomber	B-1	2005	124816690	21859	78410415

```
# ... with 47 more rows
```

Data Wrangling

Practice

Your turn!

Let's apply **dplyr** and **tidyr** functions you have learned.

- **tidyr::who** is a dataset of TB cases broken down by year, country, age, gender and diagnosis methods
- Examine the table and its documentation, and clean the data into a tidy data
- Hint:
 1. Gather columns and remove missing values
 2. Separate the values like "new_sp_m014" into columns
 3. Separate the gender and age
 4. Drop redundant columns using **dplyr::select()**
- Create a graph with the tidy data

Solution

1. Pivot lonter new_sp_m014:newrel_f65 columns and remove missing values

```
who1 <- who %>%  
  pivot_longer(  
    cols = new_sp_m014:newrel_f65,  
    names_to = "key",  
    values_to = "cases",  
    values_drop_na = TRUE  
  )  
  
## # A tibble: 76,046 x 6  
##   country    iso2 iso3   year key      cases  
##   * <chr>    <chr> <chr> <int> <chr>    <int>  
## 1 Afghanistan AF    AFG   1997 new_sp_m014    0  
## 2 Afghanistan AF    AFG   1998 new_sp_m014   30  
## 3 Afghanistan AF    AFG   1999 new_sp_m014    8  
## 4 Afghanistan AF    AFG   2000 new_sp_m014   52  
## 5 Afghanistan AF    AFG   2001 new_sp_m014  129
```

Solution

2. Clean up the key column, e.g. new_sp_m014

- **new** encodes new cases of TB
- next 3 letters describe the type of TB:
 - **rel** = relapse
 - **ep** = extrapulmonary TB
 - **sn** = could not be diagnosed by a pulmonary smear (smear negative)
 - **sp** = smear positive
- **m** = males, **f** = females
- age group: **014** = 1 - 14 years old age group

```
## # A tibble: 76,046 x 6
##   country      iso2 iso3  year key      cases
##   * <chr>      <chr> <chr> <int> <chr>    <int>
## 1 Afghanistan AF    AFG  1997 new_sp_m014      0
## 2 Afghanistan AF    AFG  1998 new_sp_m014     30
```

Solution

2. But, there is an inconsistency in naming the key.

- eg. `new_rel` and `newrel`
- use `stringr::str_replace()` function to replace `newrel` with `new_rel`


```
who2 <- who1 %>%  
  mutate(key = stringr::str_replace(key, "newrel", "new_rel"))
```

	country <chr>	iso2 <chr>	iso3 <chr>	year <int>	key <chr>	cases <int>
1	Afghanistan	AF	AFG	1997	new_sp_m014	0
2	Afghanistan	AF	AFG	1998	new_sp_m014	30
3	Afghanistan	AF	AFG	1999	new_sp_m014	8
4	Afghanistan	AF	AFG	2000	new_sp_m014	52
5	Afghanistan	AF	AFG	2001	new_sp_m014	129

Solution

2. Now use `separate()` to split key into 3 columns

```
who3 <- who2 %>%  
  separate(key, into = c("new", "type", "sex-age"), sep = "_")
```



	key	##	country	iso2	iso3	year	new	type	`sex-age`	cases
	<chr>	##	<chr>	<chr>	<chr>	<int>	<chr>	<chr>	<chr>	<int>
1	new_sp_m014	##	1 Afghanistan	AF	AFG	1997	new	sp	m014	0
2	new_sp_m1524	##	2 Afghanistan	AF	AFG	1998	new	sp	m014	30
3	new_sp_m2534	##	3 Afghanistan	AF	AFG	1999	new	sp	m014	8
4	new_sp_m3544	##	4 Afghanistan	AF	AFG	2000	new	sp	m014	52
5	new_sp_m4554	##	5 Afghanistan	AF	AFG	2001	new	sp	m014	129

Solution

3. Separate sex-age into 2 columns

```
who4 <- who3 %>%  
  separate(`sex-age`, into = c("sex", "age"), sep = 1)
```



`sex-age` <chr>	##	country <chr>	iso2 <chr>	iso3 <chr>	year <int>	new <chr>	type <chr>	sex <chr>	age <chr>	cases <int>
1 m014	##	1 Afghanistan	AF	AFG	1997	new	sp	m	014	0
2 m1524	##	2 Afghanistan	AF	AFG	1998	new	sp	m	014	30
3 m2534	##	3 Afghanistan	AF	AFG	1999	new	sp	m	014	8
4 m3544	##	4 Afghanistan	AF	AFG	2000	new	sp	m	014	52
5 m4554	##	5 Afghanistan	AF	AFG	2001	new	sp	m	014	129

Solution

4. Drop redundant columns:

```
who5 <- who4 %>%  
  select(-new, -iso2, -iso3)
```

	country <chr>	year <int>	var <chr>	sex <chr>	age <chr>	value <int>
1	Afghanistan	1997	sp	m	014	0
2	Afghanistan	1998	sp	m	014	30
3	Afghanistan	1999	sp	m	014	8
4	Afghanistan	2000	sp	m	014	52
5	Afghanistan	2001	sp	m	014	129

Solution

You can build up a complete piped operations as shown below:

```
who %>%  
  pivot_longer(cols = new_sp_m014:newrel_f65, names_to="key",  
               values_to = "cases", values_drop_na = TRUE) %>%  
  mutate(key = stringr::str_replace(key, "newrel", "new_rel")) %>%  
  separate(key, c("new", "var", "sexage")) %>%  
  separate(sexage, c("sex", "age"), sep = 1) %>%  
  select(-new, -iso2, -iso3)
```

Importing data

readr package

Tabular data

- to import tabular data, use functions from **readr** package (**tidyverse**)
 - It is much faster than base R function such as `read.csv()`
 - It loads as **tibble** instead of R's traditional `data.frame`

Function	Description
<code>read_csv()</code>	comma separated (csv) files
<code>read_tsv()</code>	tab separated files
<code>read_delim()</code>	general delimited files
<code>read_fwf()</code>	fixed width files
<code>read_table()</code>	tabular files where columns are separated by white-space

example

```
car_data <- read_csv(file = "data/mtcars.csv")
```

```
## Parsed with column specification:  
## cols(  
##   mpg = col_double(),  
##   cyl = col_integer(),  
##   disp = col_double(),  
##   hp = col_integer(),  
##   drat = col_double(),  
##   wt = col_double(),  
##   qsec = col_double(),  
##   vs = col_integer(),  
##   am = col_integer(),  
##   gear = col_integer(),  
##   carb = col_integer()  
## )
```

example

```
car_cata <- read_csv(file = "data/mtcars.csv", col_types =  
  cols(  
    mpg = col_double(),  
    cyl = col_integer(),  
    disp = col_double(),  
    hp = col_integer(),  
    drat = col_double(),  
    vs = col_integer(),  
    wt = col_double(),  
    qsec = col_double(),  
    am = col_integer(),  
    gear = col_integer(),  
    carb = col_integer()  
  )  
)
```


Excel file

- to load *.xls* or *.xlsx* files, use **read_excel()** function in the **readxl** package.
 - You can specify the sheet of of an Excel spreadsheet.
 - Read its documentation for more details.
 - Installed as a part of **tidyverse**, but not a core **tidyverse** package

```
# Need to load the package explicitly
```

```
library(readxl)
```

```
# Loading the first sheet of an Excel file
```

```
data <- read_excel("data/datasets.xlsx", 1)
```

Exporting images

ggsave

- call `ggsave()` after calling `ggplot()` to save the last ggplot object:

```
ggplot(mtcars, mapping = aes(x = wt, y = mpg))+  
  geom_point(shape = 1)  
# save PDF  
ggsave("myplot.pdf", width = 8, height = 8, units = "cm", useDingbats = F)  
# save PNG  
ggsave("myplot.png", width = 8, height = 8, units = "cm")
```

create a PDF

- call `ggplot()` functions between `pdf()` and `dev.off()`
- Each plot on a sperate page

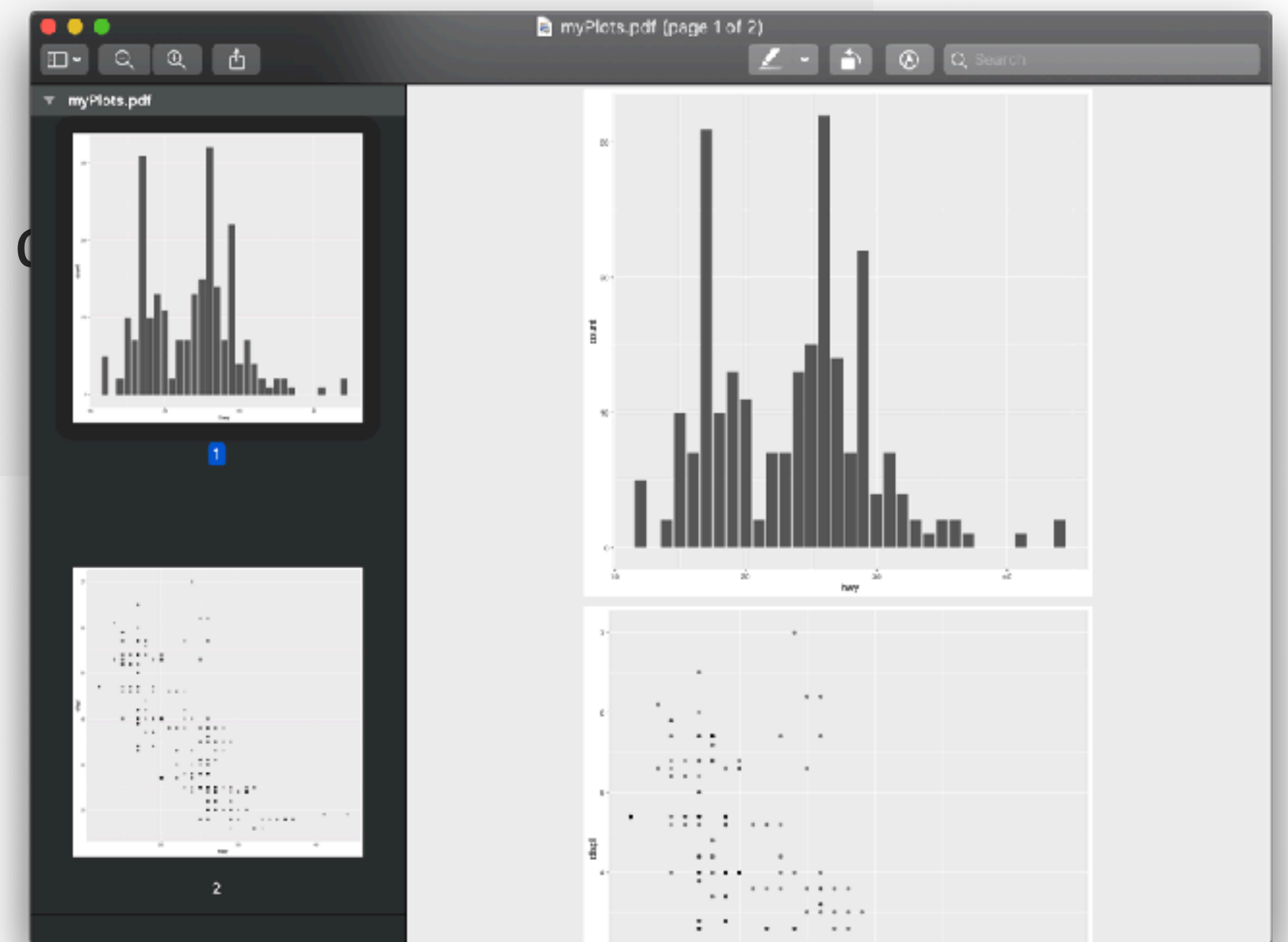
```
pdf("myPlots.pdf", width = 8, height = 8)
```

```
# plot your graphs here
```

```
ggplot(data = mpg) + geom_bar(aes(x = hwy))
```

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

```
dev.off()
```



Axes/graph labels

ggplot2

Labels

ggtitle("Title")	Add a main title above the plot
xlab("x-axis")	Change the label on the x-axis
ylab("y-axis")	Change the label on the y-axis
labs(...)	title = "main title above the plot", subtitle = "subtitle below title", caption = "caption below plot", x = "x-axis", y = "y-axis"
annotate()	annotate(geom="text", x = 10, y = 10, label="A")

Scales

ggplot2

Scales

- Scales map data values to the visual values of an aesthetic.
- To change a mapping, add a new scale:

```
ggplot(data = mpg) +  
  geom_bar(mapping = aes(x = drv, fill = drv)) +  
  scale_fill_manual(  
    values = c("skyblue", "royalblue", "navy"),  
    limits = c("4", "f", "r"),  
    breaks = c("4", "f", "r"),  
    name = "drv",  
    labels = c("4 (4wd)", "f (fwd)", "r (rwd)")  
  )
```

*Title to use in
legend/axis*

aesthetic to adjust

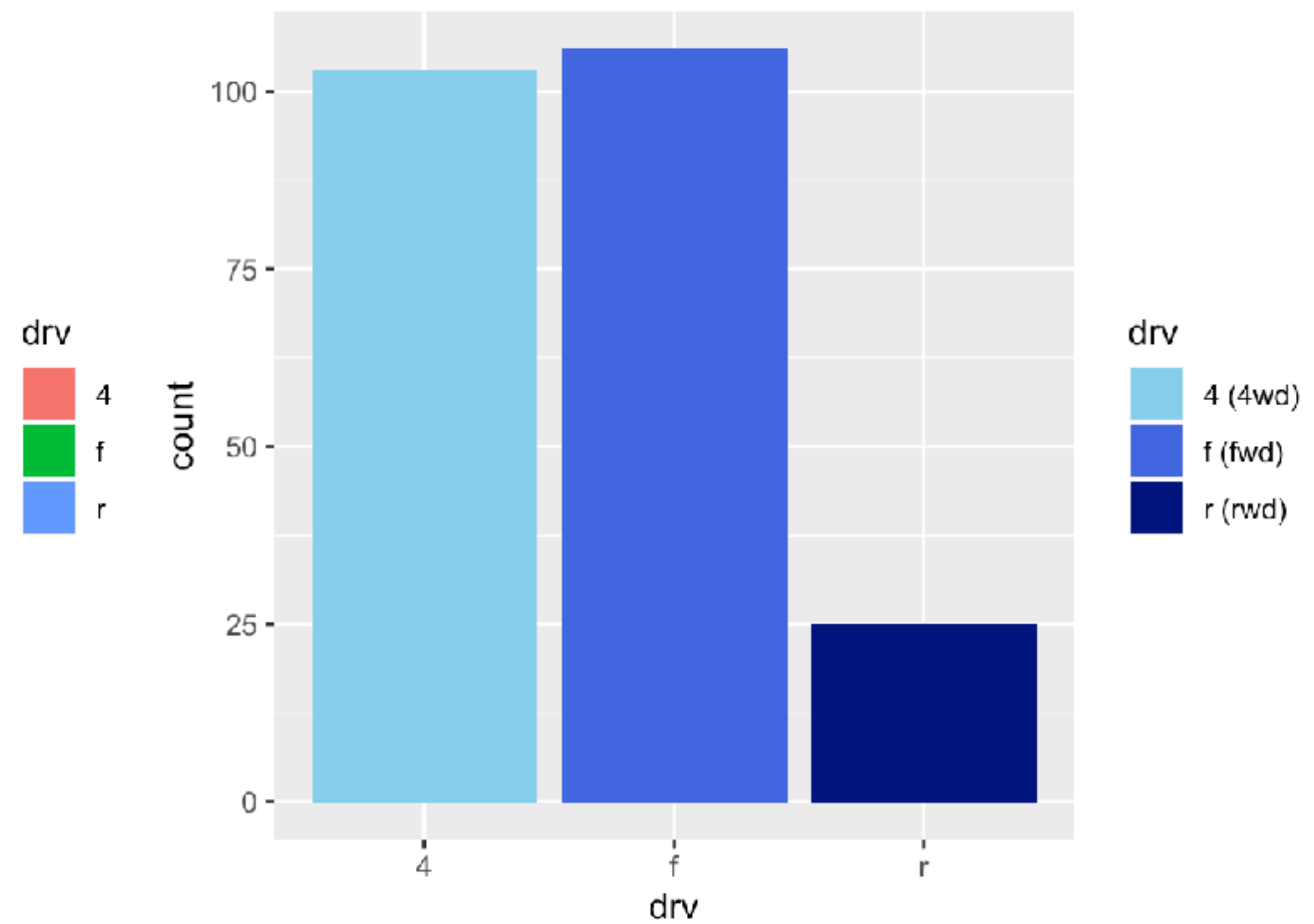
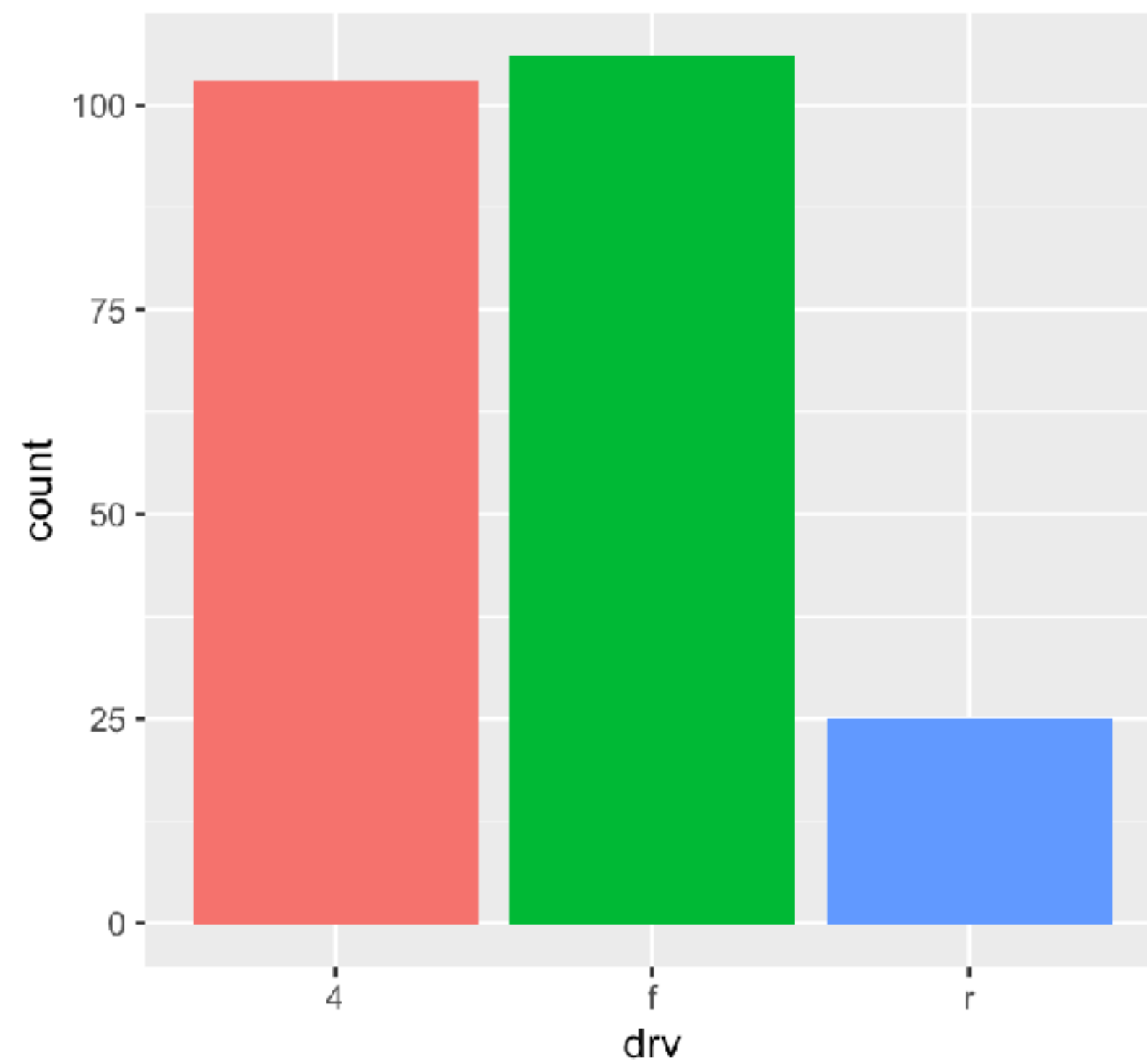
= types of scale

scale-specific arguments

range of values to include in mapping

breaks to use in legend/axis

labels to use in legend/axis



```
p1 <- ggplot(data = mpg)+
  geom_bar( mapping = aes(x = drv, fill = drv))

p2 <- ggplot(data = mpg)+
  geom_bar( mapping = aes(x = drv, fill = drv)) +
  scale_fill_manual(
    values = c("skyblue", "royalblue", "navy"),
    limits = c("4", "f", "r"),
    breaks = c("4", "f", "r"),
    name = "drv",
    labels = c("4 (4wd)", "f (fwd)", "r (rwd)")
  )
grid.arrange(p1, p2, ncol = 2)
```

General purpose scales

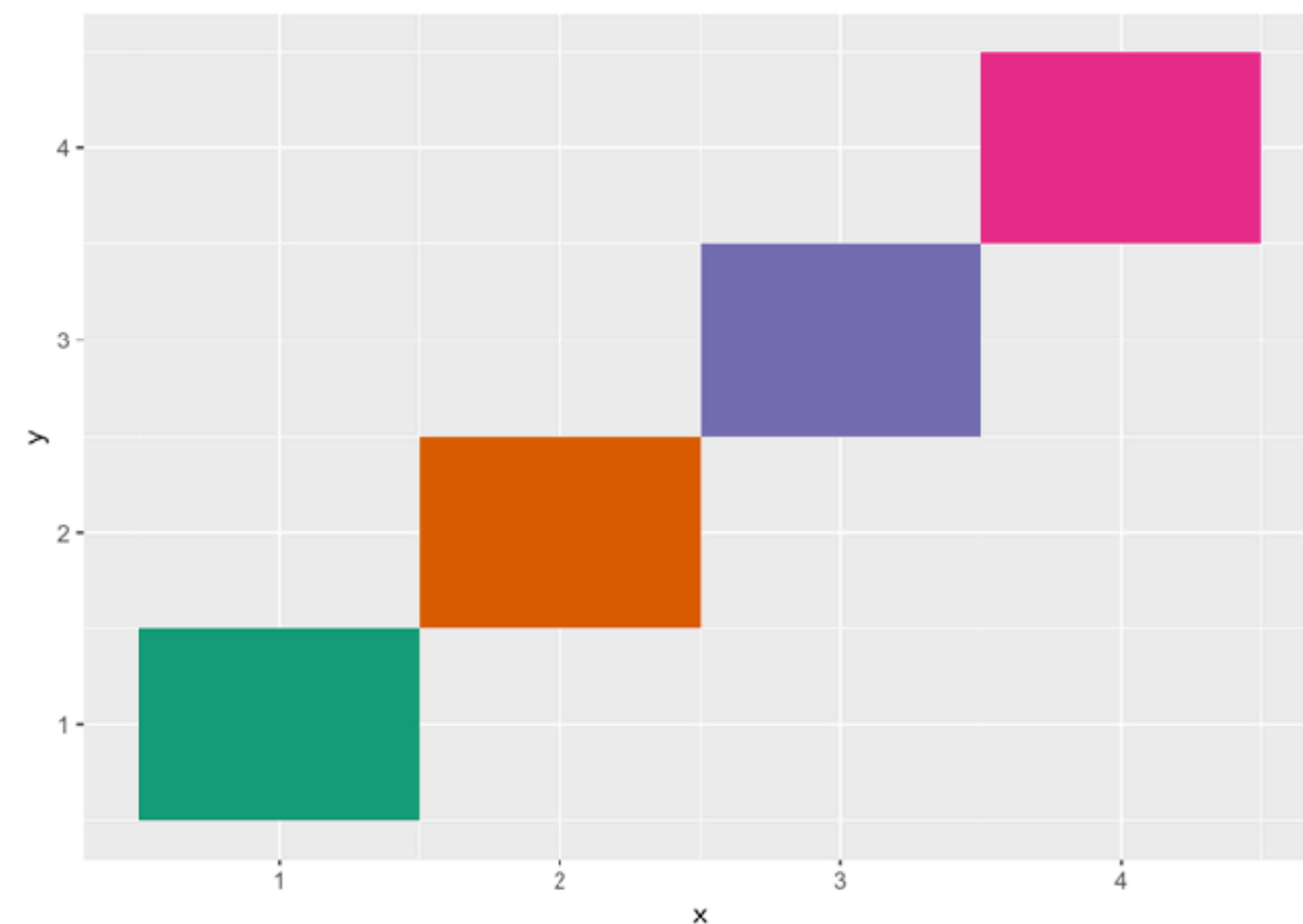
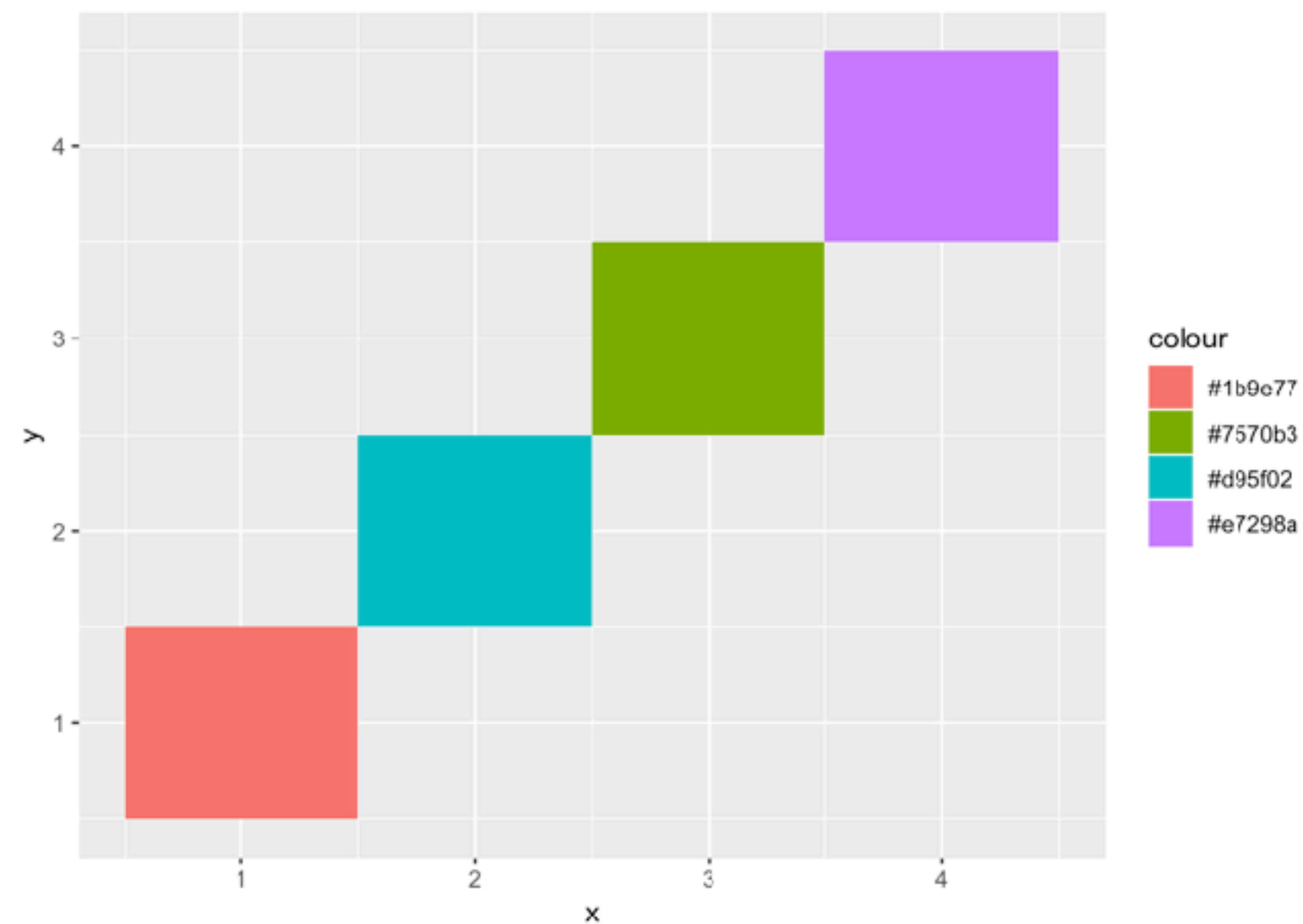
<code>scale*_continuous()</code>	map continuous values to visual scale
<code>scale*_discrete()</code>	map discrete values to visual scale
<code>scale*_identity()</code>	use data values as visual ones
<code>scale*_manual(values = c())</code>	map discrete values to manually chosen visual ones
<code>scale*_date(date_labels = “%m/%d”, date_breaks = “2 weeks”)</code>	treat data values as dates (class = Date)
<code>scale*_datetime()</code>	treat data x values as date times. (class = POSIXct)

*** = alpha, color, fill, linetype, shape, size, x, y**

scale_*_identity()

```
df <- data.frame(  
  x = 1:4,  
  y = 1:4,  
  colour = c('#1b9e77', '#d95f02', '#7570b3', '#e7298a')  
)  
ggplot(df, aes(x, y)) + geom_tile(aes(fill = colour))  
  
ggplot(df, aes(x, y)) + geom_tile(aes(fill = colour)) +  
  scale_fill_identity()
```

	x	y	colour
1	1	1	#1b9e77
2	2	2	#d95f02
3	3	3	#7570b3
4	4	4	#e7298a

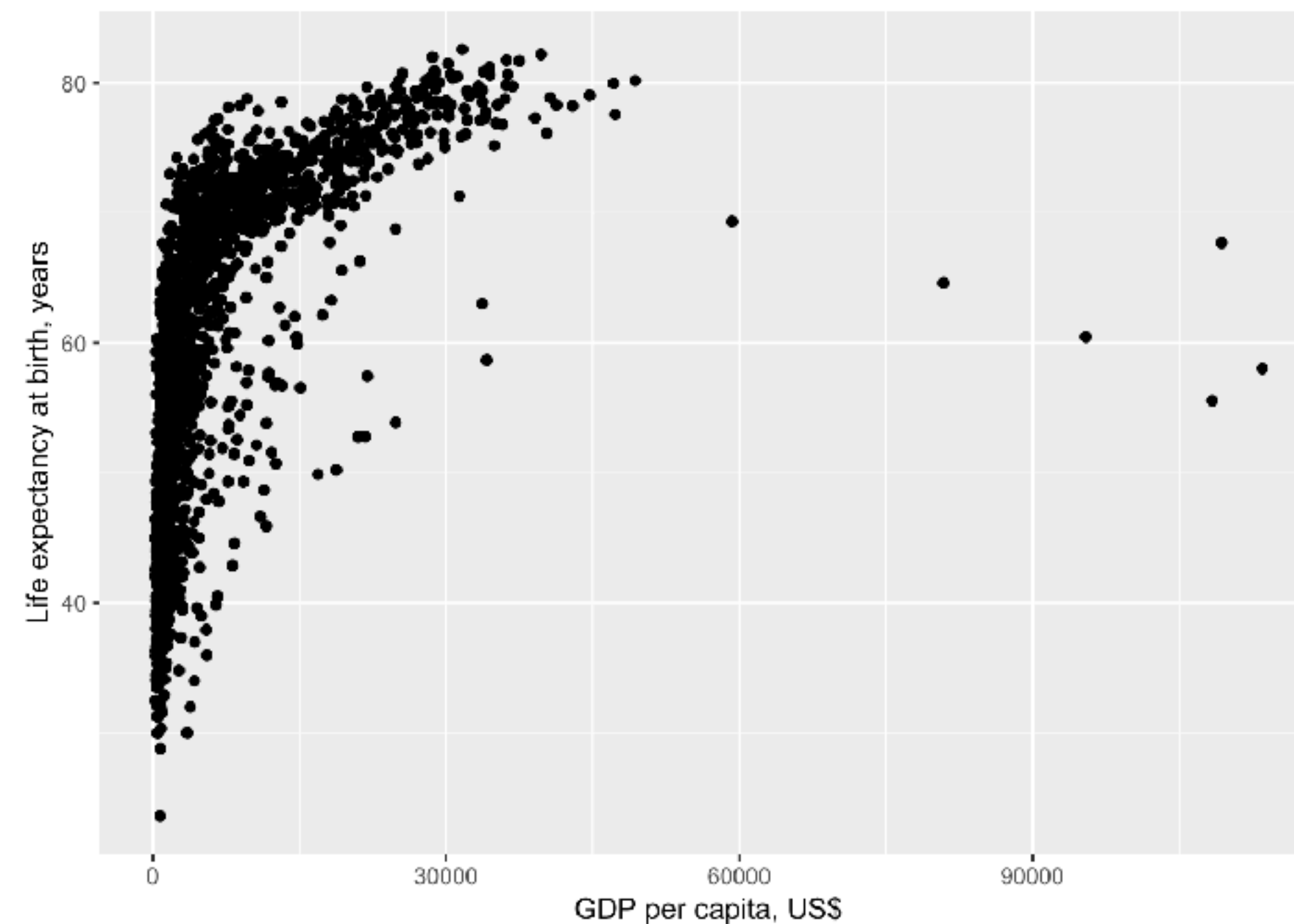


X and Y location scales

scale_x_log10()	plot x on log10 scale
scale_x_reverse()	reverse direction of x-axis
scale_x_sqrt()	plot x on square root scale

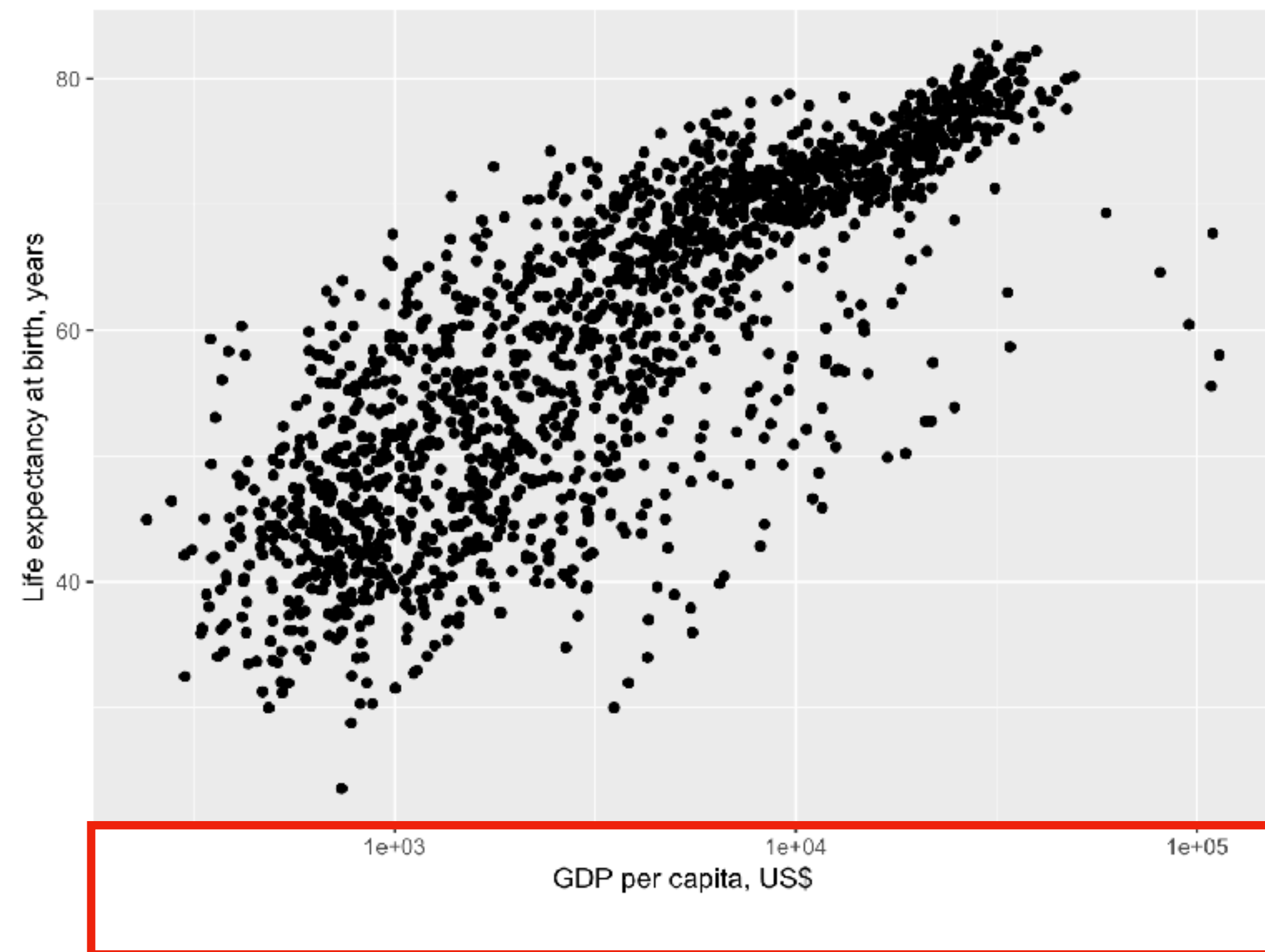
gapminder

```
# install.packages("gapminder")  
library(gapminder)  
  
# Looking up ?gapminder to see the documentation  
  
# scatter plot of GDP per capita and Life expectancy  
ggplot(data = gapminder, mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point() +  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```



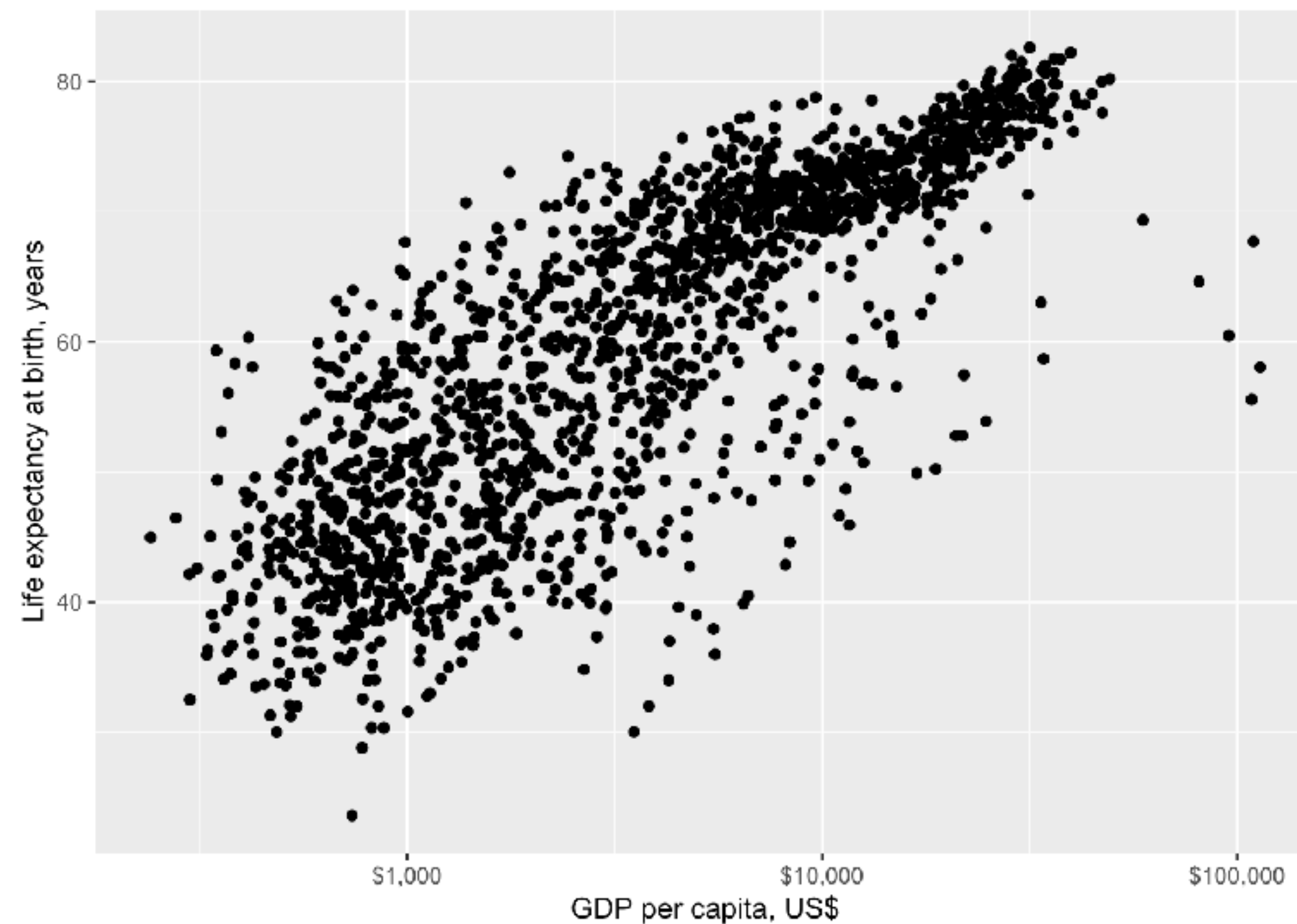
gapminder

```
ggplot(data = gapminder, mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point() +  
  scale_x_log10() +  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```



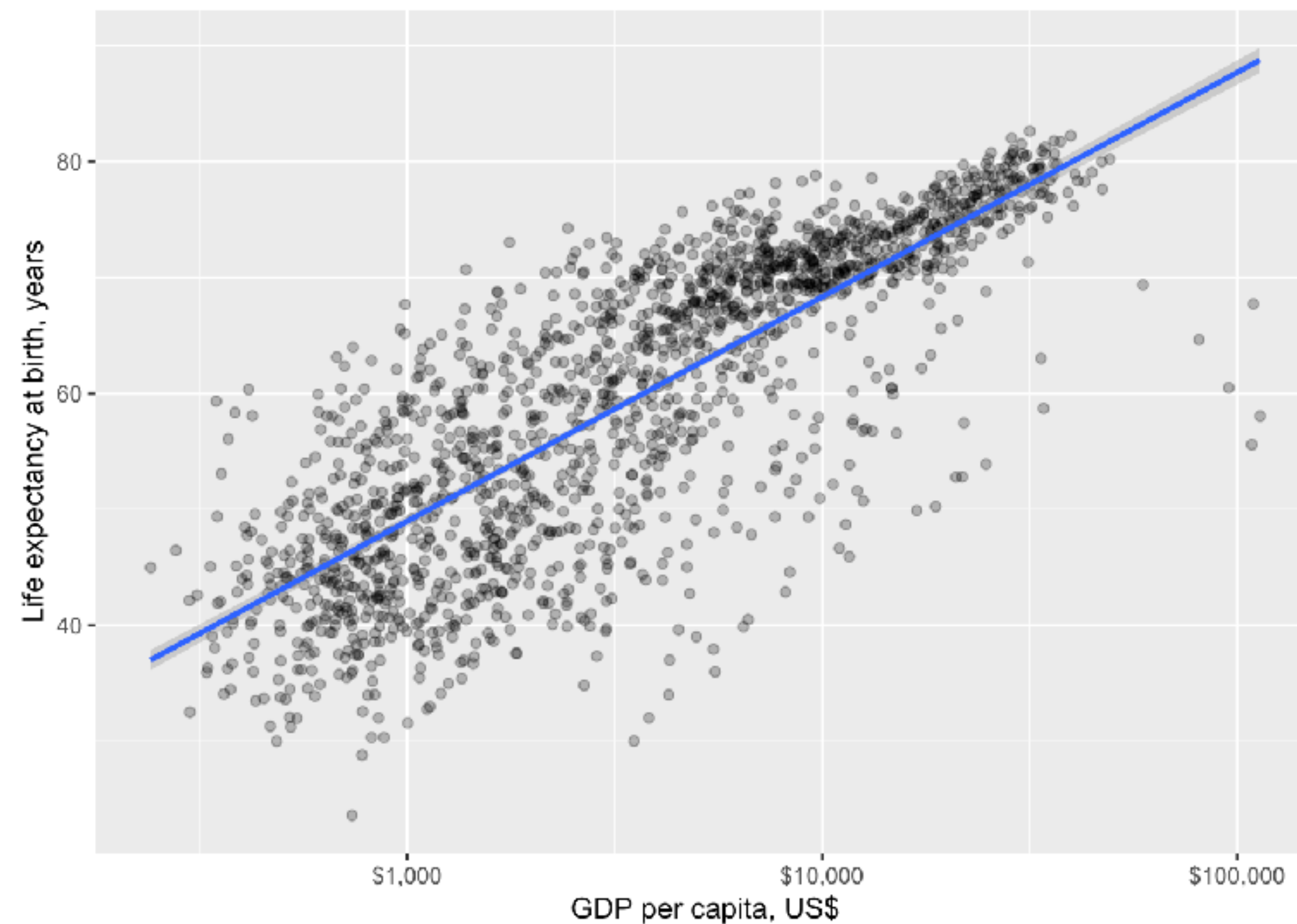
gapminder

```
# handy function scale::dollar, scale::comma by passing a function  
ggplot(data = gapminder, mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point() +  
  scale_x_log10(labels = scales::dollar) +  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```



gapminder

```
ggplot(data = gapminder, mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point(alpha = 0.3) +  
  geom_smooth(method = "gam") +  
  scale_x_log10(labels = scales::dollar) +  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```



Colour scales

ggplot2

Colour scale

- Is your data **discrete** or **continuous**?
- Colourblind-friendly palette
 - Package: **viridis** (comes with ggplot2 v3.0.0)
 - <https://cran.r-project.org/web/packages/viridis/vignettes/intro-to-viridis.html>
 - Package: **RColorBrewer**
 - `display.brewer.all()`
 - `display.brewer.pal(8, "Blues")`

viridis

viridis



magma



plasma



inferno



cividis



Comparisons

Base R: rainbow.colours

rainbow



Base R: heat.colours

heat



ggplot2 default

ggplot default



colorbrewer blues

brewer blues



*colorbrewer
yellow-green-blue*

brewer yellow-green-blue



viridis



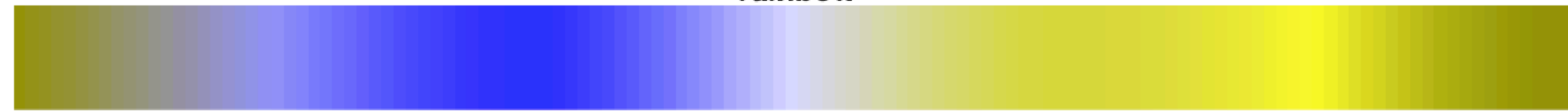
magma



Green-blind (Deuteranopia)

Comparisons

rainbow



heat



ggplot default



brewer blues



brewer yellow-green-blue



viridis



magma



Red-blind (Protanopia)

Comparisons

rainbow



heat



ggplot default



brewer blues



brewer yellow-green-blue



viridis



magma



Blue-blind (Tritanopia)

Comparisons

rainbow



heat



ggplot default



brewer blues



brewer yellow-green-blue



viridis



magma



Desaturated

Comparisons

rainbow



heat



ggplot default



brewer blues



brewer yellow-green-blue



viridis



magma

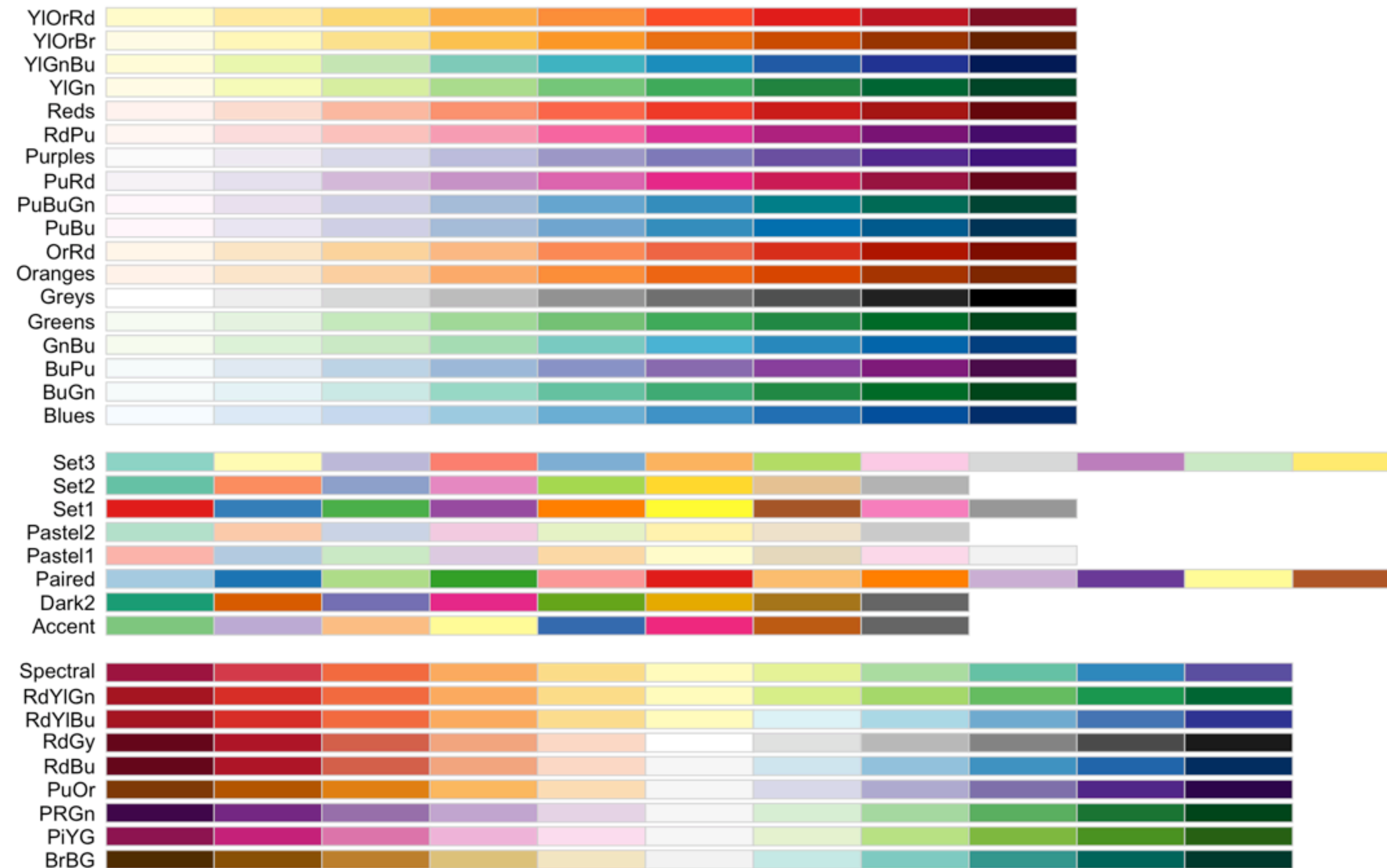


Fill & colour scale - Discrete

scale*_discrete()	Colours evenly spaced around the colour wheel (same as hue)
scale*_hue()	Colours evenly spread around the colour wheel
scale*_grey()	Greyscale palette
scale*_viridis_d()	Viridis palettes
scale*_brewer()	ColorBrewer palettes
scale*_manual()	Manually specified colours

ColorBrewer

```
display.brewer.all()
```



Fill & colour scale - Continuous

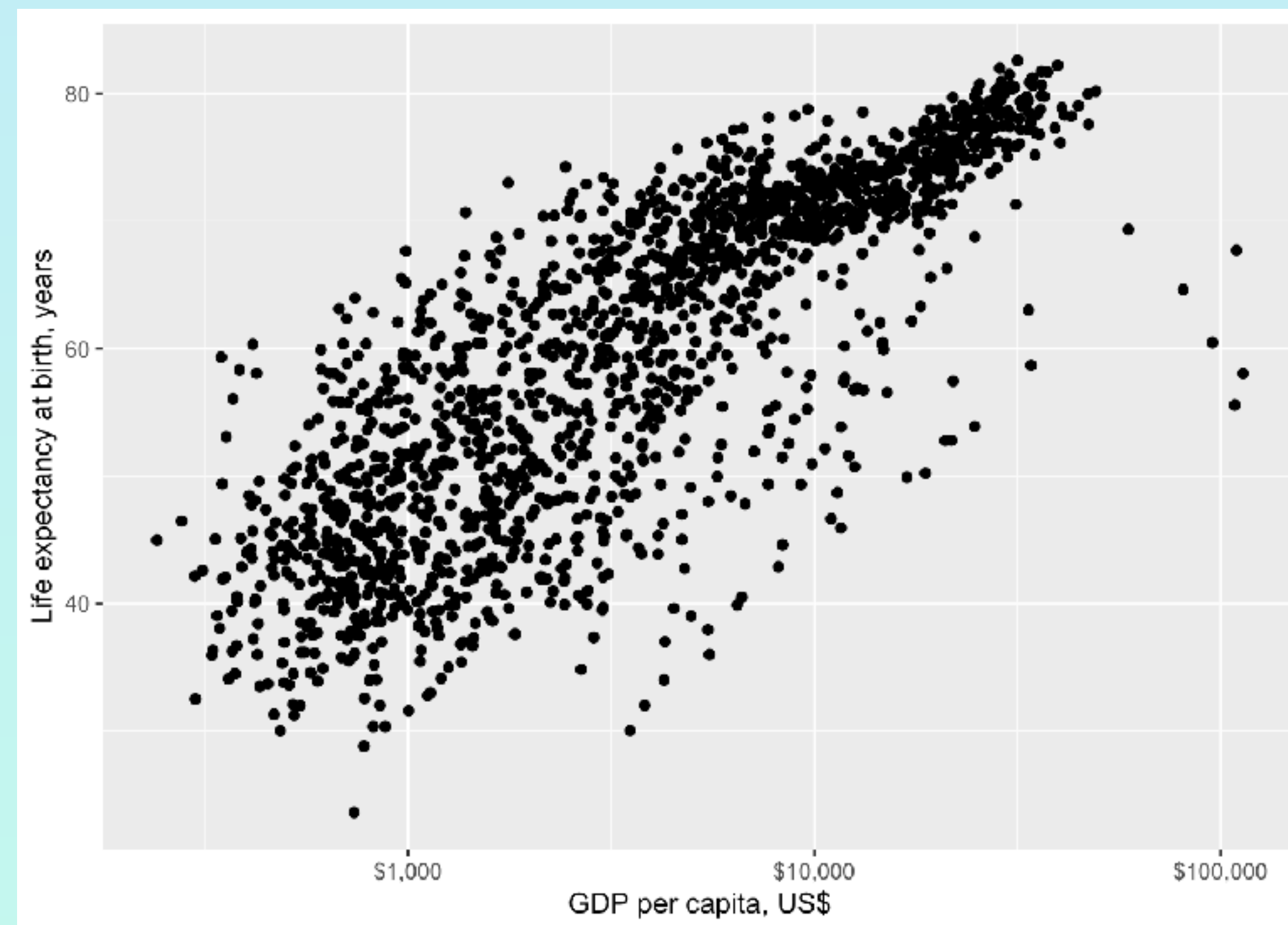
scale*_gradient()	Two-colour gradient
scale*_gradient2()	Gradient with a middle colour and tow colours that diverge from it
scale*_gradientn()	Gradient with n colour, equally spaced
scale*_viridis_c()	Viridis palettes

- Do you really need a continuous scale?
 - Can you bin the data to use discrete scale?

Your turn!

- 1) Modify the code to map **pop** variable (population) to colour aesthetics.
- 2) Did you map colour using a linear scale or log scale?

```
ggplot(data = gapminder, mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point() +  
  scale_x_log10(labels = scales::dollar) +  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```



Solution

```
ggplot(data = gapminder %>% arrange(pop), mapping = aes(x = gdpPercap, y = lifeExp)) +  
  geom_point(mapping = aes(colour = log(pop))) +  
  geom_smooth(method = "gam") +  
  scale_x_log10(labels = scales::dollar) +  
  scale_color_viridis_c(direction = -1)+  
  labs(x = "GDP per capita, US$", y = "Life expectancy at birth, years" )
```

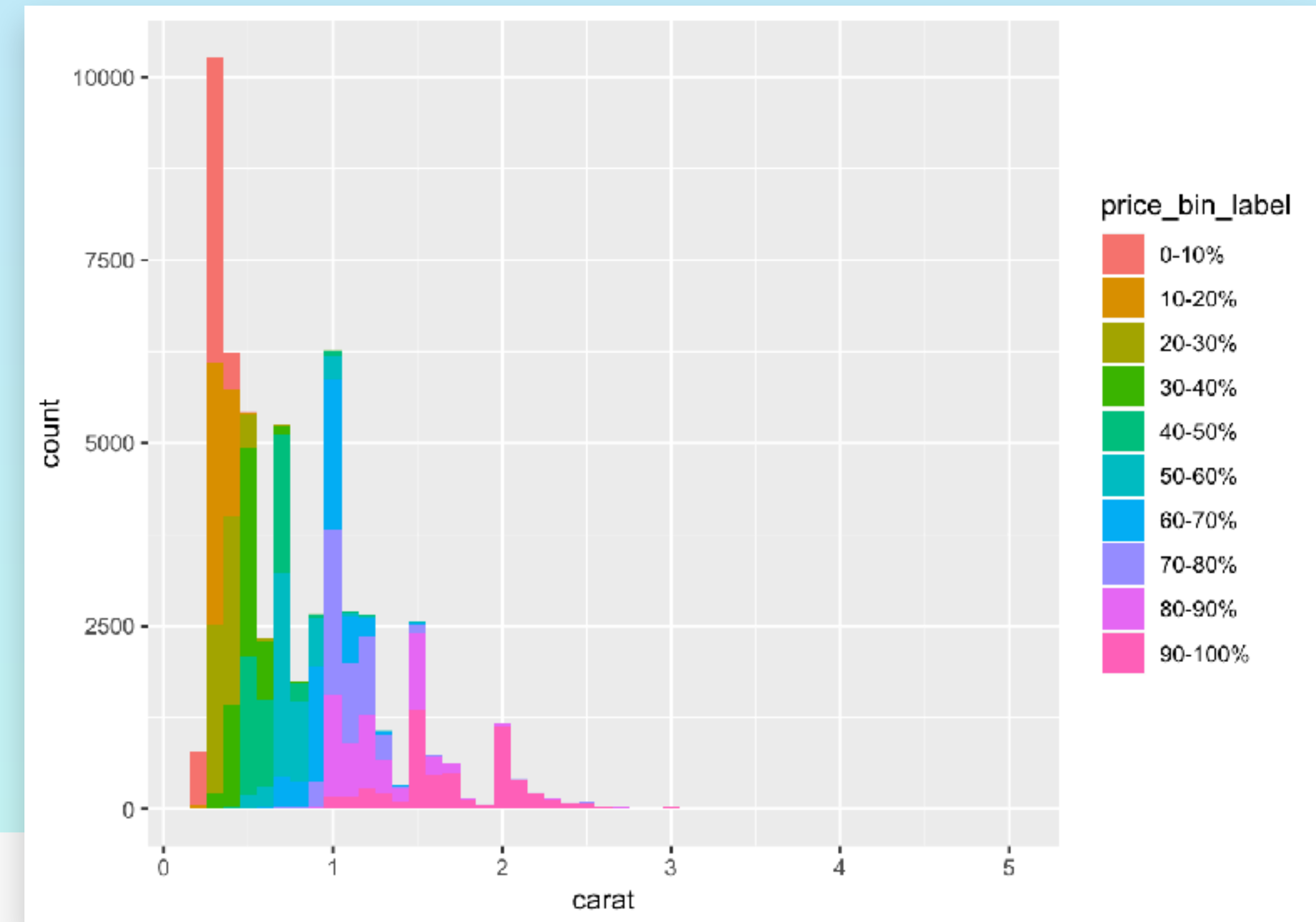


Your turn!

- 1) Apply `scale*_viridis_d()` to the following code:
- 2) Can you use other viridis colour palettes?

```
# create a new variable by binning by price
df <- diamonds %>%
  mutate(price_bin = ntile(price, 10),
         price_bin_label = factor(price_bin, levels = c(1:10),
                                labels = c("0-10%", "10-20%", "20-30%", "30-40%", "40-50%",
                                           "50-60%", "60-70%", "70-80%", "80-90%", "90-100%")))

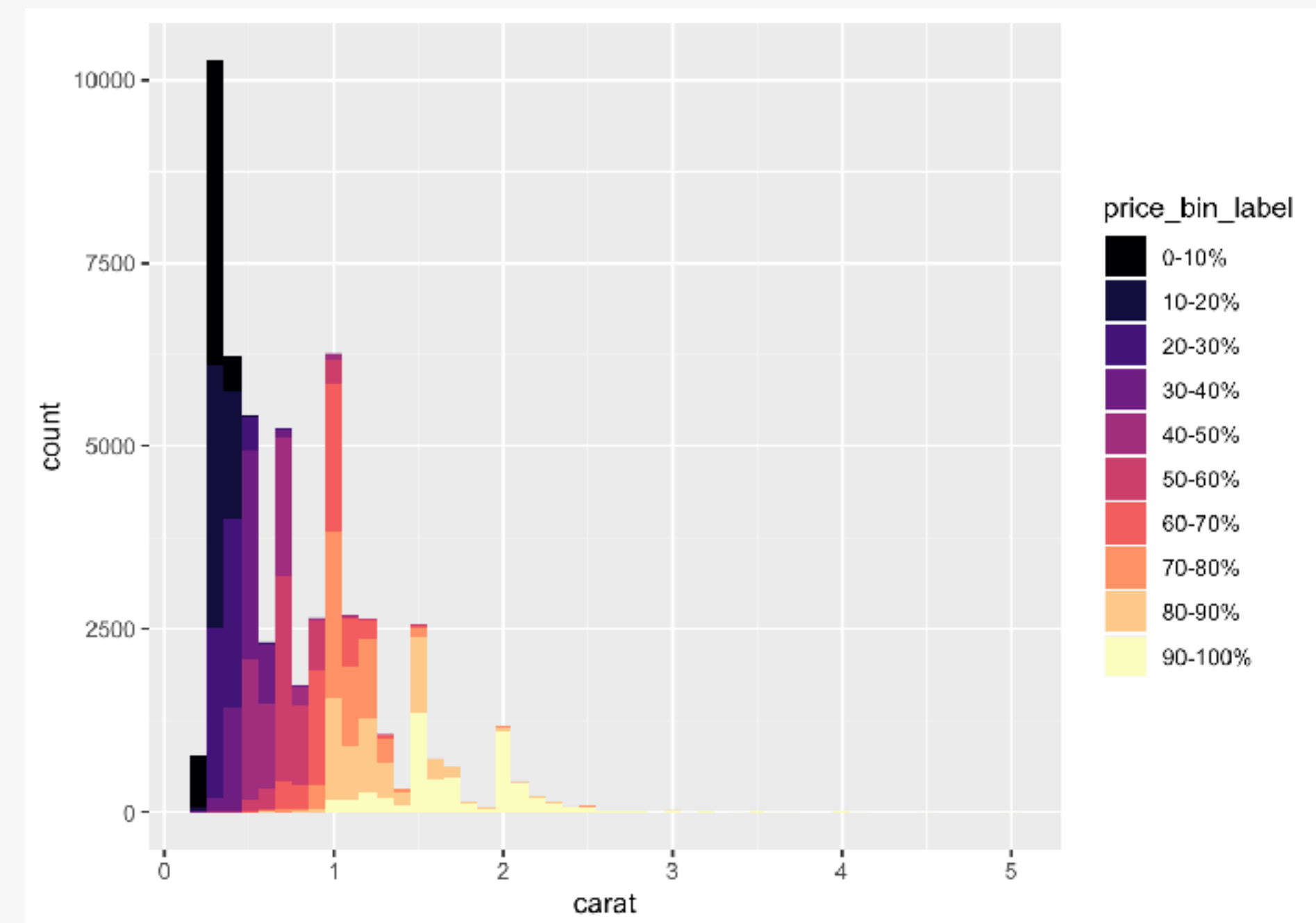
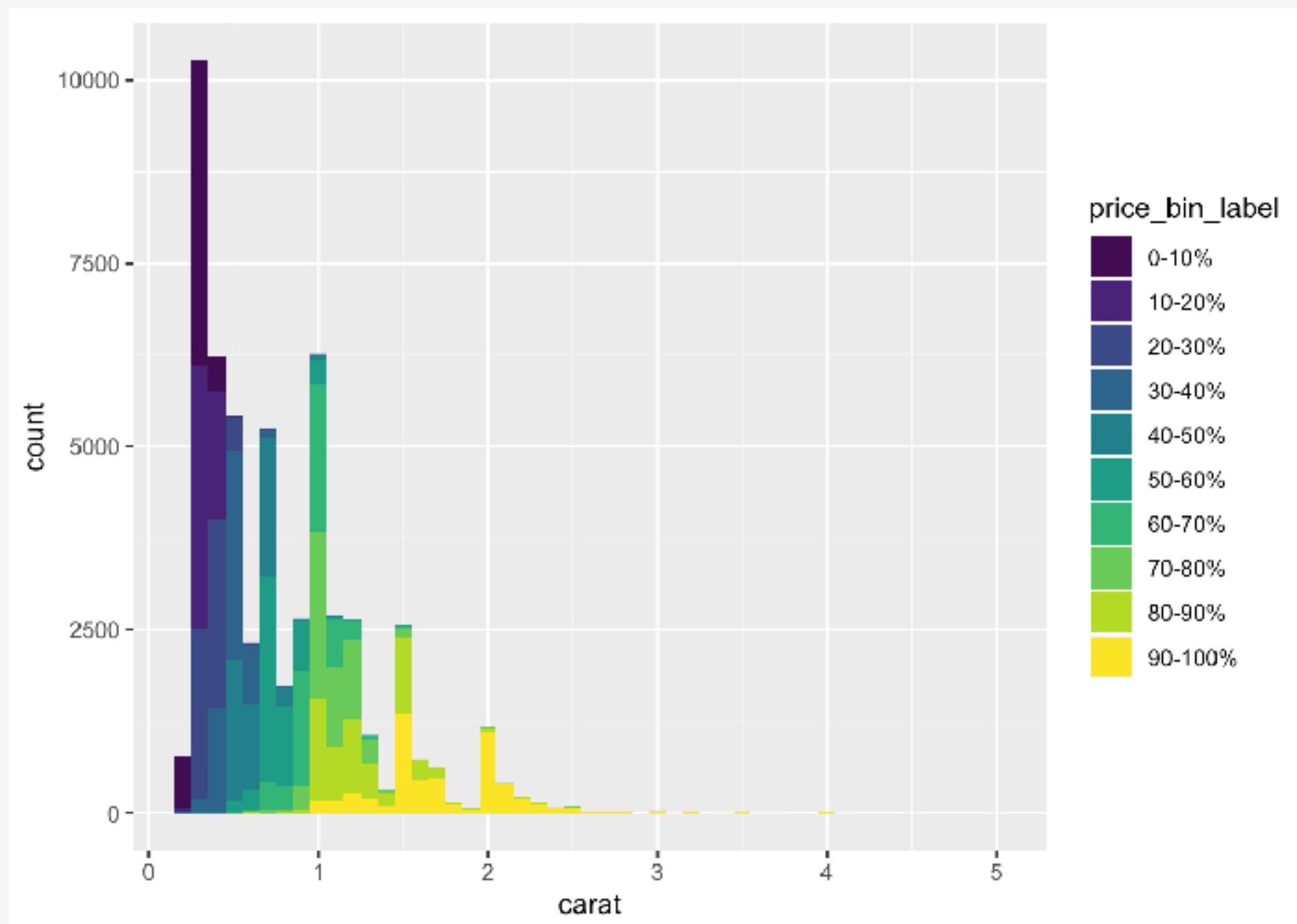
ggplot(data = df) +
  geom_histogram(mapping = aes(x = carat, fill = price_bin_label), binwidth = 0.1)
```



Solution

```
ggplot(data = df) +  
  geom_histogram(mapping = aes(x = carat, fill = price_bin_label), binwidth = 0.1)+  
  scale_fill_viridis_d()
```

```
ggplot(data = df) +  
  geom_histogram(mapping = aes(x = carat, fill = price_bin_label), binwidth = 0.1)+  
  scale_fill_viridis_d(option="magma")
```



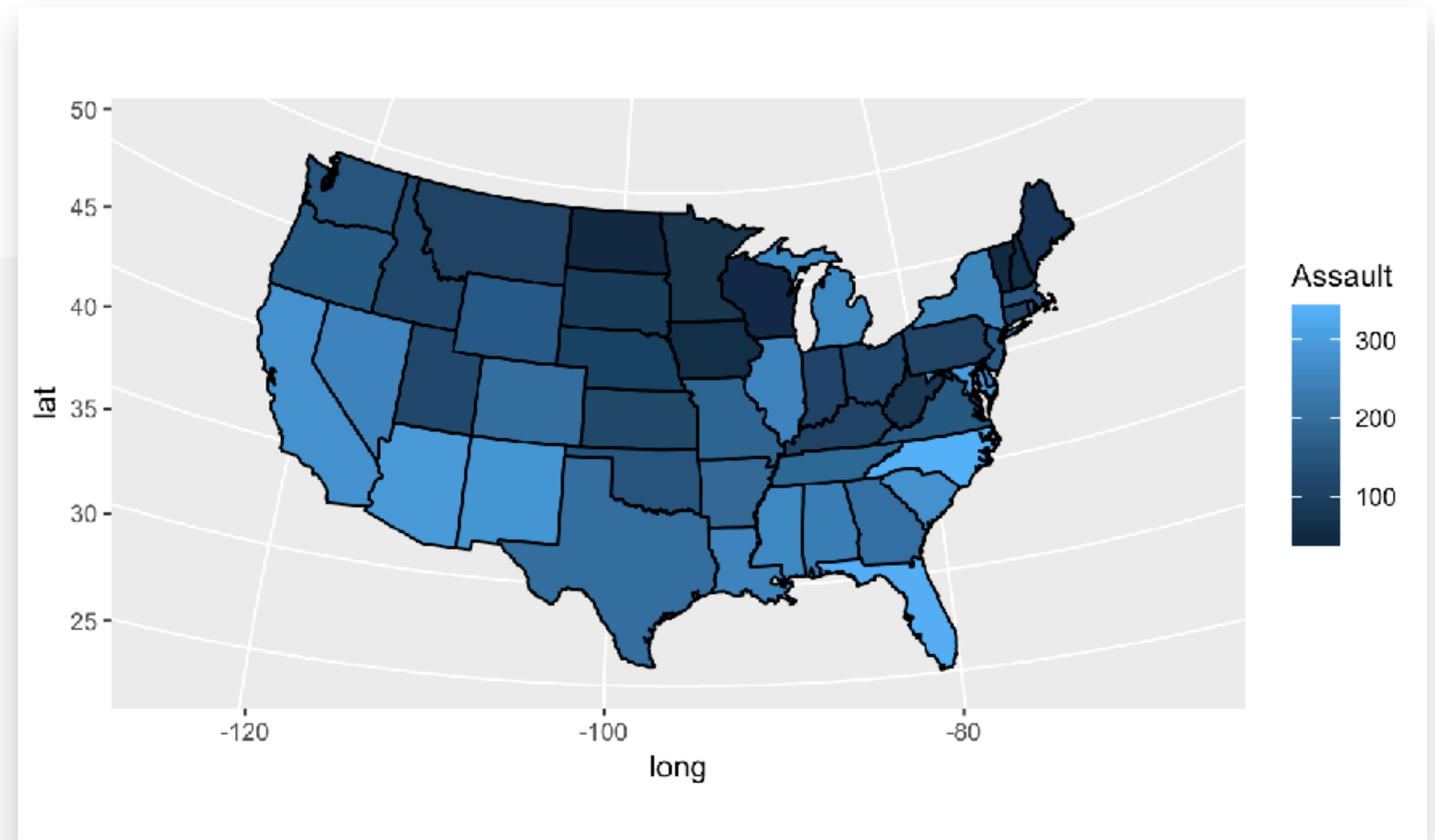
Choropleth

Choropleth Map

```
# using USArrests datasets
crimes <- data.frame(USArrests) %>%
  mutate(state = tolower(rownames(USArrests)))

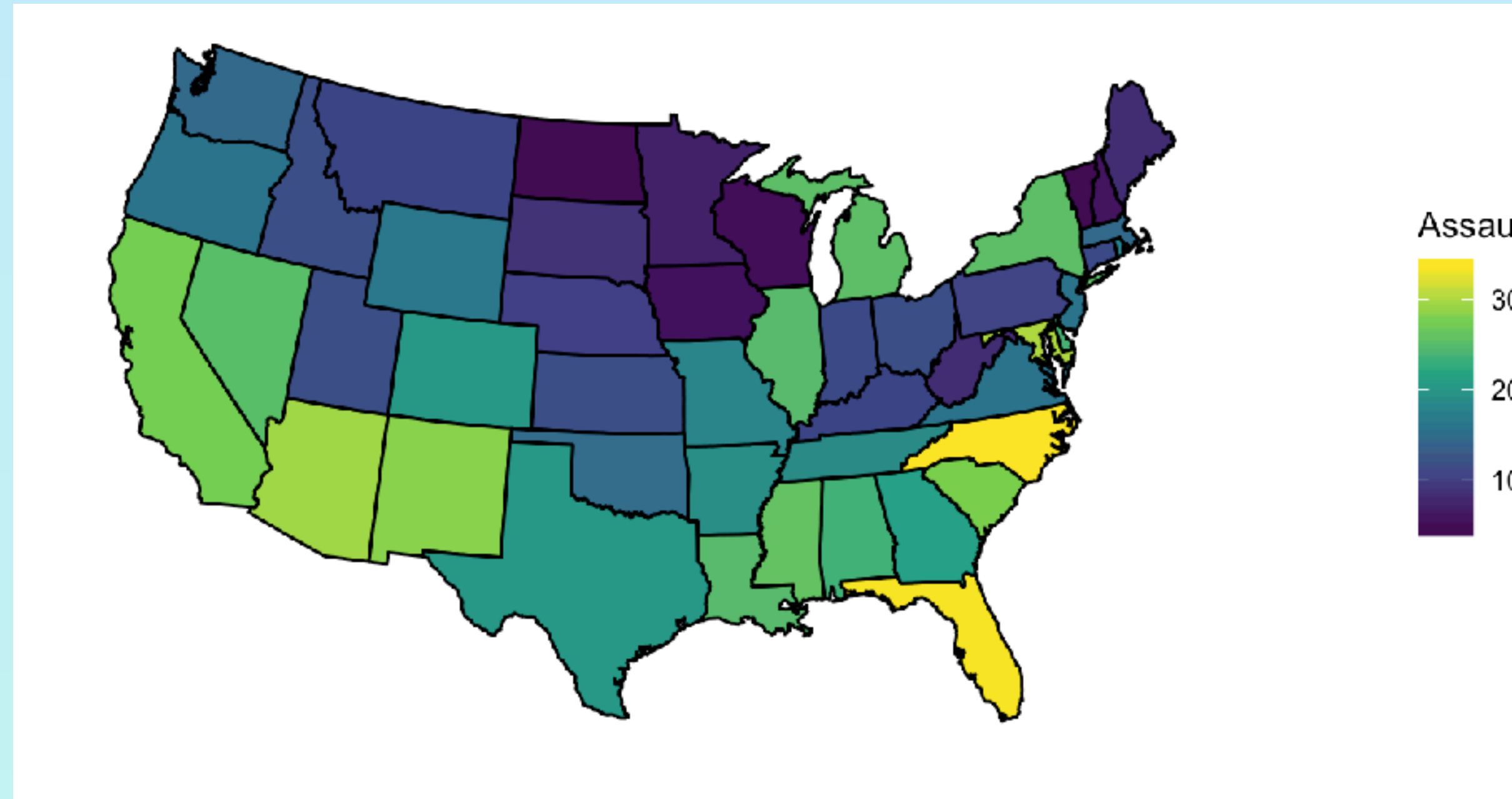
# using maps library
library(maps)
library(ggplot2)
states_map <- map_data("state")
# merge datasets together
crime_map <- states_map %>%
  left_join(crimes, by = c("region" = "state")) %>%
  arrange(group, order)

ggplot(data = crime_map, aes(x=long, y=lat, group = group, fill = Assault))+
  geom_polygon(colour = "black")+
  coord_map("polyconic")
```

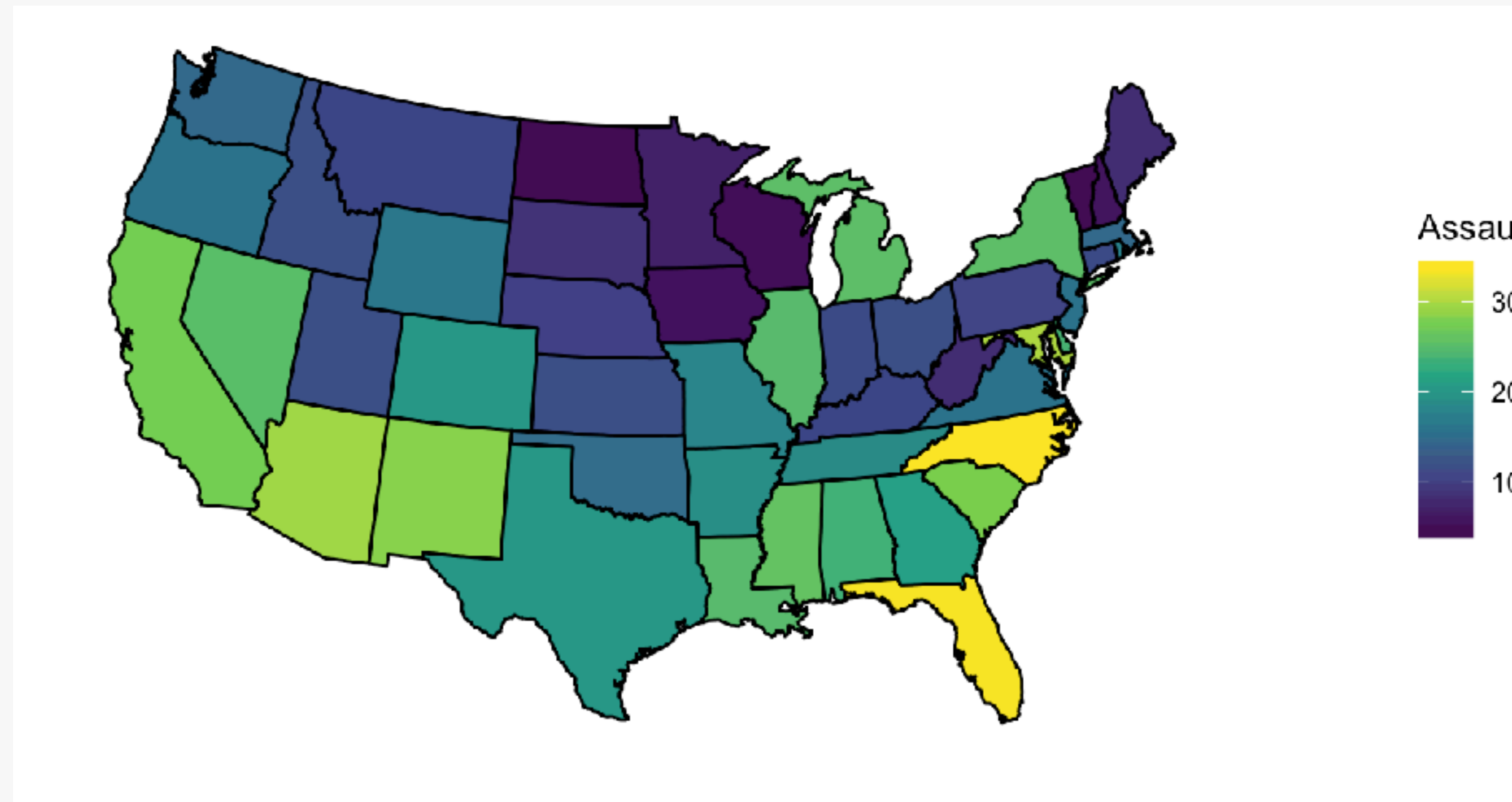


Your turn!

1) Apply `scale_fill_viridis_c()`



Solution

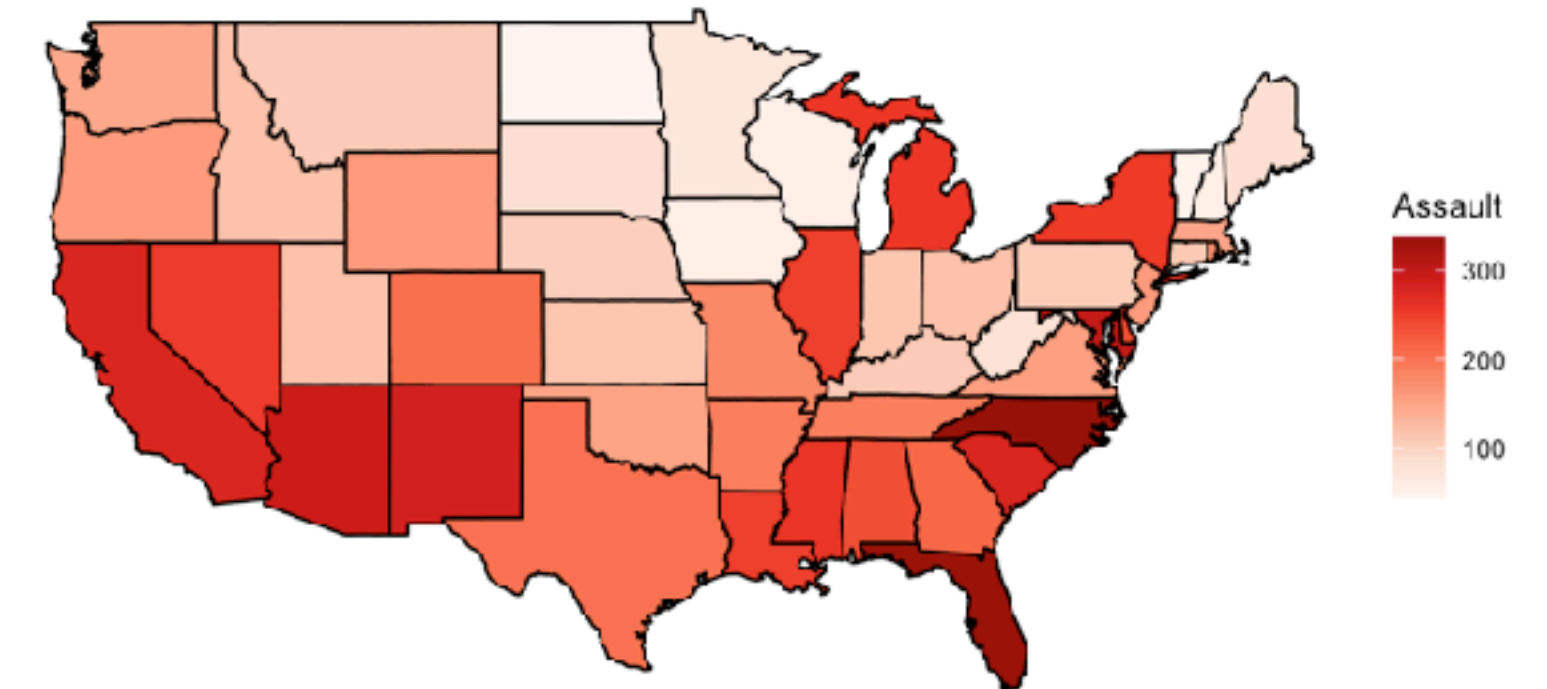


```
ggplot(data = crime_map, aes(x=long, y=lat, group = group, fill = Assault))+  
  geom_polygon(colour = "black")+  
  coord_map("polyconic")+  
  scale_fill_viridis_c()+  
  theme_void()
```

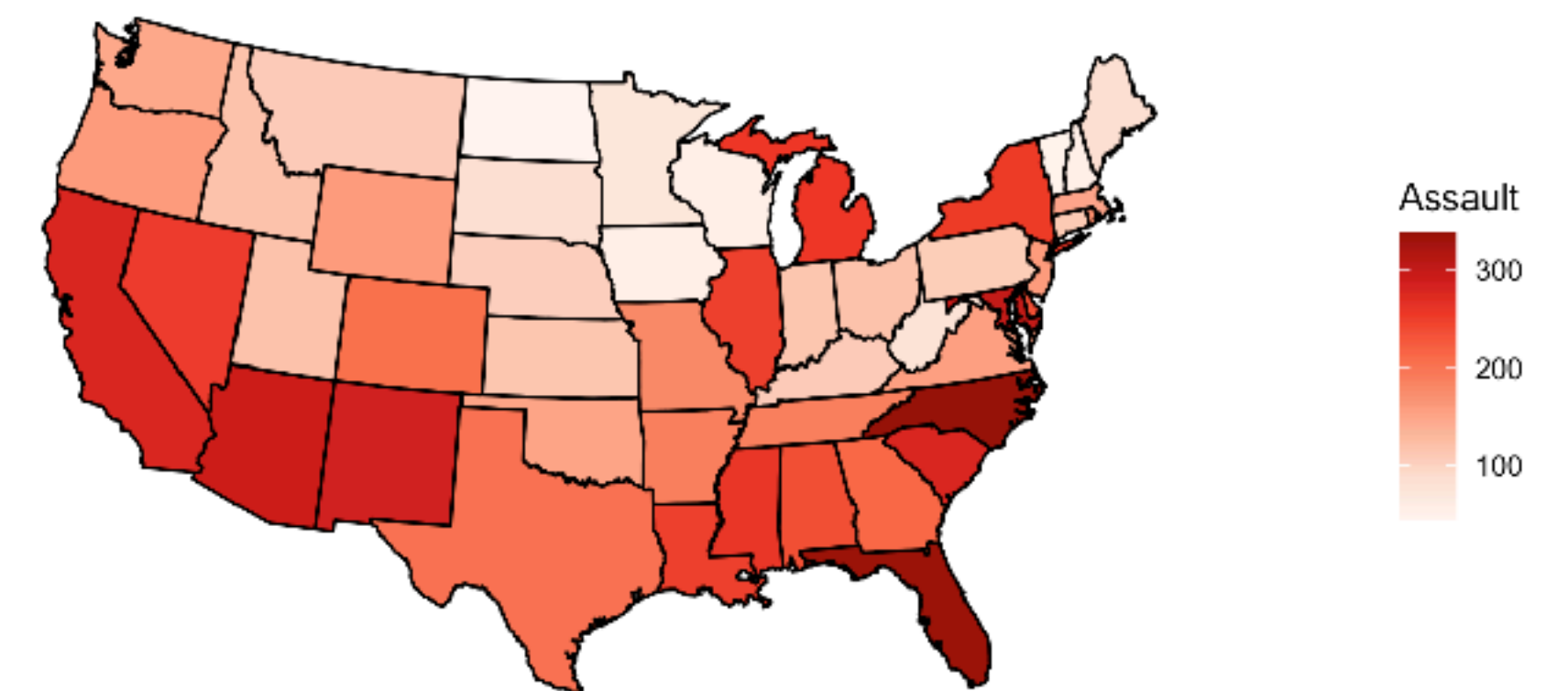
Choice of Projections

```
library(RColorBrewer)

ggplot(data = crime_map, aes(x=long, y=lat, group = group,
fill = Assault))+
  geom_polygon(colour = "black")+
  scale_fill_gradientn(colors = brewer.pal(8,"Reds")) +
  coord_map(projection = "mercator")+
  theme_void()
```



```
ggplot(data = crime_map, aes(x=long, y=lat, group = group,
fill = Assault))+
  geom_polygon(colour = "black")+
  scale_fill_gradientn(colors = brewer.pal(8,"Reds")) +
  coord_map("albers", lat0=30, lat1=40) +
  theme_void()
```



Map of EU countries

```
# Map of Europe  
# Get map data for world  
world_map <- map_data("world")  
  
# Check $regions  
sort(unique(world_map$region))  
  
##      [1] "Afghanistan"  
##      [2] "Albania"  
##      [3] "Algeria"  
##      [4] "American Samoa"  
##      [5] "Andorra"  
##      [6] "Angola"  
##      [7] "Anguilla"  
##      [8] "Antarctica"  
##      [9] "Antigua"  
##     [10] "Argentina"
```

Map of EU countries

```
# A vector of countries
europeanUnion <- c("Austria", "Belgium", "Bulgaria", "Croatia", "Cyprus",
                  "Czech Republic", "Denmark", "Estonia", "Finland", "France",
                  "Germany", "Greece", "Hungary", "Ireland", "Italy", "Latvia",
                  "Lithuania", "Luxembourg", "Malta", "Netherlands", "Poland",
                  "Portugal", "Romania", "Slovakia", "Slovenia", "Spain",
                  "Sweden")

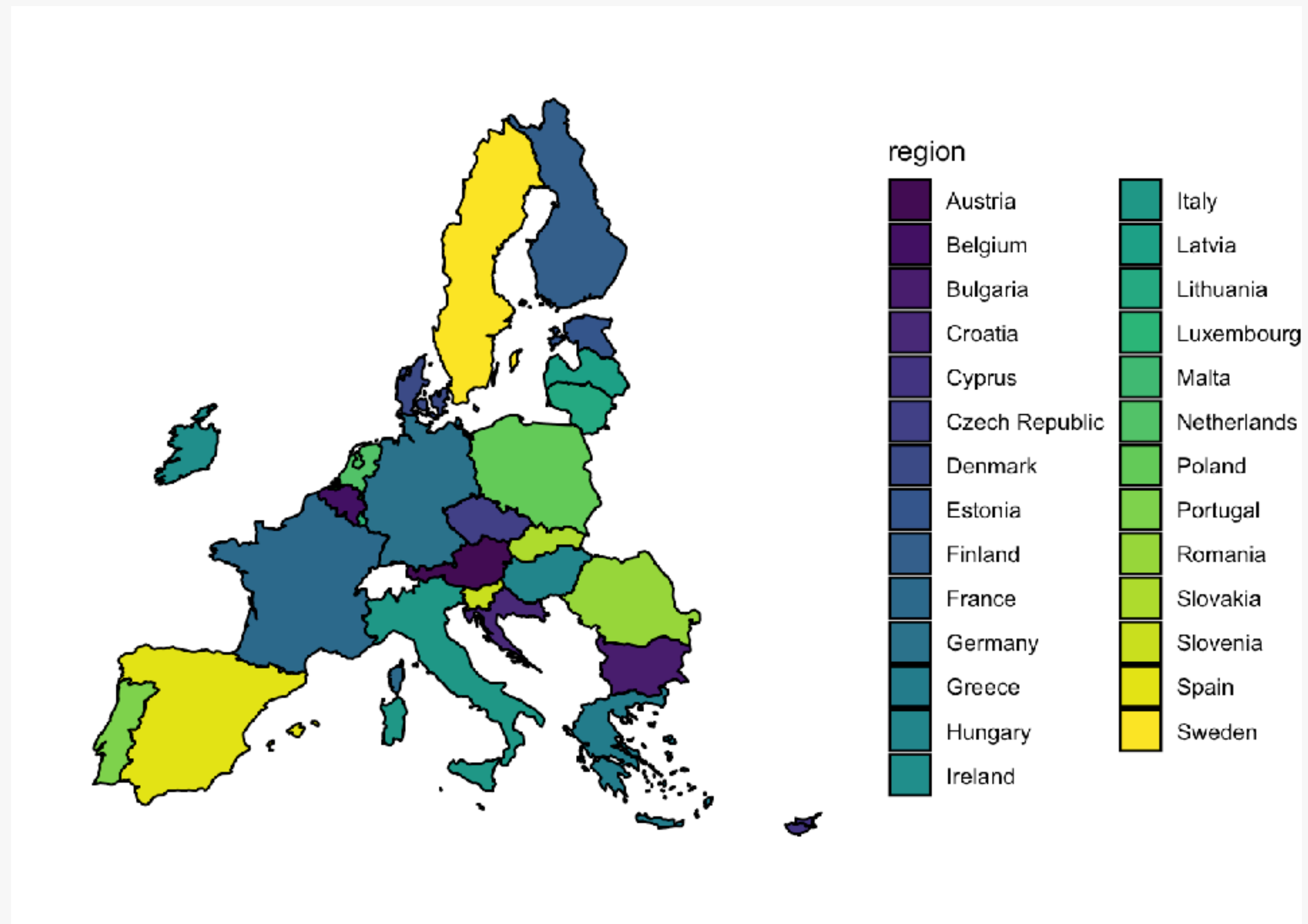
# "UK"
# get the map data
eu_map <- map_data("world", region = europeanUnion)
```

Your turn!

- 1) Create a map of EU countries with `eu_map` data.

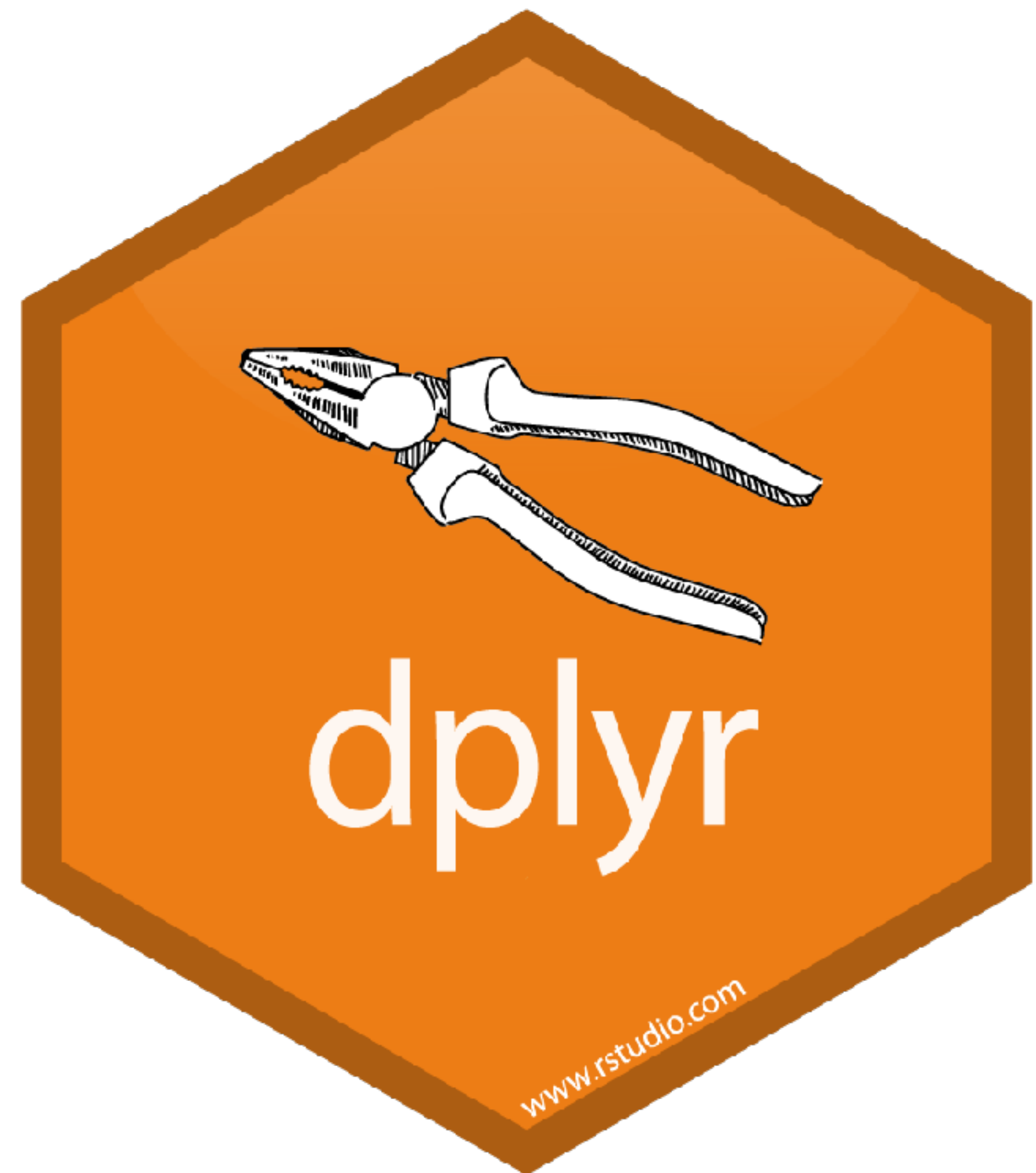
Solution

```
ggplot(eu_map, mapping = aes(x = long, y = lat, group = group, fill = region)) +  
  geom_polygon(colour = "black")+  
  scale_fill_viridis_d() +  
  coord_map("sinusoidal")+  
  theme_void()
```

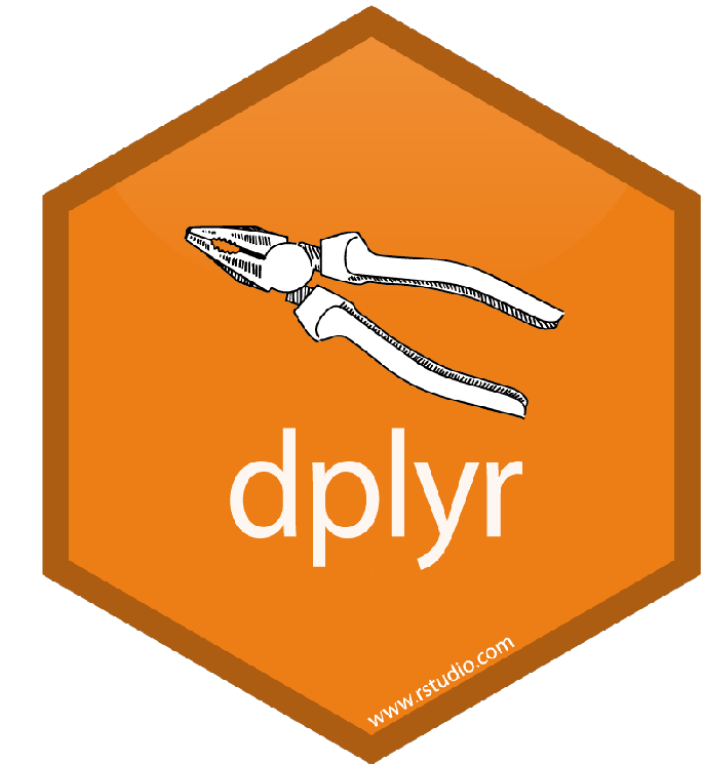


Lecture 8 - Summary

- Data transformation
 - `dplyr` package
- Tidy data
- Importing data
- Exporting images
- Labels
- Scales
- Choropleth map



Recap



- Data transformation with dplyr

Function	Description	Equivalent SQL
<code>select()</code>	selecting columns	SELECT
<code>filter()</code>	filtering rows / subsetting	WHERE
<code>group_by()</code>	grouping data	GROUP BY
<code>summarise()</code>	summarising / aggregating data	-
<code>arrange()</code>	sorting data	ORDER BY
<code>join()</code>	joining data tables	JOIN
<code>mutate()</code>	creating new columns	COLUMN ALIAS

Recap - tidy data

- Each variable must have its own column
- Each observation must have its own row
- Each value must have its own cell

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	1280425583

variables

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	1280425583

observations

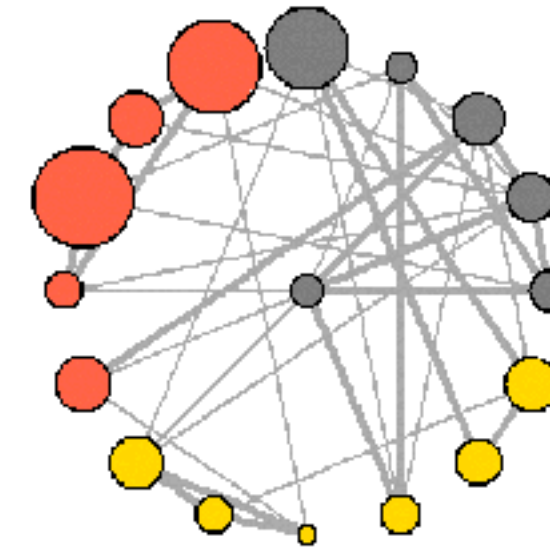
country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280425583

values

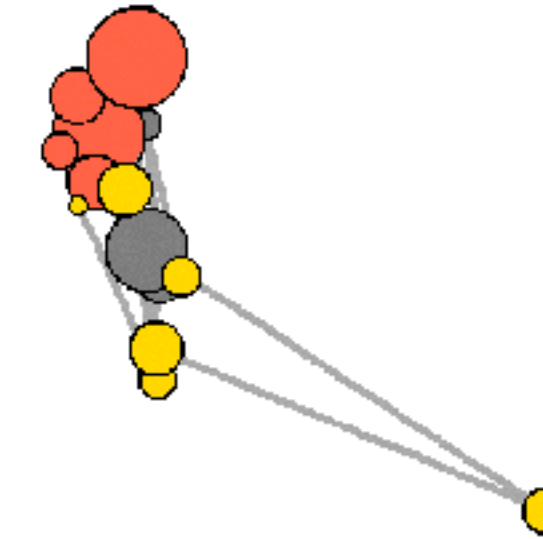
Lecture 9 - preview

- Network visualisation
- Exploratory data analyses

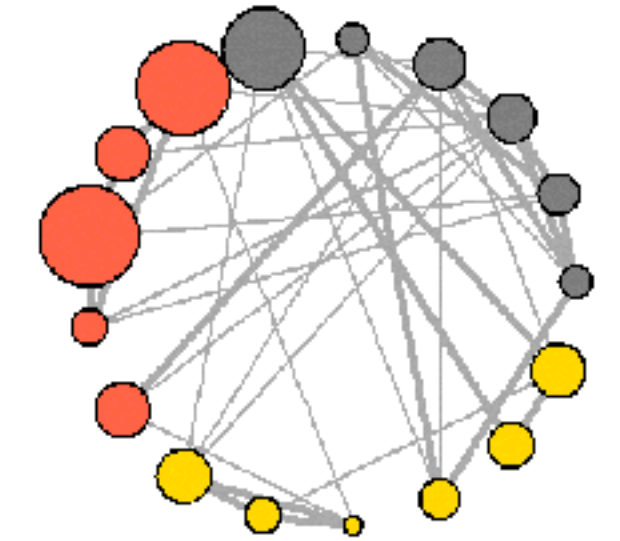
layout_as_star



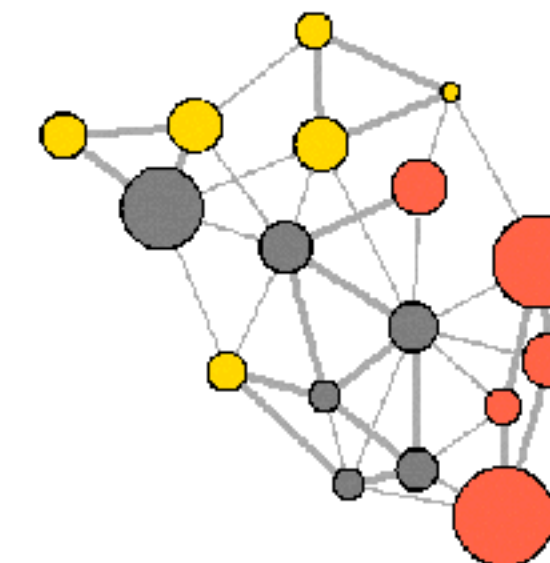
layout_components



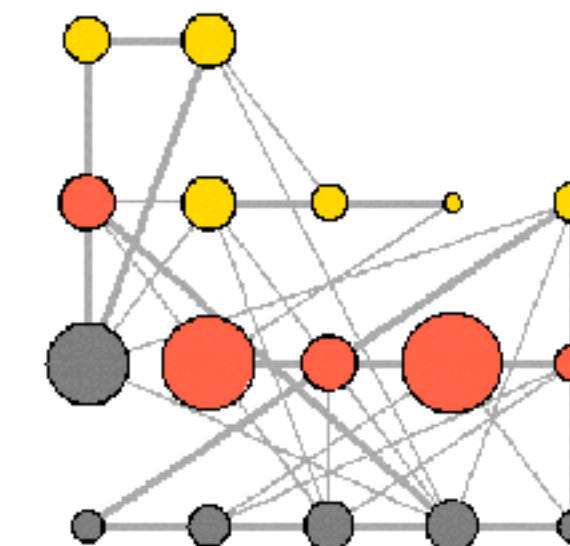
layout_in_circle



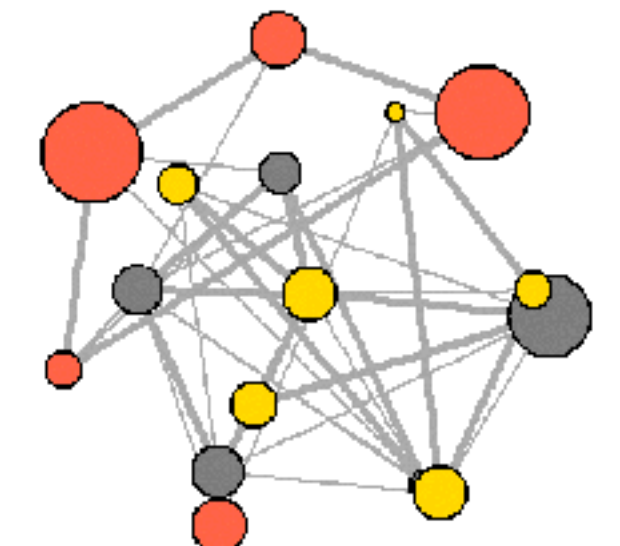
layout_nicely



layout_on_grid



layout_on_sphere



In-class exercise

- **Instruction:**

- Go to Insendi and download the markdown:
- Work together with your classmates in the breakout room
- If you have a question, send a message to the instructor
 - You may be pulled out of breakout room if there is a common question
 - Also, check the forum to see answers to FAQs
- Submit the HTML output individually, via Insendi by **the next day 7am (UK time)!**
- **You should now be able to resubmit within the deadline**