



Data Visualization

Kylie A. Bem

Northeastern University
Khoury College of Computer Science



Northeastern University

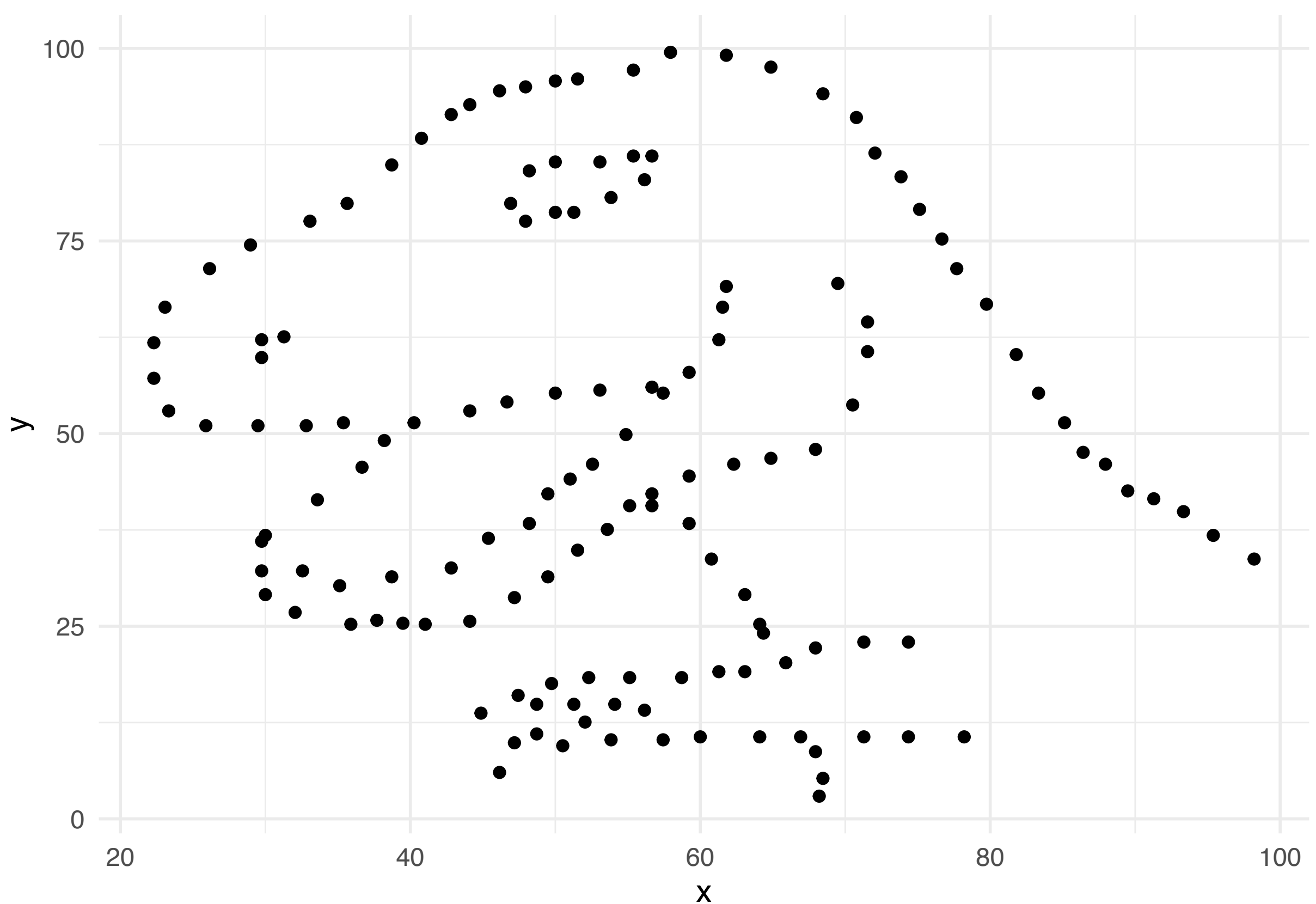
Learning goals

- What are common statistical visualizations
- How to look at data
- Key ingredients of useful data visualizations
- A grammar of graphics

STATISTICAL G HOW TO LOOK

Why do we look

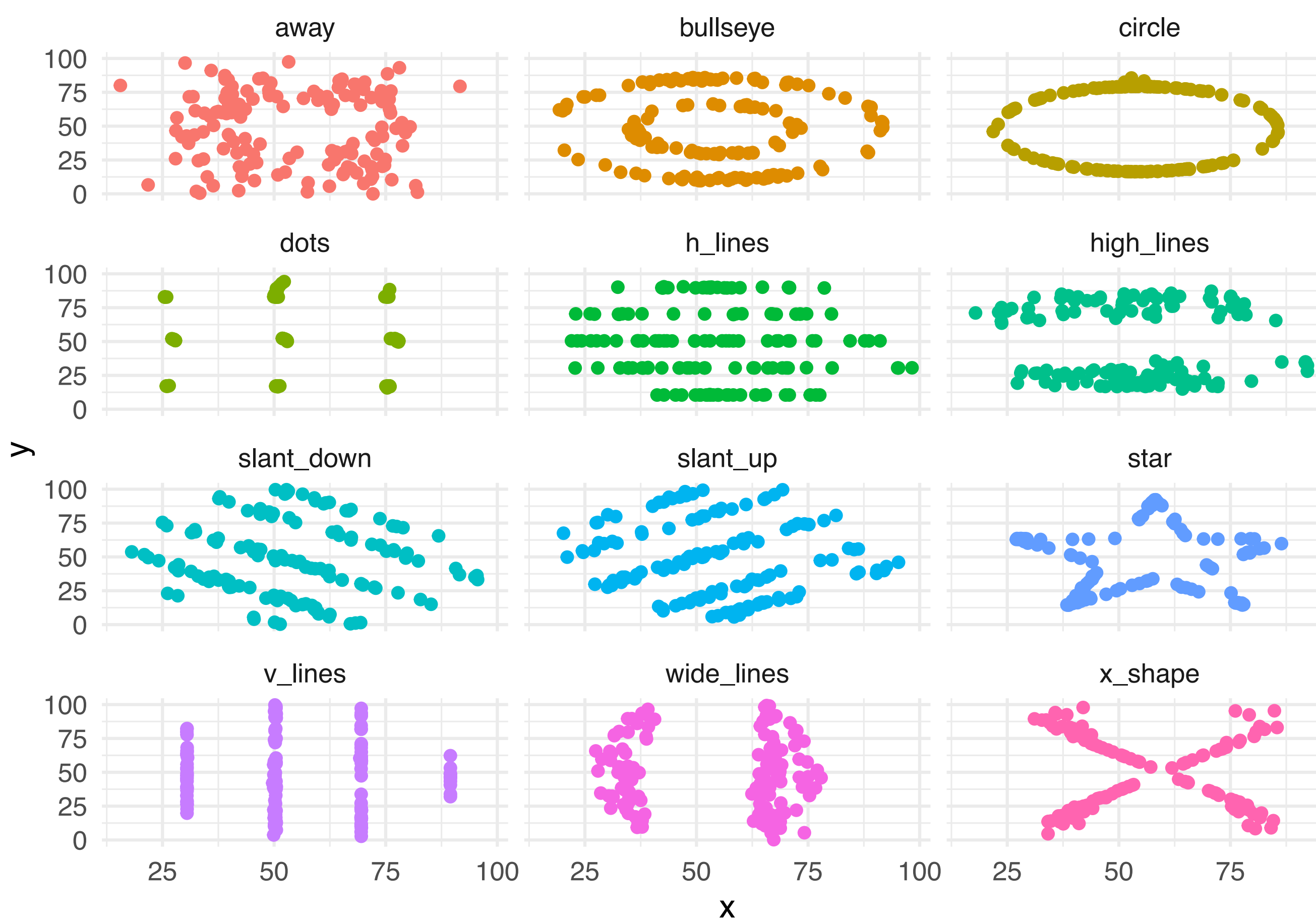
Dataset #1:



How similar are the other 2 datasets?

Why we look at

Datasets #2-13:



“The Datasaurus Dozen”: <https://www.autodeskresearch.com/publications/datasaurusdozen>

Looking at data is

- Summary statistics don't
- Easily find patterns
- Identify potential outliers
- Check model assumptions
- Intuitively display results

What are some common ways

Some common statis

- Scatter plot
- Line plot
- Box-and-whisker plot
- Histogram
- Bar plot

Roles of statistical

One variable

- Histogram
- Bar plot
- Box plot
- Pie chart

Roles of statistical

Distributions

- Histogram
- Bar plot
- Box plot
- Pie chart

SINGLE-VARIABLE

Example data: Ga

Life expectancy, GDP per capita, an

```
gapminder

## #      tibble: 1,704 x 6
##      country      continent  year  lif
##      <fct>        <fct>      <int>  <
##  1  fghanistan    sia        1952
##  2  fghanistan    sia        1957
##  3  fghanistan    sia        1962
##  4  fghanistan    sia        1967
```

<http://www.gapminder.org>

Example data: Fuel

Fuel economy on 38 popular models of

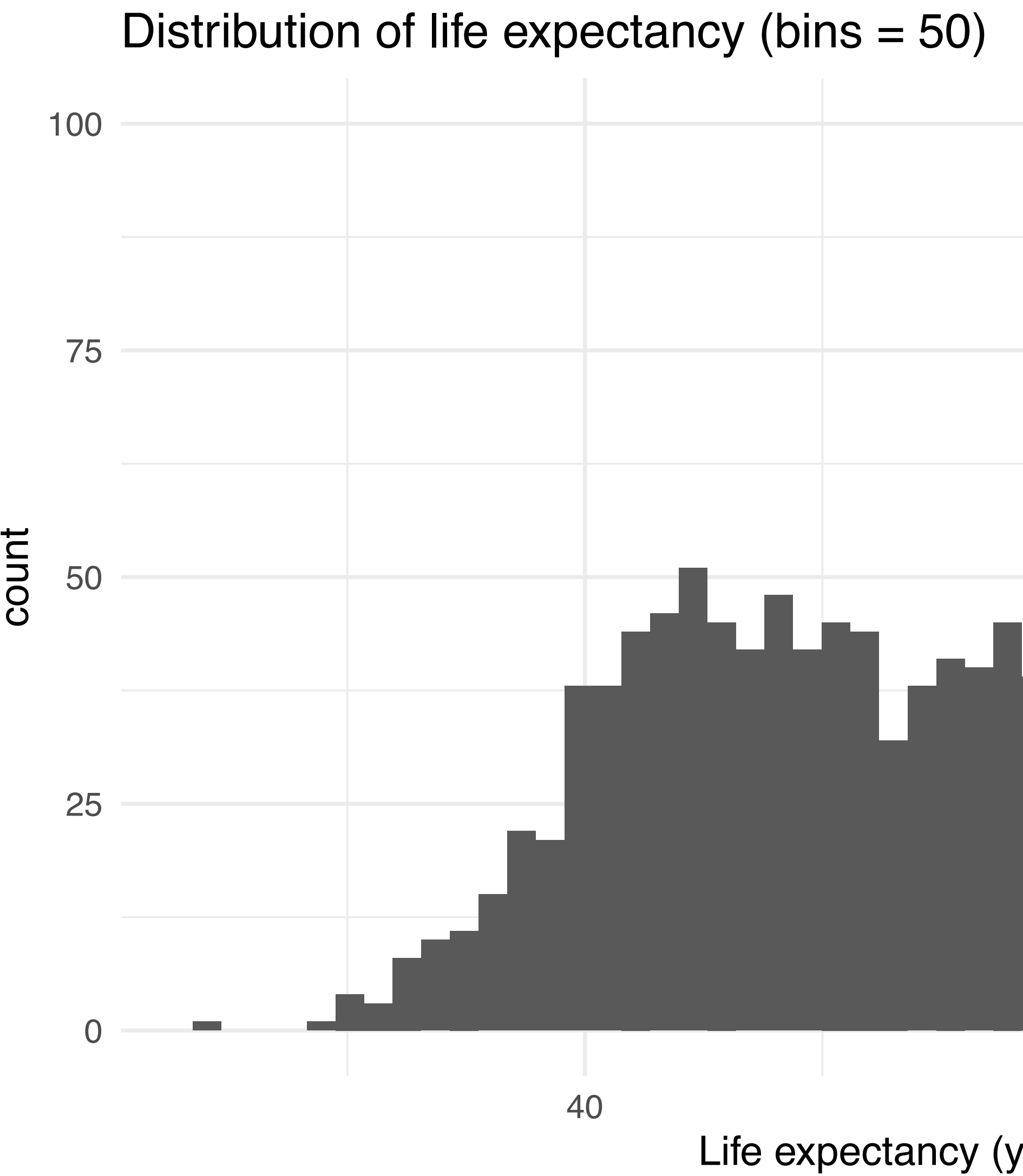
mpg							
##	#	tibble: 234 x 11					
##		manufacturer	model	displ	year	cyl	trans
##		<chr>	<chr>	<dbl>	<int>	<int>	<chr>
##	1	audi	a4	1.8	1999	4	auto
##	2	audi	a4	1.8	1999	4	manu
##	3	audi	a4	2	2008	4	manu
##	4	audi	a4	2	2008	4	auto
##	5	audi	a4	2.8	1999	6	auto
##	6	audi	a4	2.8	1999	6	manu

<http://fueleconomy.gov>

Looking at a single

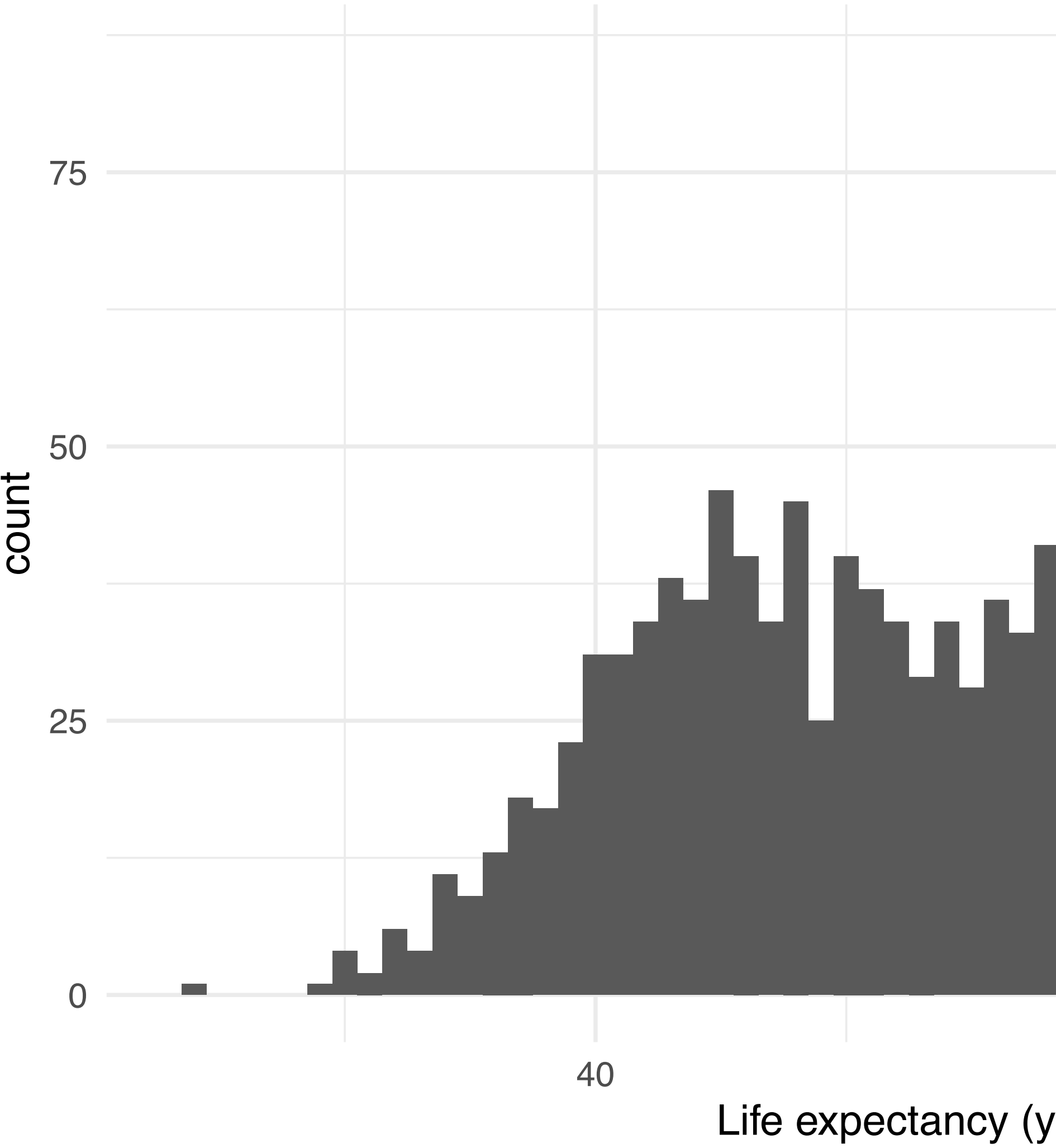
- What is the distribution?
 - ◆ Location - e.g., mean, median, mode
 - ◆ Spread - e.g., variance
 - ◆ Shape - symmetric vs. skewed
- Are there outliers?
- What is notable about the

Histogram

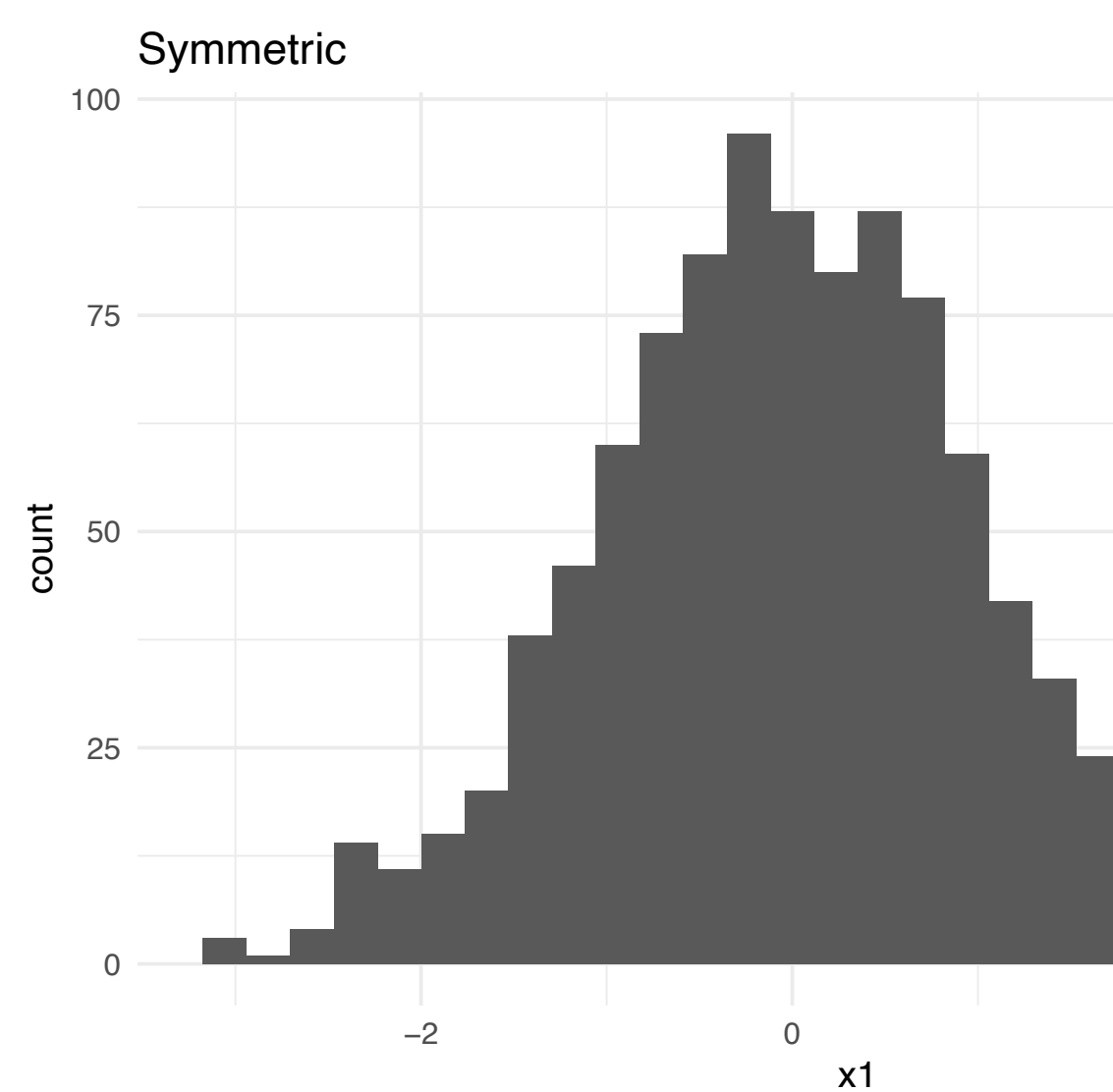
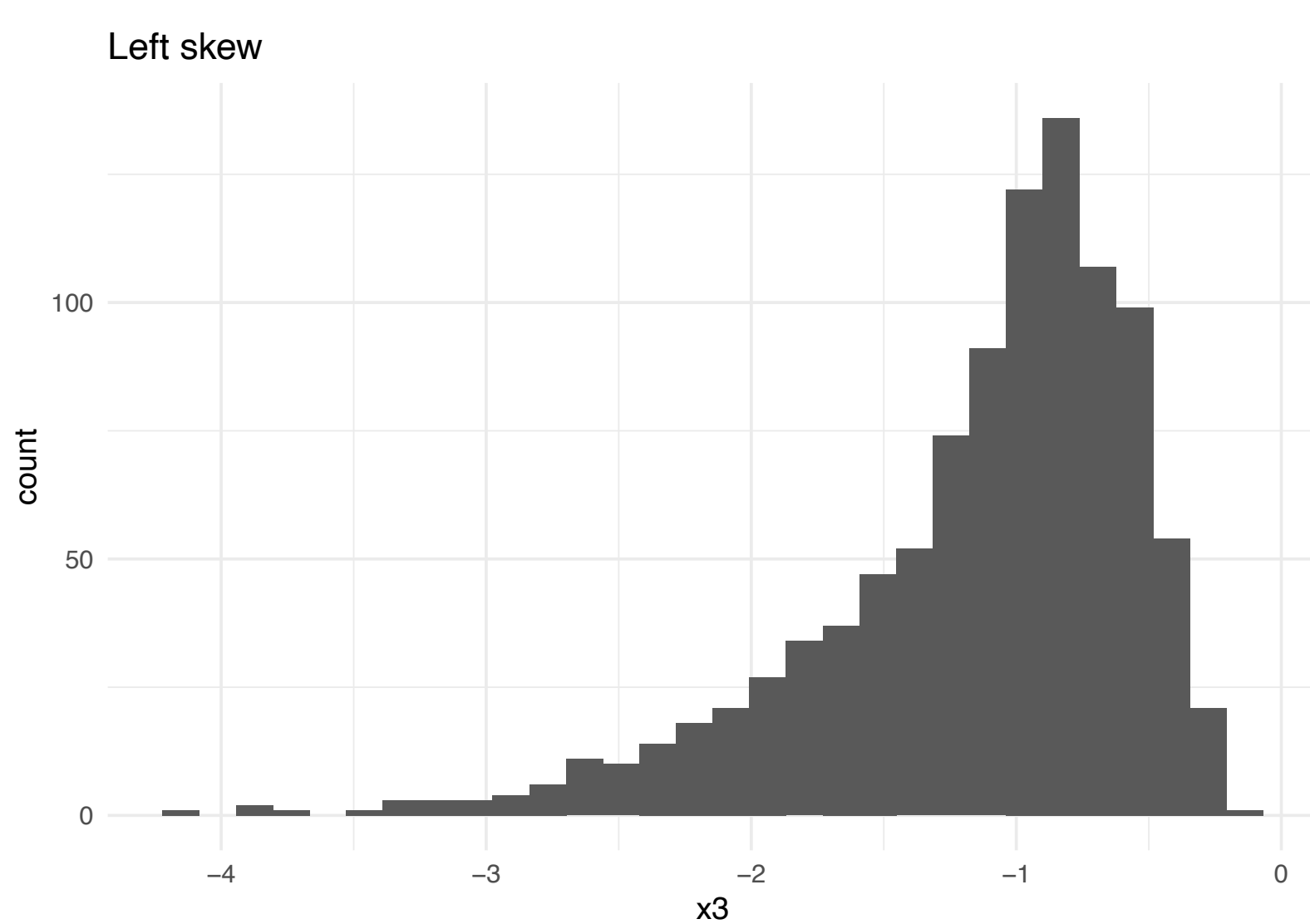


Histogram

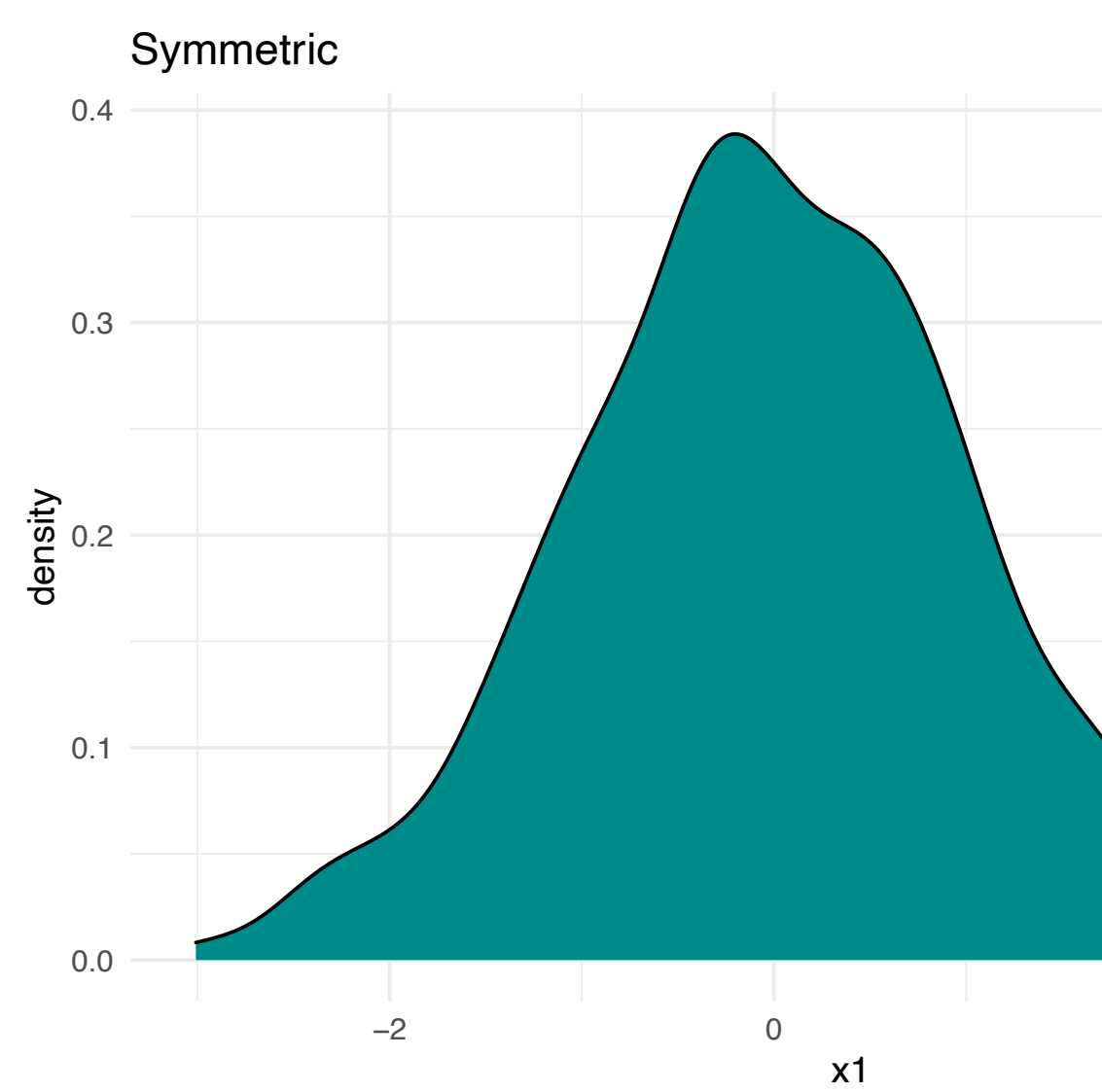
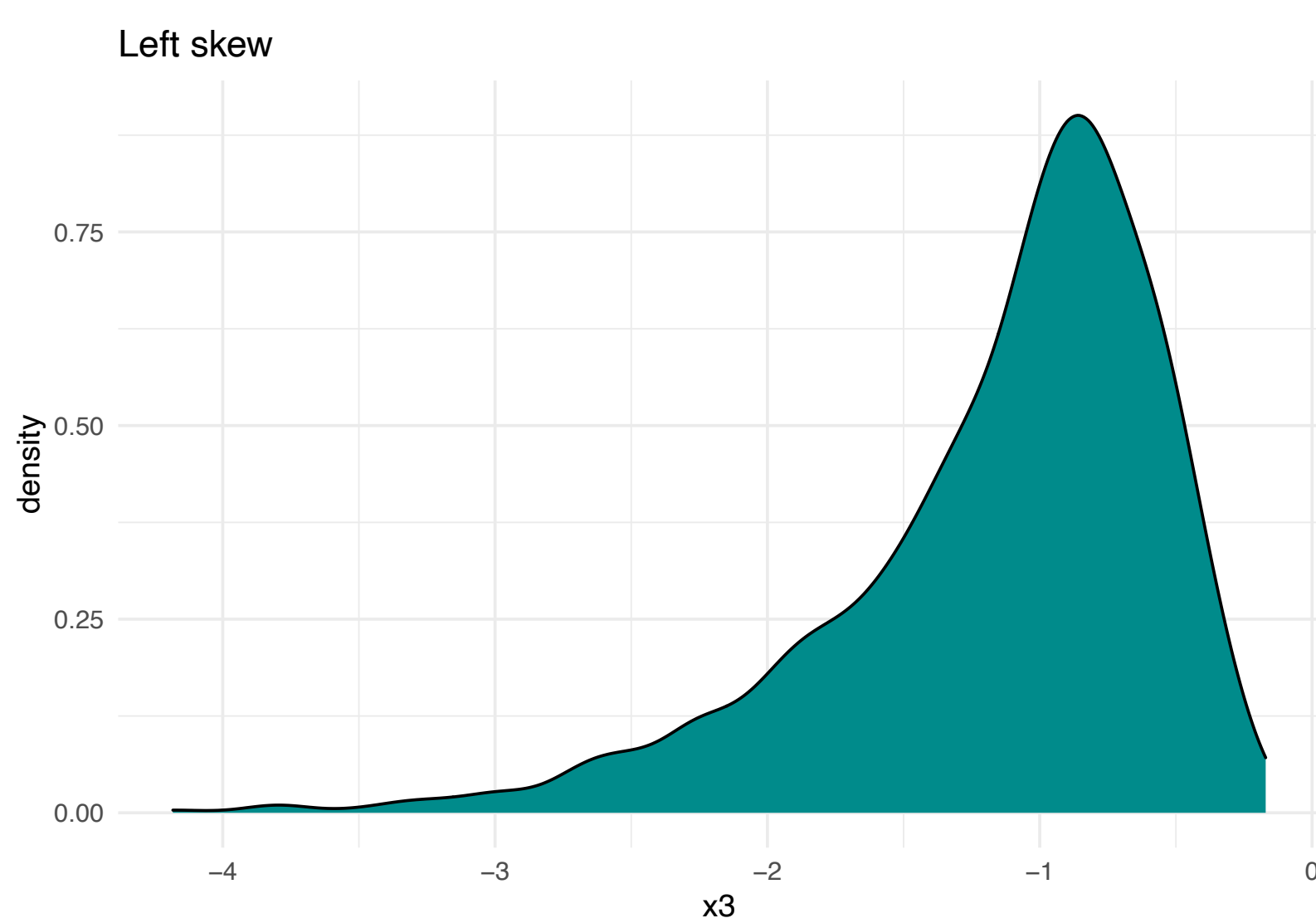
Distribution of life expectancy (binwidth = 1)



Histogram

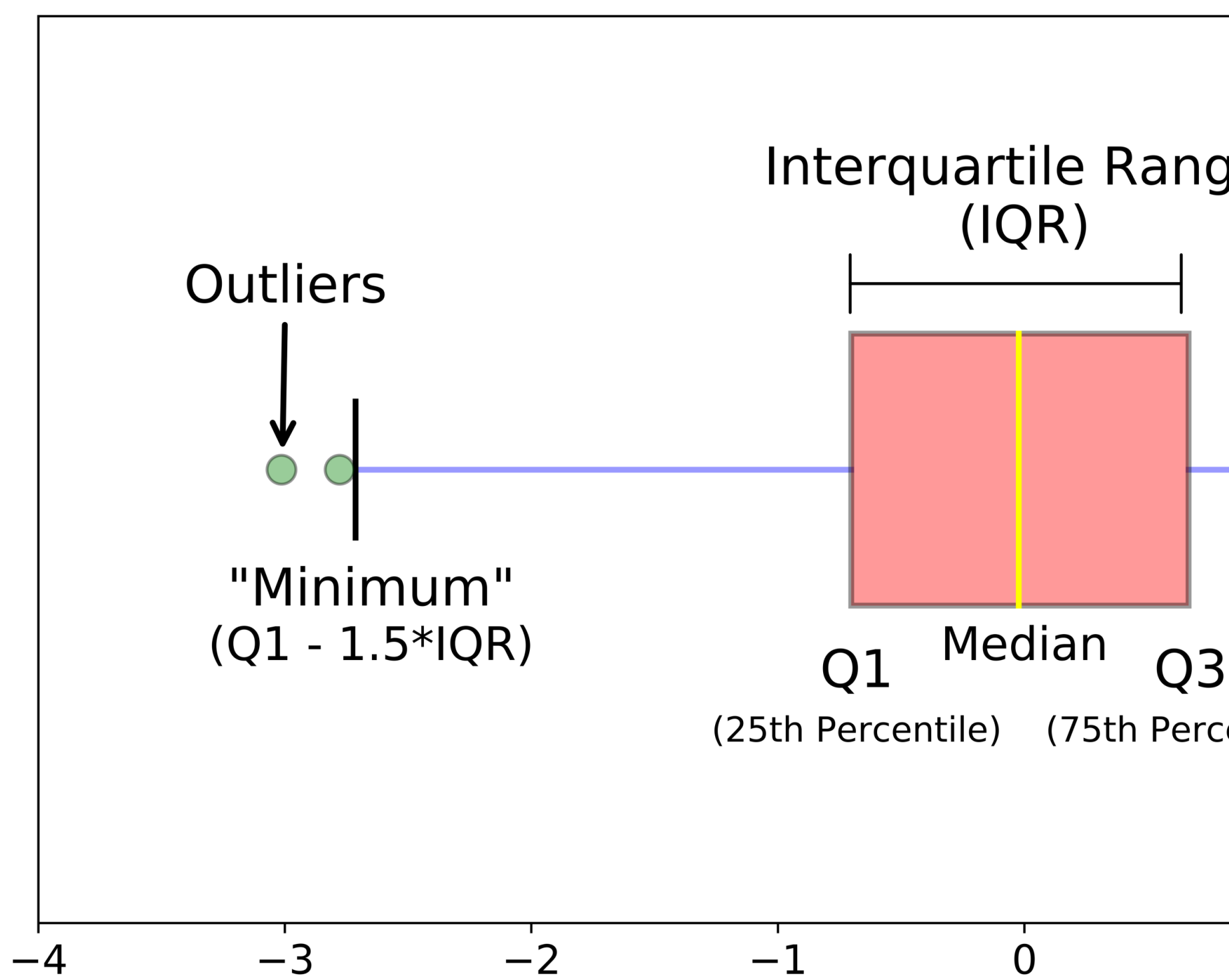


Density plot

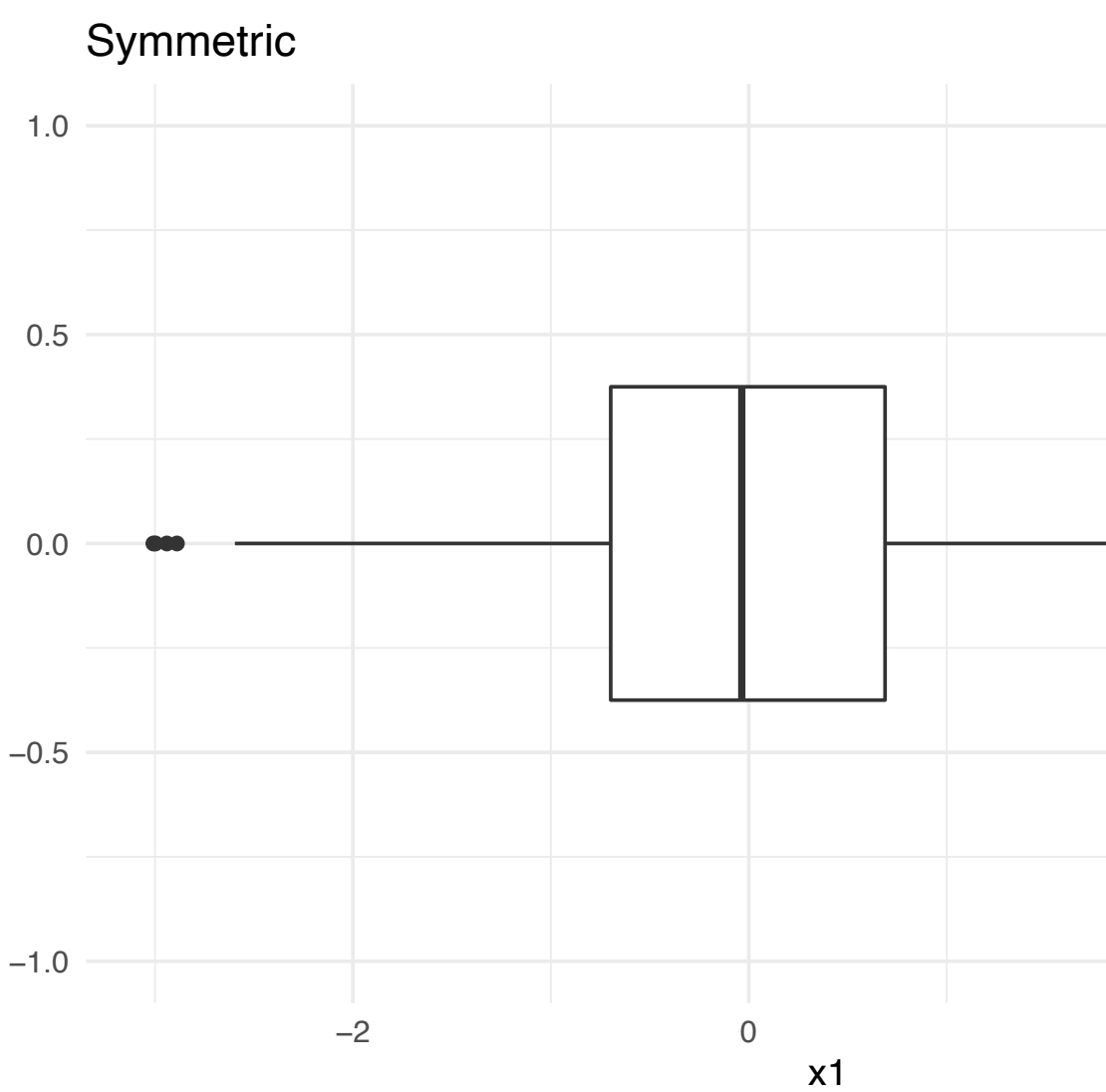
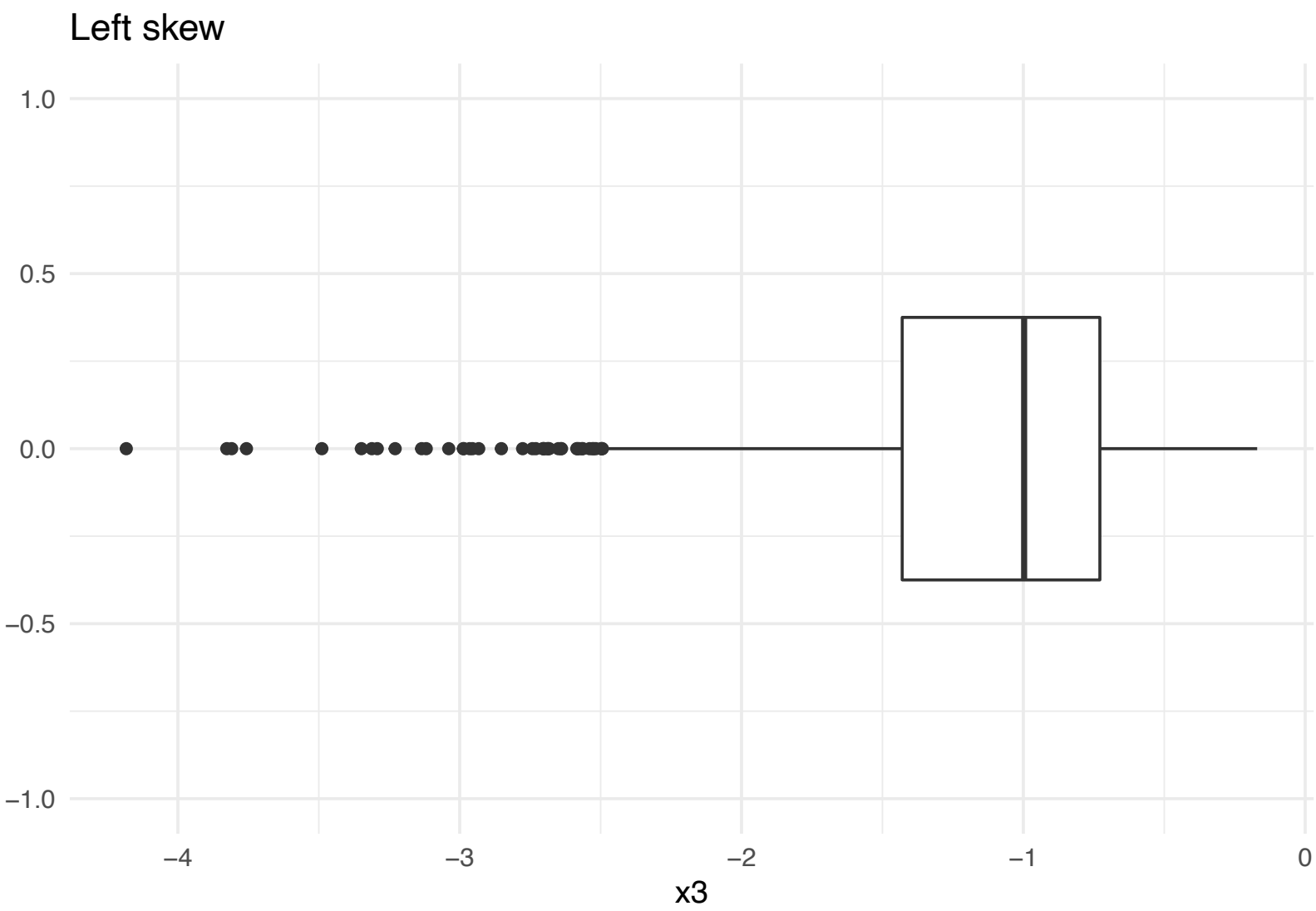


simulated data
with different distributions

Box plot

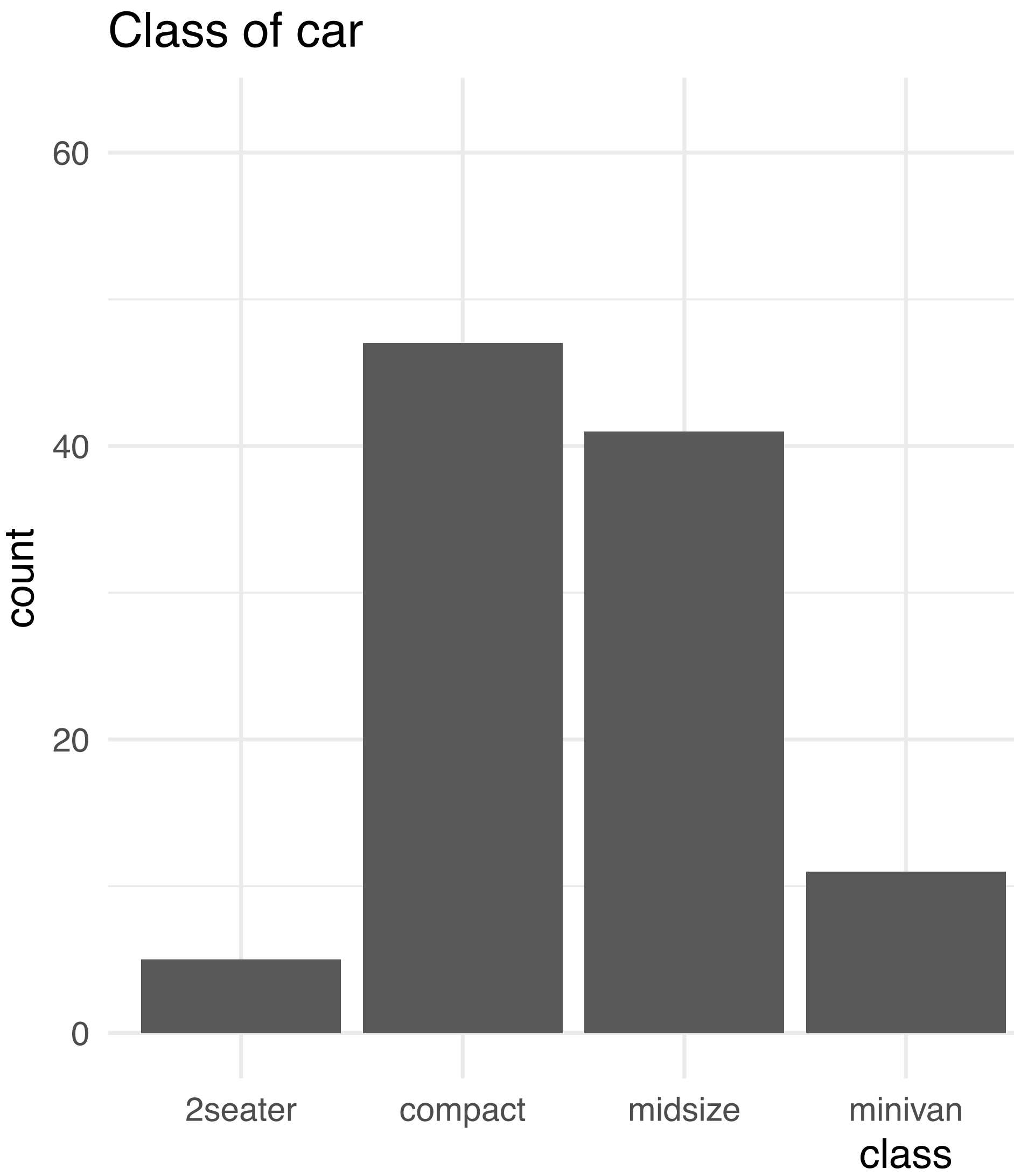


Box plots

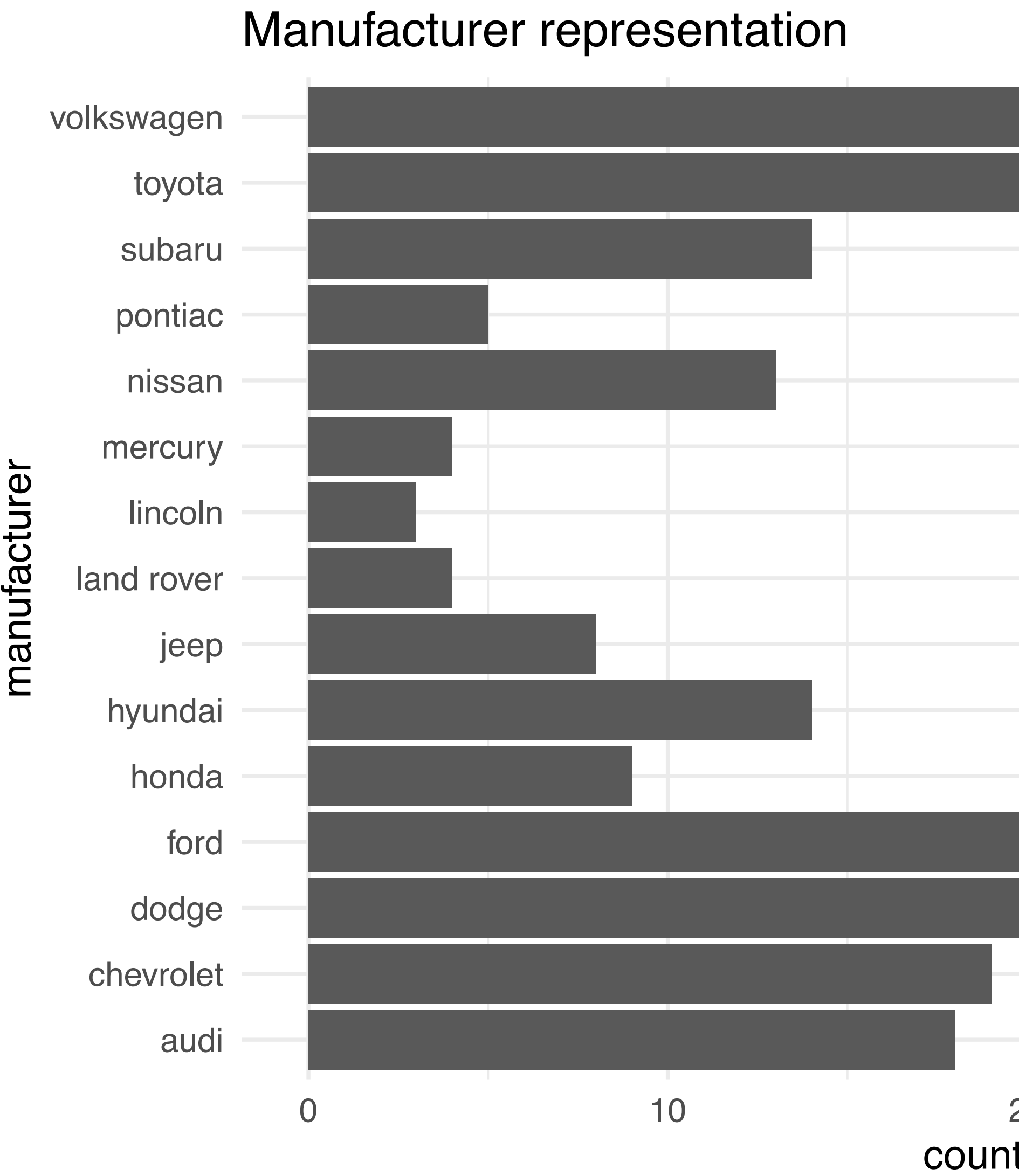


simulated data
with different distributions

Bar plot

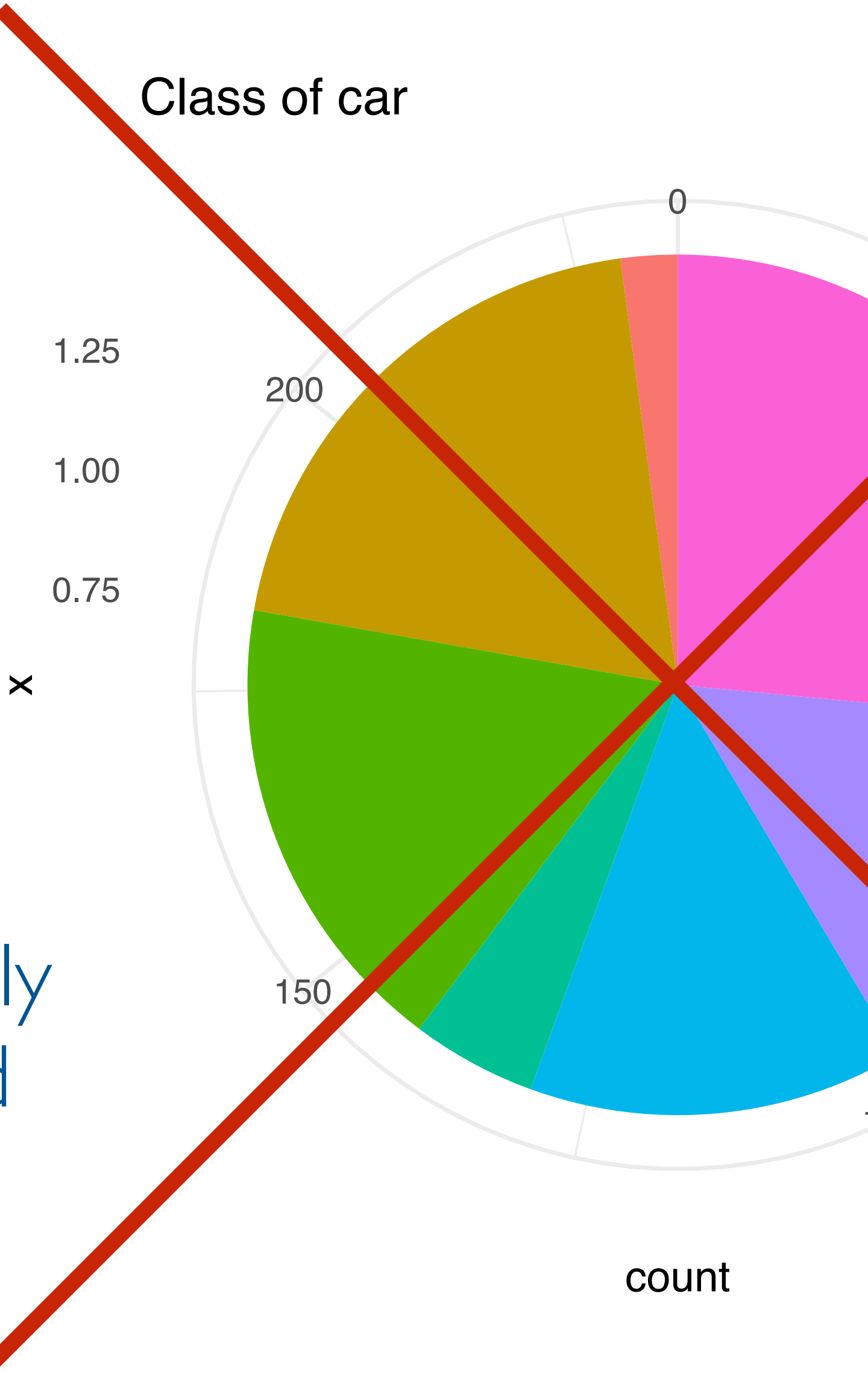


Bar plot



Pie chart

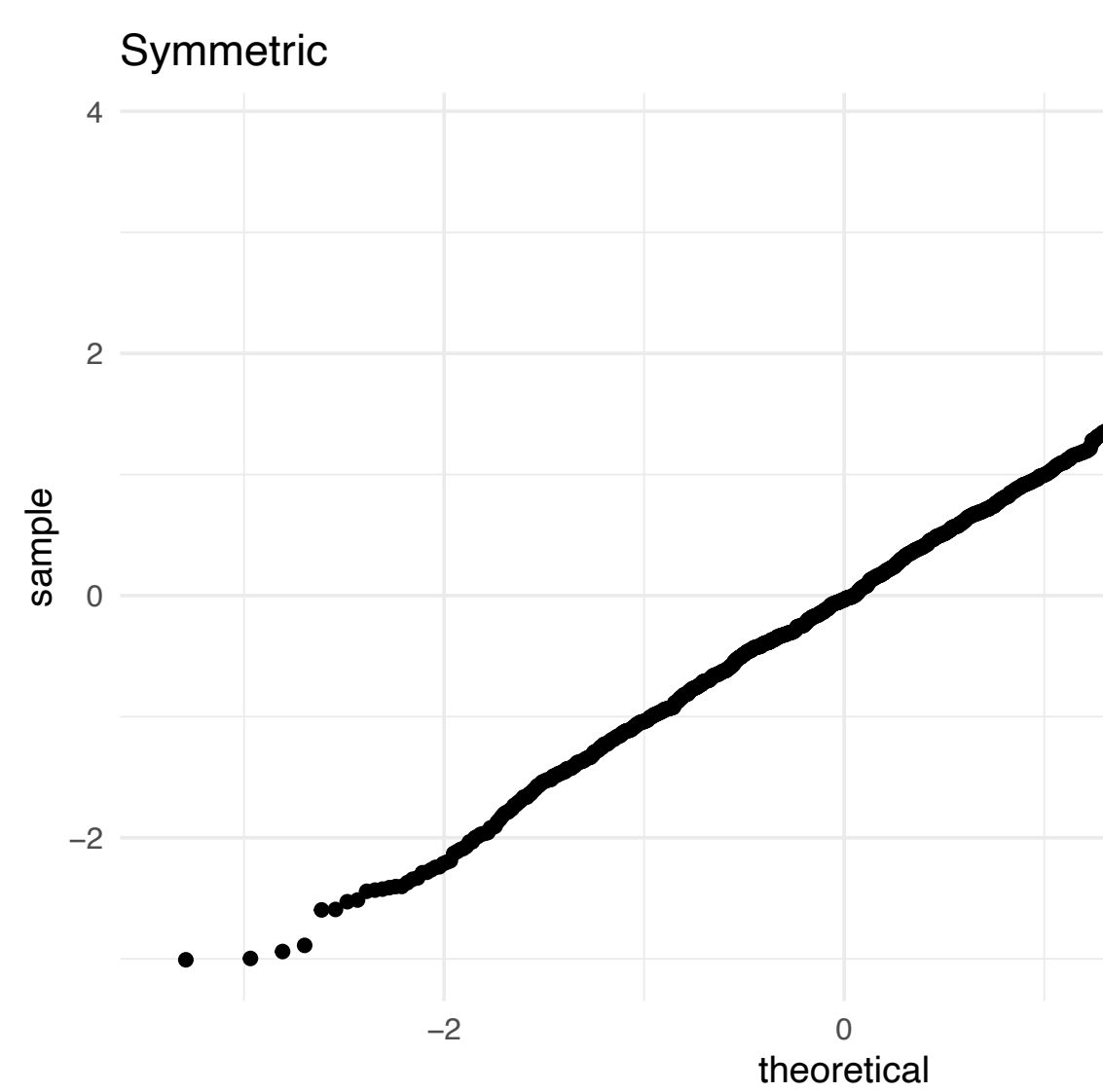
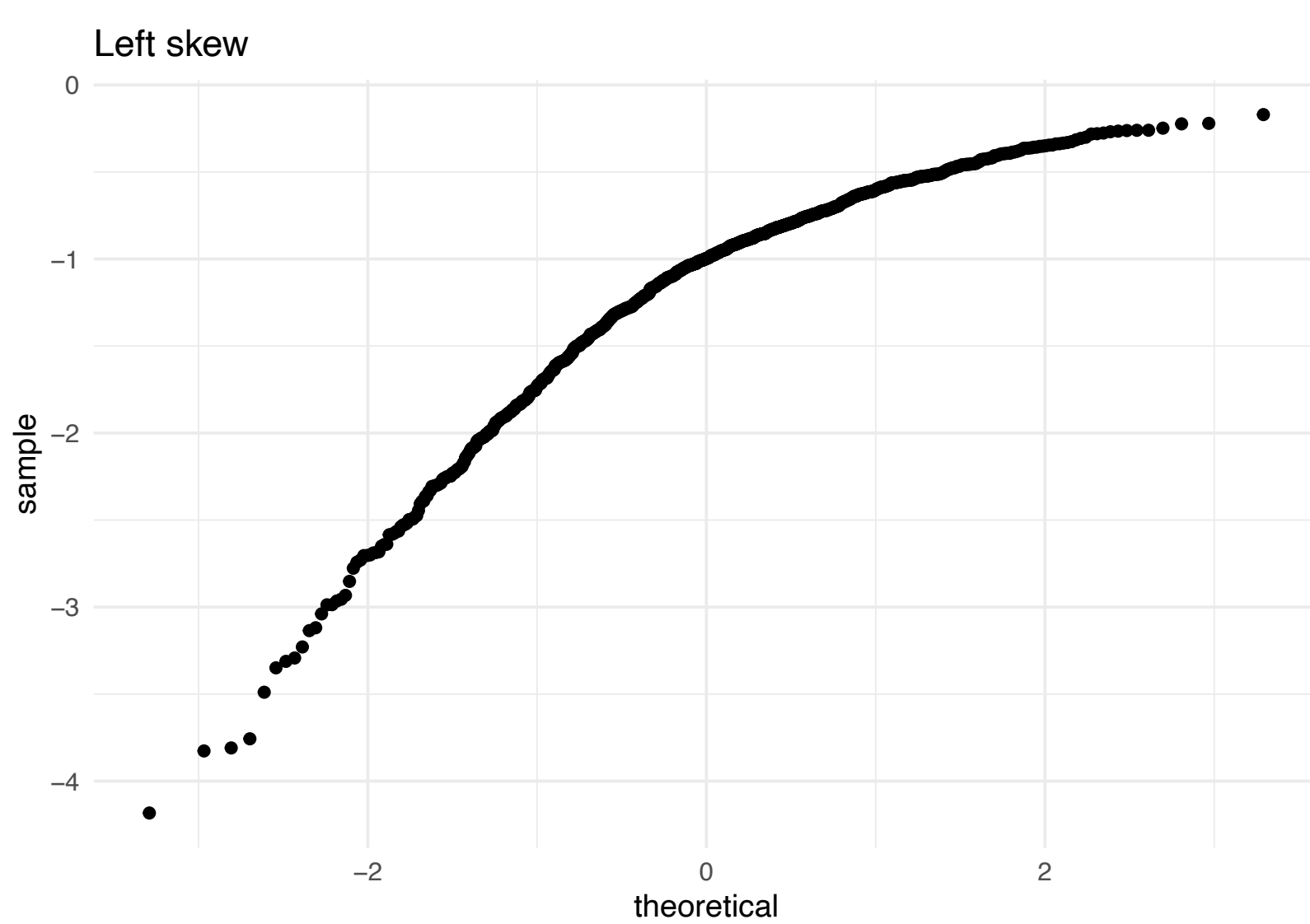
Class of car



difficult to read easily
not recommended

Q-Q plot

sample quantiles versus theoretical quantile



simulated data
with different distributions

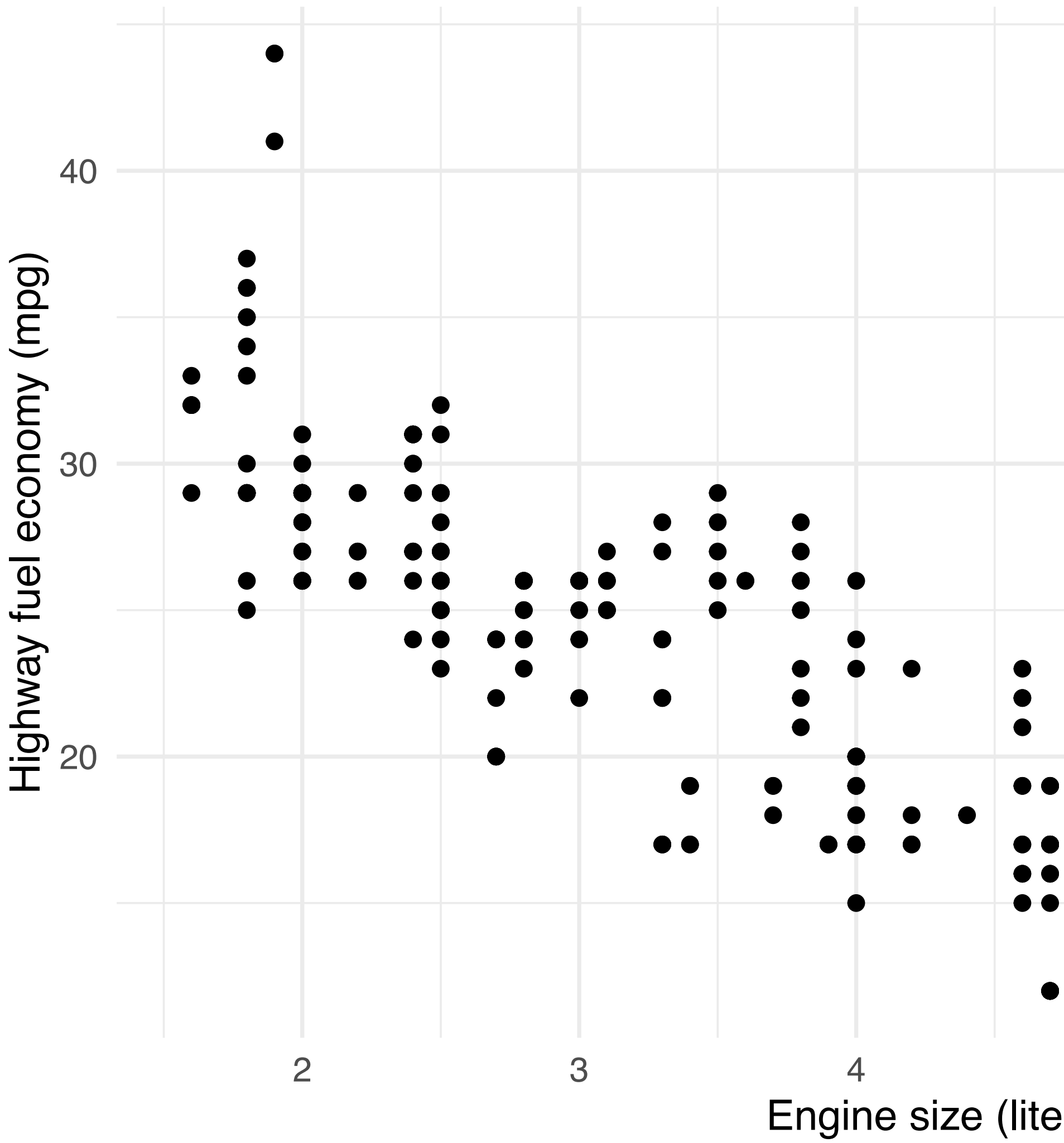
MULTI-VARIABLE

Looking at multiple

- How are the variables related?
 - ◆ Is there a relationship?
 - ◆ Type of relationship (e.g., linear)
 - ◆ Direction (positive vs. negative)
- Are there outliers?
- Does the relationship change?

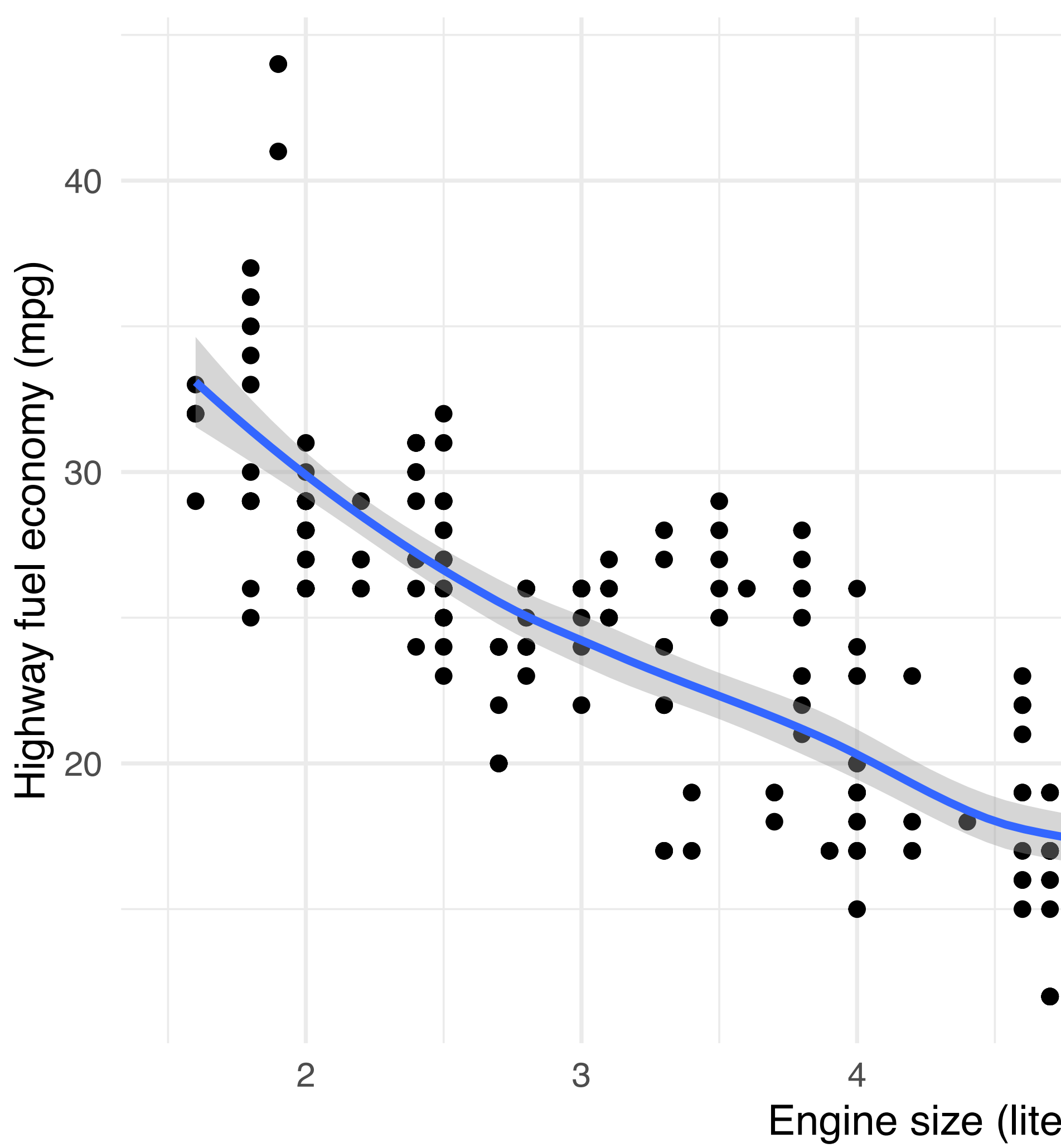
Scatter plot

Larger engines are less efficient



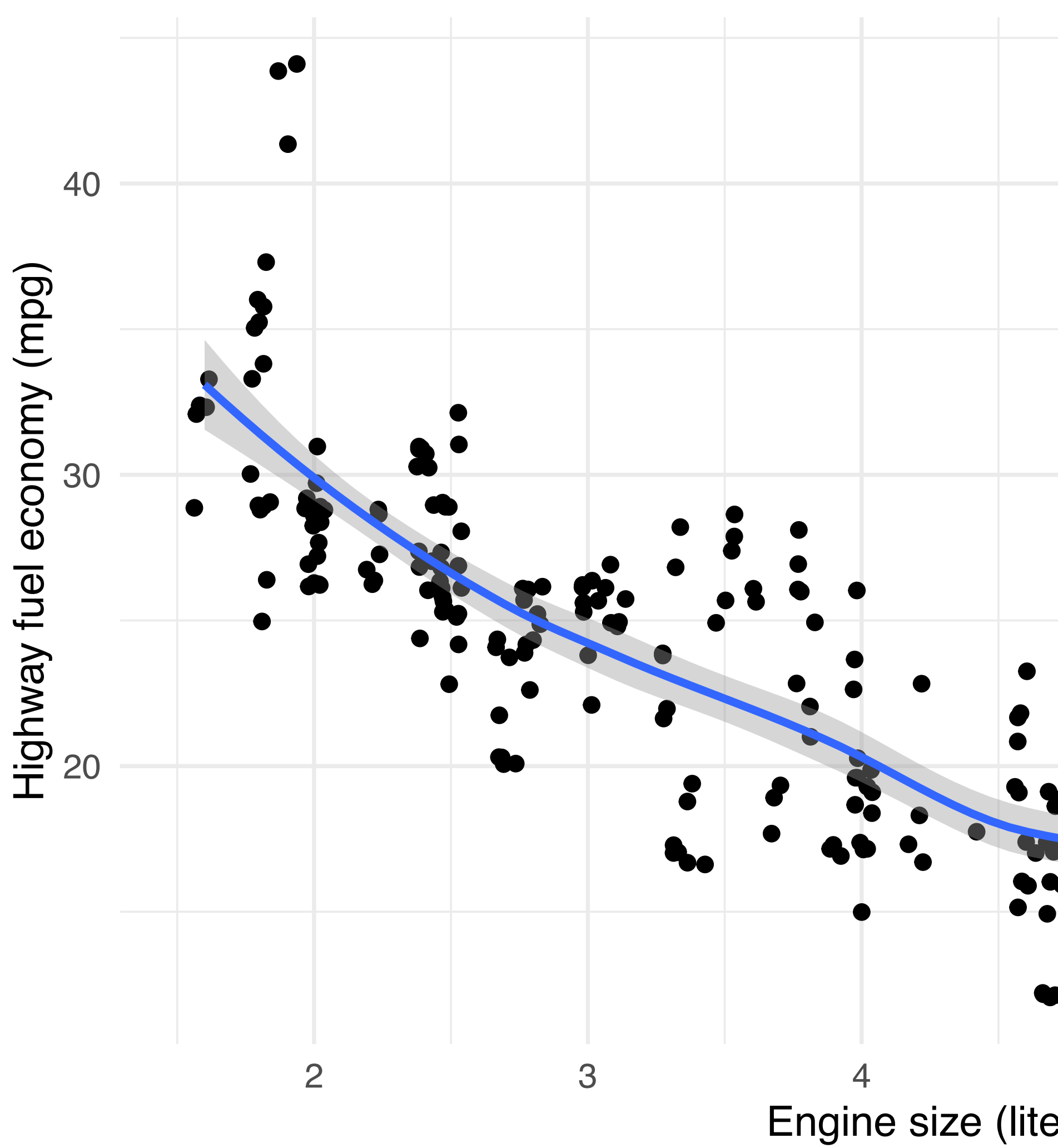
Scatter plot

Larger engines are less efficient



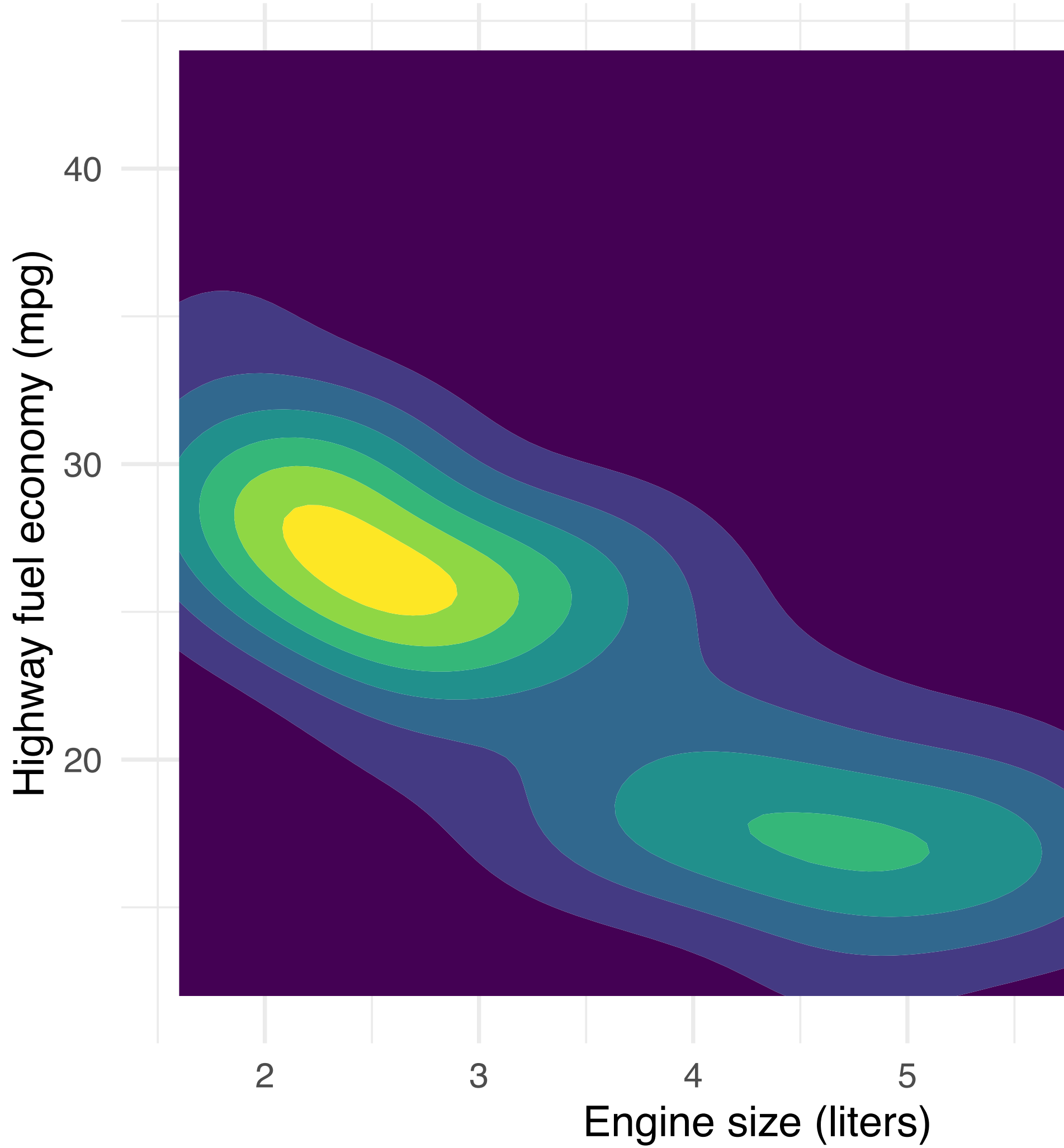
Scatter plot

Larger engines are less efficient

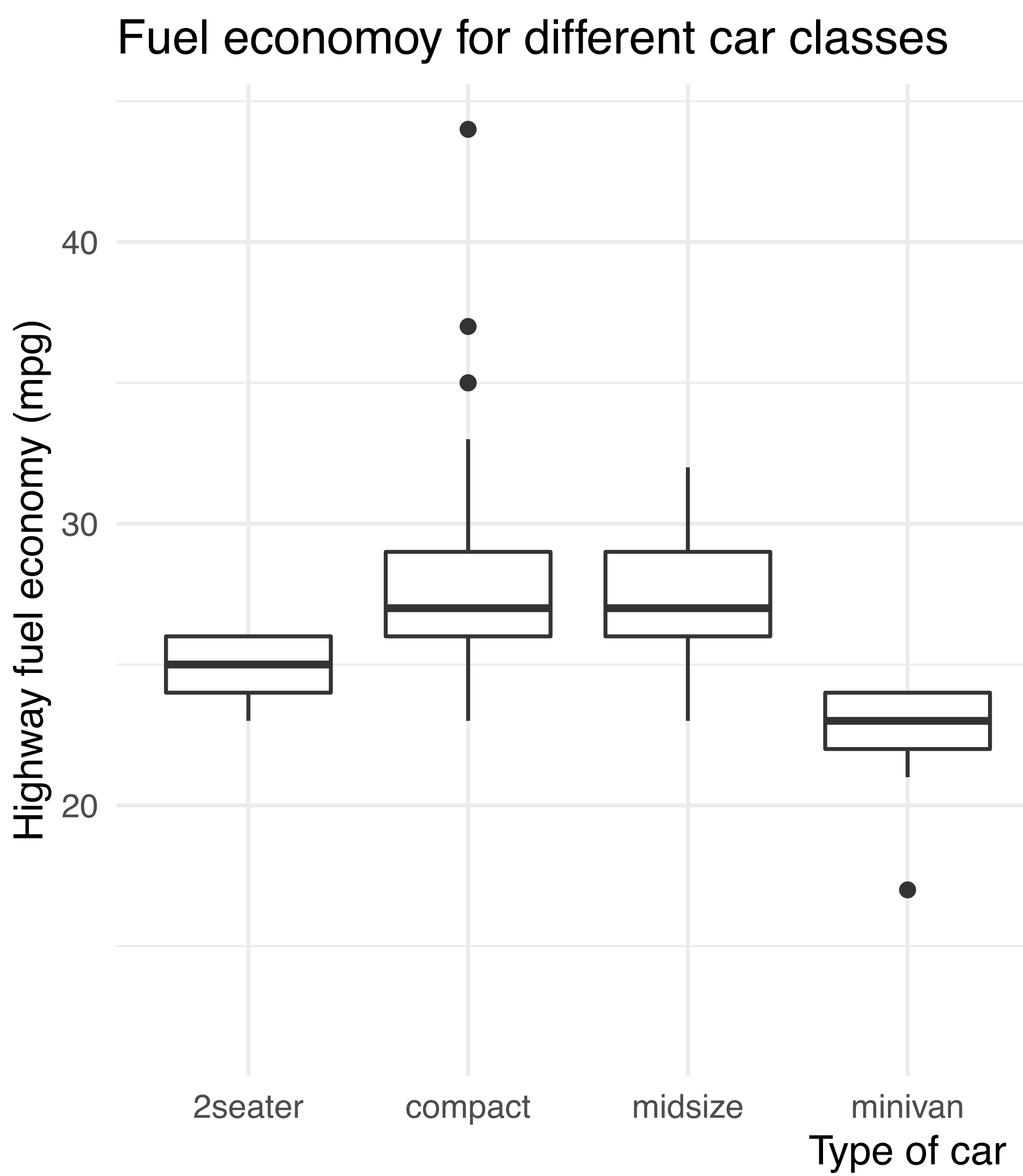


2D density plot

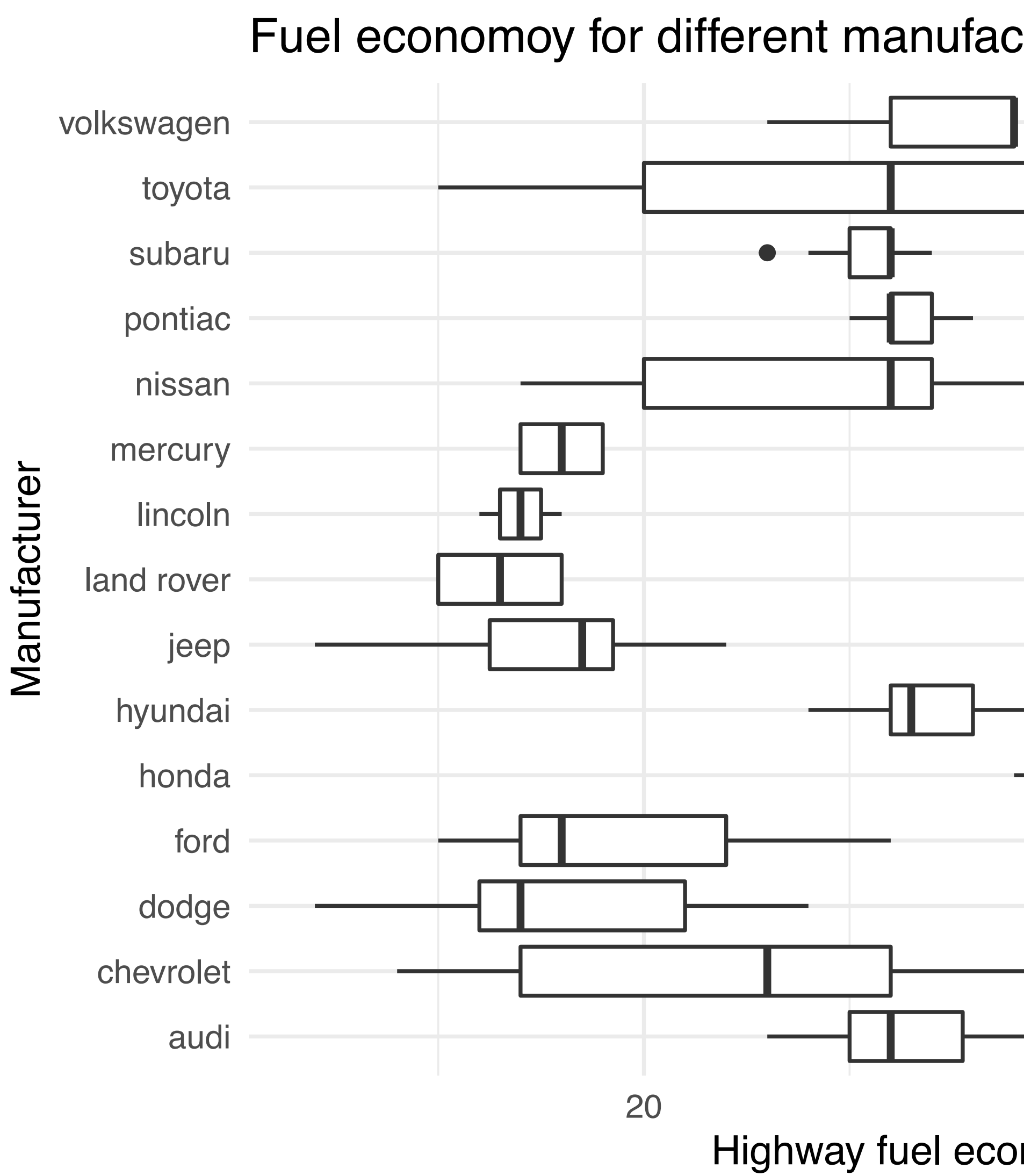
Larger engines are less efficient



Box plots

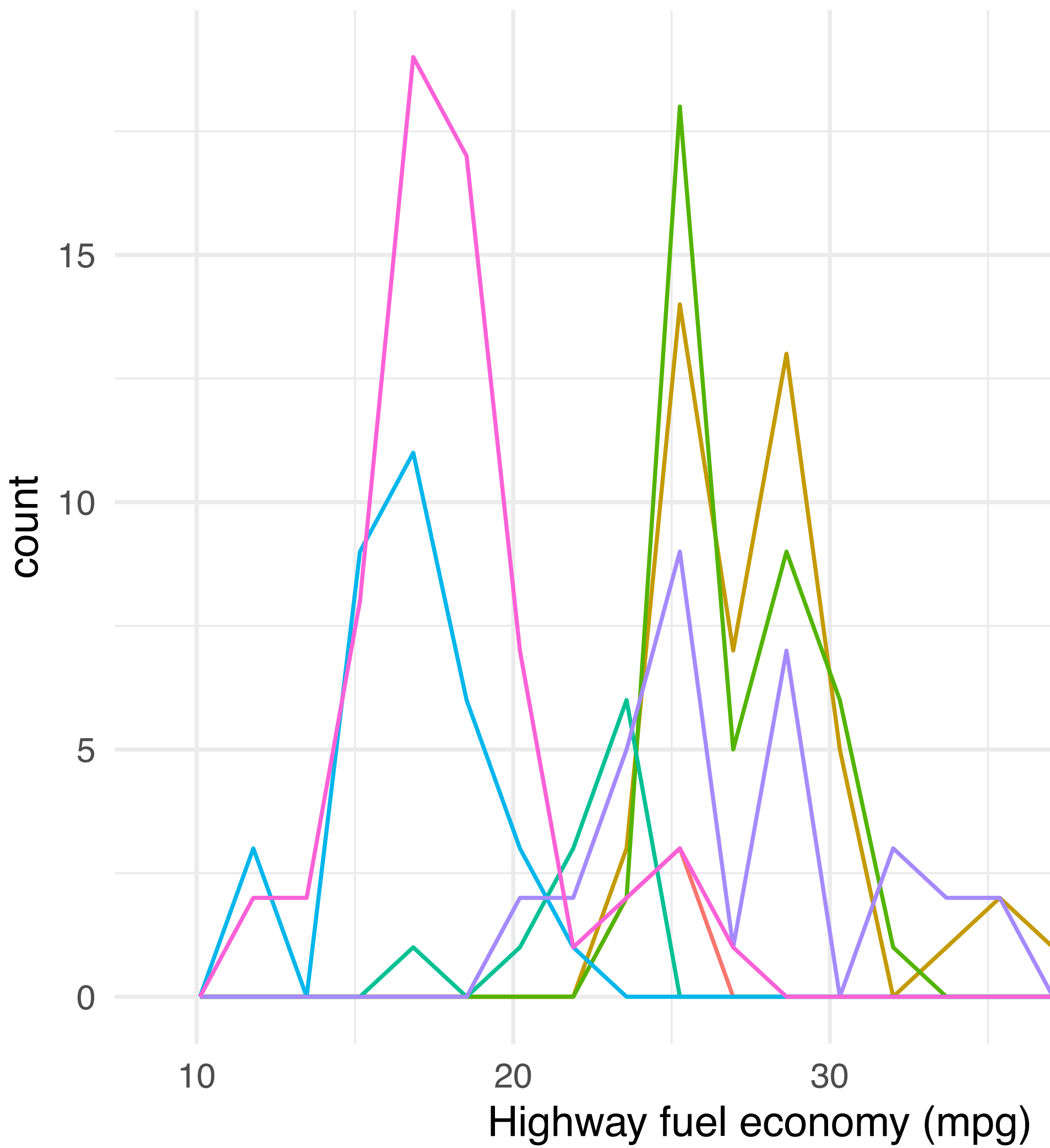


Box plots



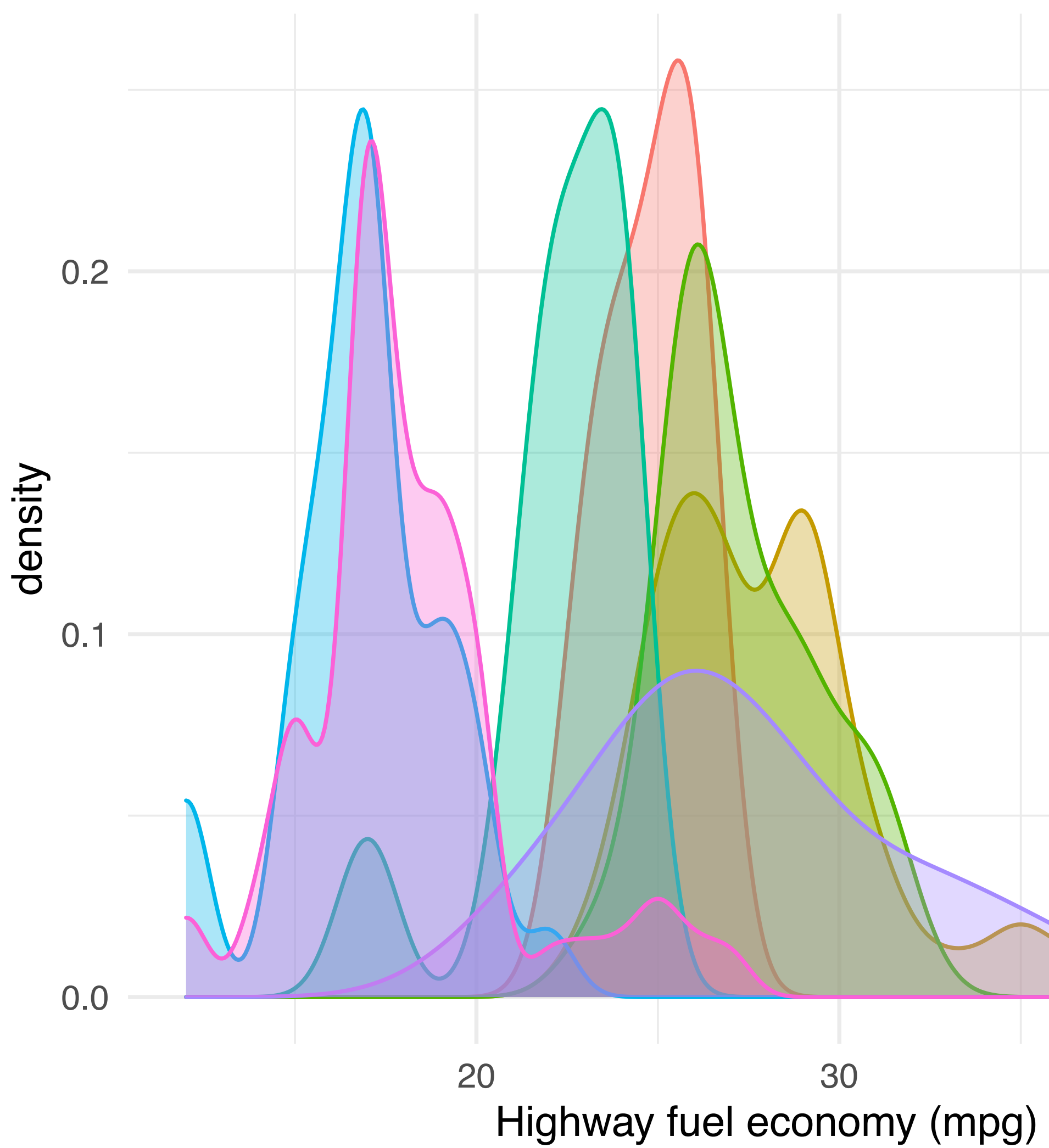
Histogram

Fuel economoy for different car classes



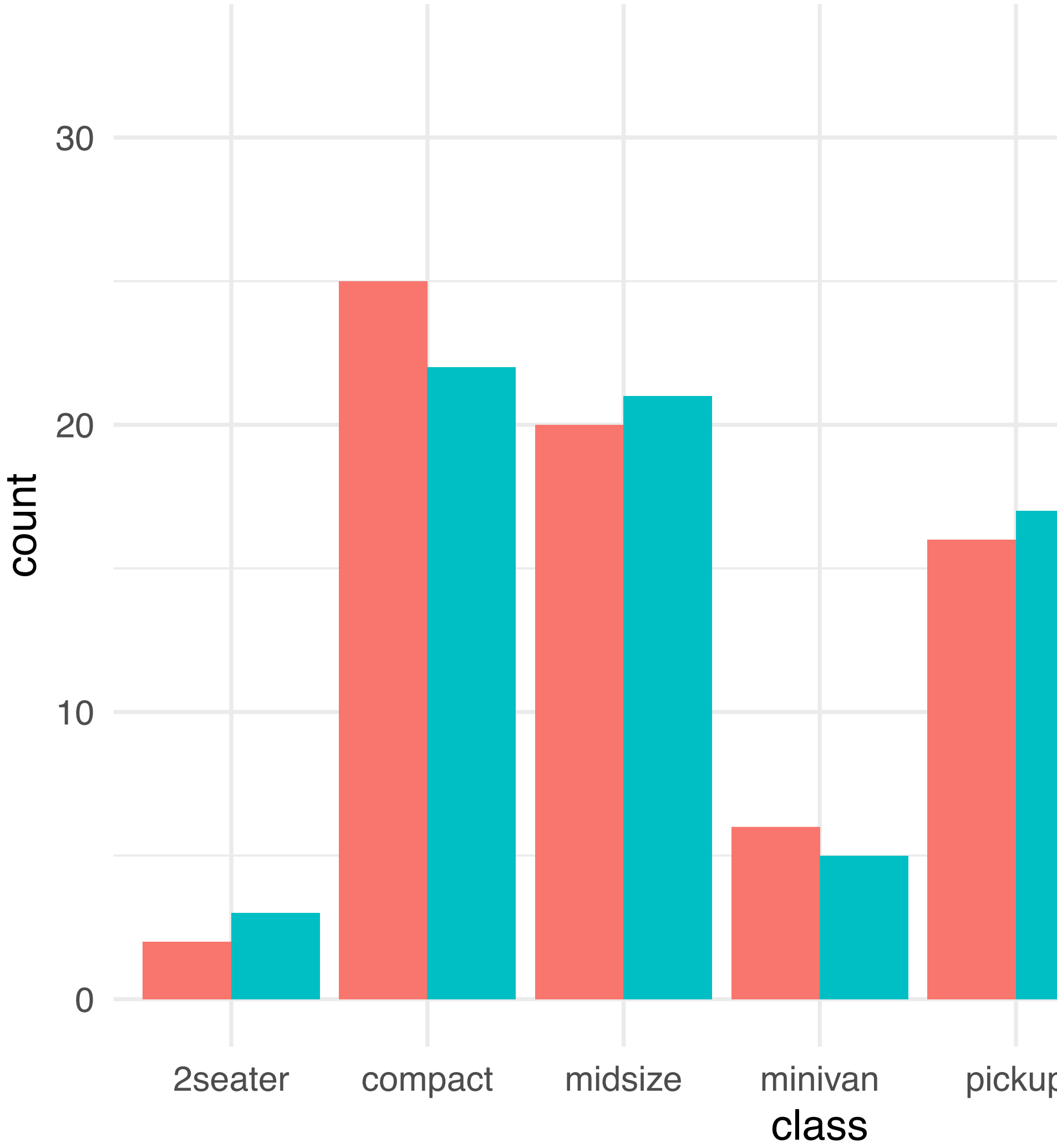
Density plot

Fuel economy for different car classes



Bar plots

Breakdown of car types by year



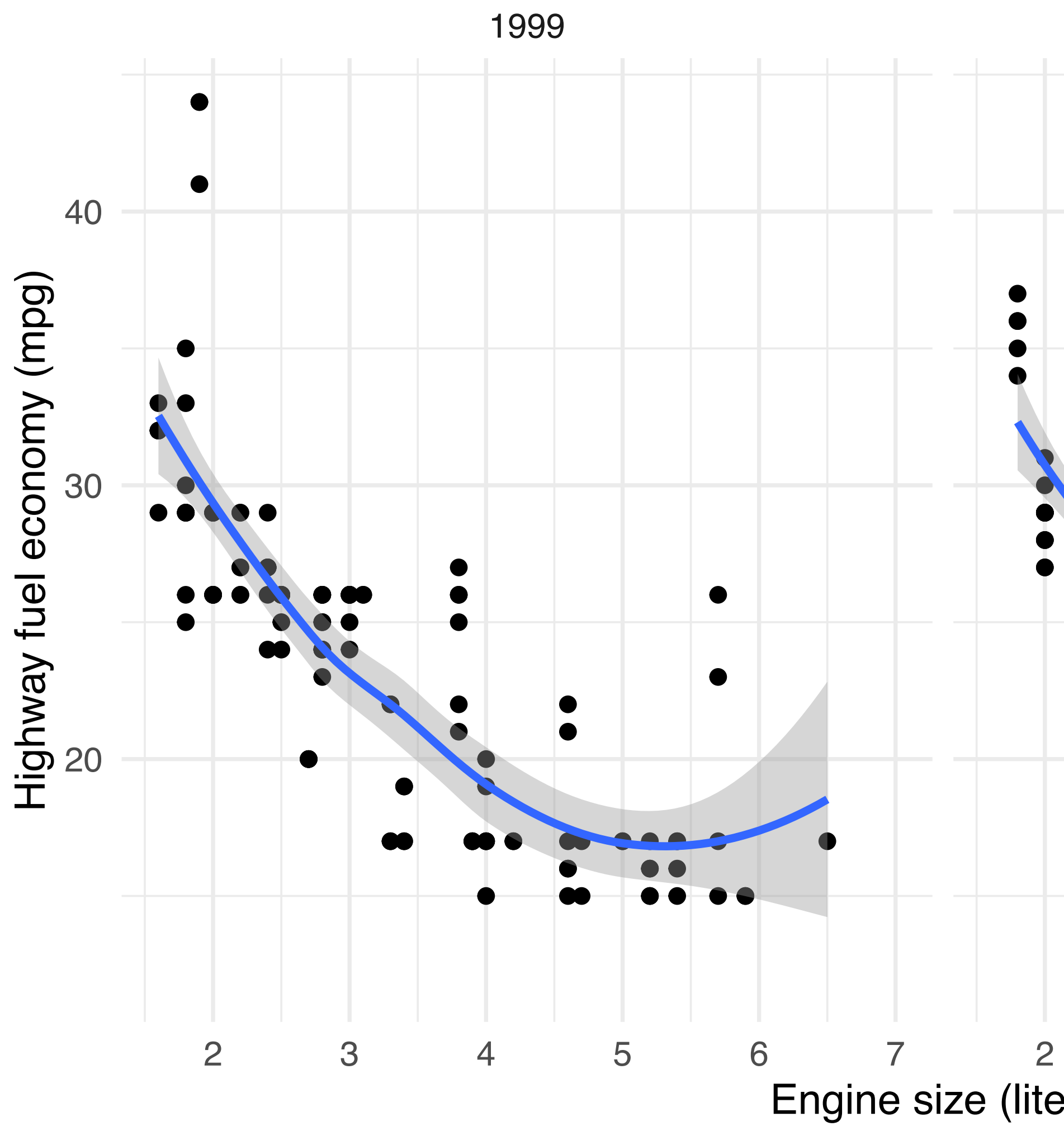
FACETIN

What is faceting

- Based on idea of “small multiples”
- Condition on levels of a variable
- Split data into subsets based on condition
- Create sub-plots for each subset
 - ◆ Sub-plots share same scales and axes
 - ◆ Easily compare between sub-plots

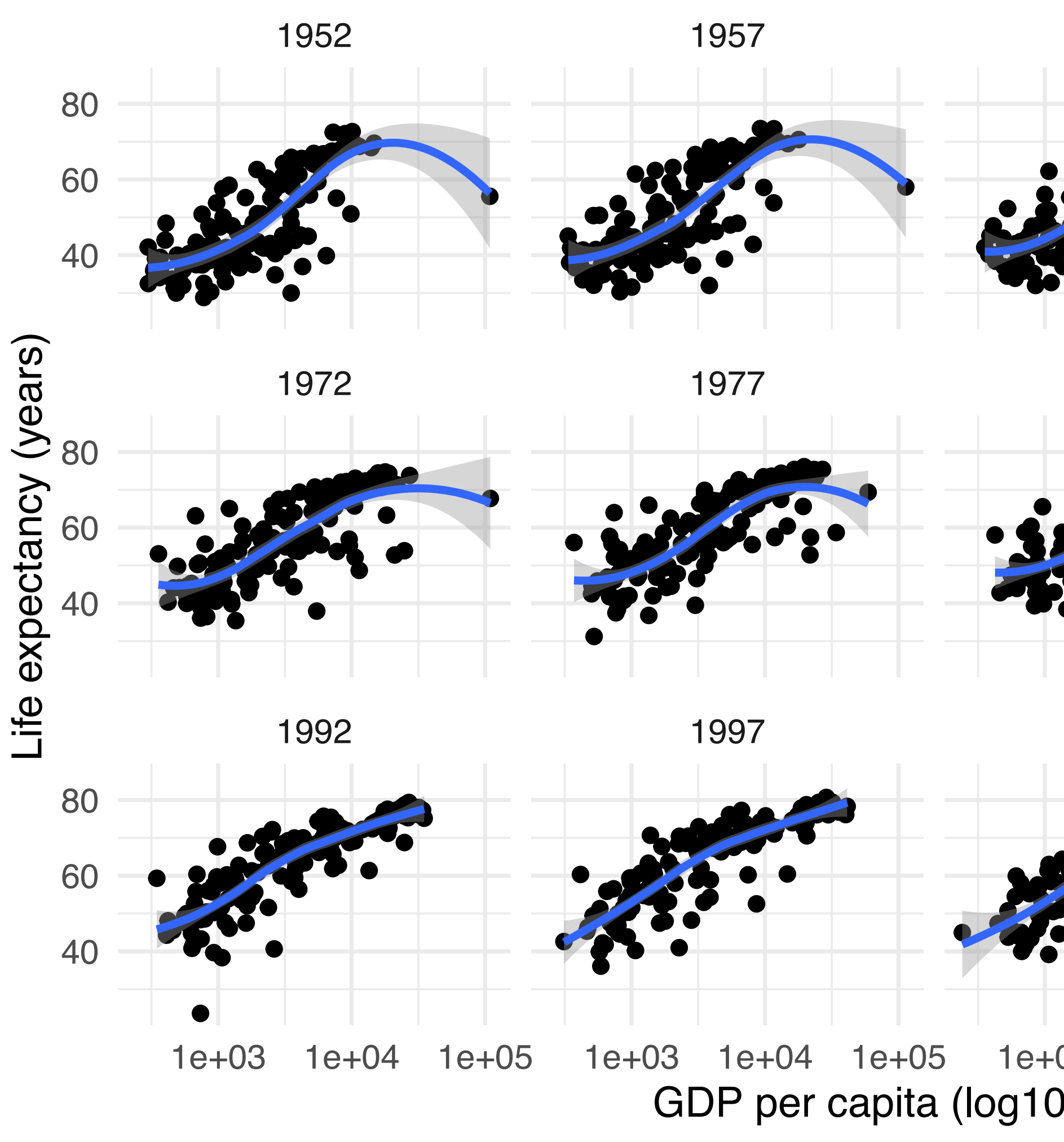
Faceting by one

Fuel efficiency vs Engine size by Year

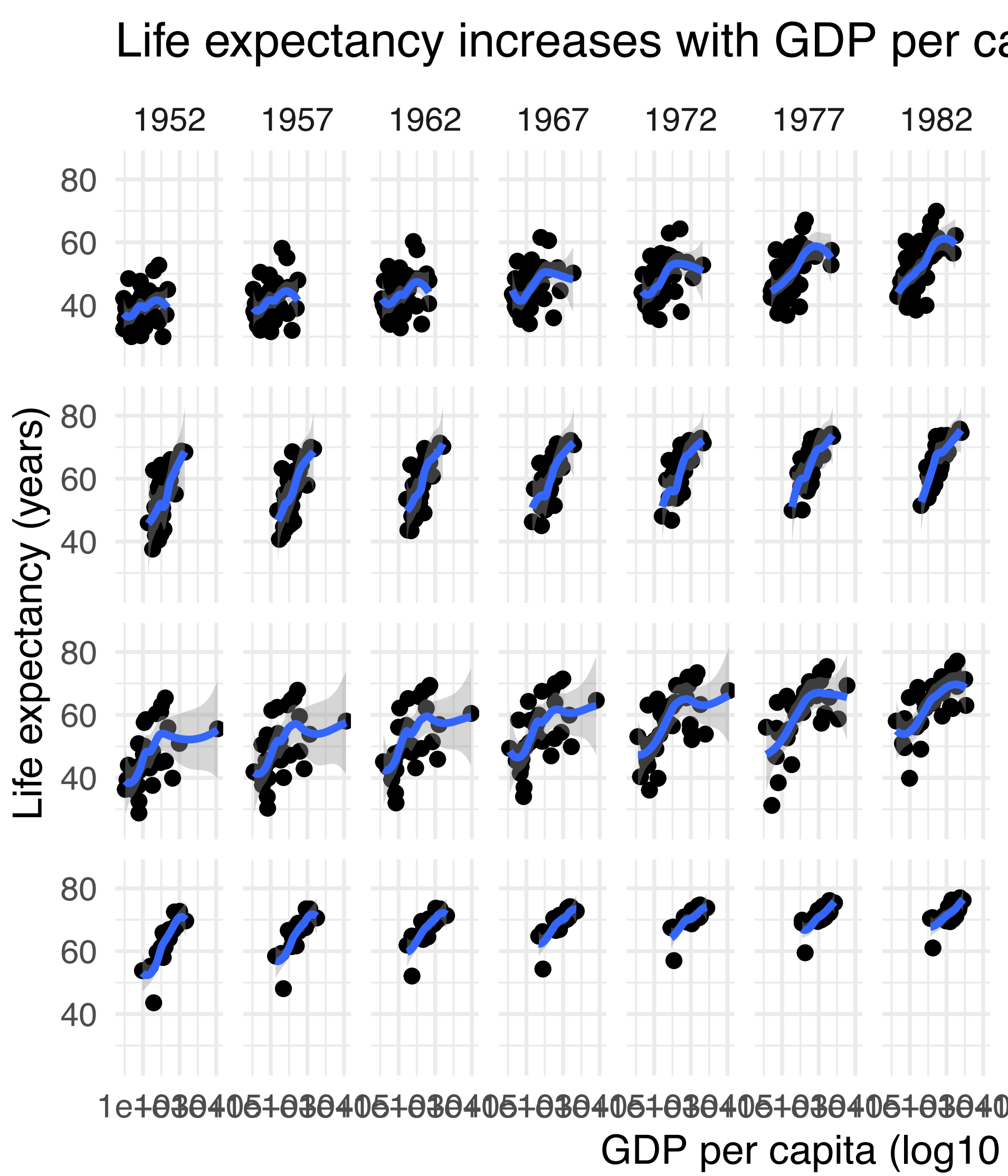


Faceting by one

Life expectancy increases with GDP per capita



Faceting by two v



A GRAMMAR OF

How do we plot?

- By using the “name” of `plot()`
 - ◆ Scatter plot
 - ◆ Box plot
 - ◆ Histogram
- Using “base” R (and similar)
 - ◆ `plot()` - scatter plot
 - ◆ `boxplot()` - box plot
 - ◆ `hist()` - histogram

What are some common ingredien

Recipes for common sta

- Scatter plot
 - ◆ Maps variables to x- and y- axes
 - ◆ Uses points to represent observations
- Line plot
 - ◆ Maps variables to x- and y- axes
 - ◆ Uses lines to connect observations
- Box plot
 - ◆ Maps 5-number summary to x- or y-axis
 - ◆ Uses boxes and whispers to show this
- Histogram
 - ◆ Maps frequency to height
 - ◆ Uses bars to represent frequency
- Bar chart
 - ◆ Maps frequency to height
 - ◆ Uses bars to represent frequency
- Pie chart
 - ◆ Maps proportion to area
 - ◆ Uses slices to represent proportion

Key ingredients for statistical graphics

- Some kind of data
- Encodings from data to visual properties
 - ◆ Marks (“**geometric objects**”, e.g., points, lines, areas)
 - ◆ Channels (“**aesthetics**”, e.g., color, size, shape)
- Statistical transformation
- Coordinate system
- Scales and annotations

Building a p

Consider a simple dataset:

<i>A</i>	<i>B</i>	<i>C</i>
2	3	4
1	2	1
4	5	15
9	10	80

<http://vita.had.co.nz/papers/layered-plot.pdf>

How do we create a scatter plot of A versus C

Wh

Building a p

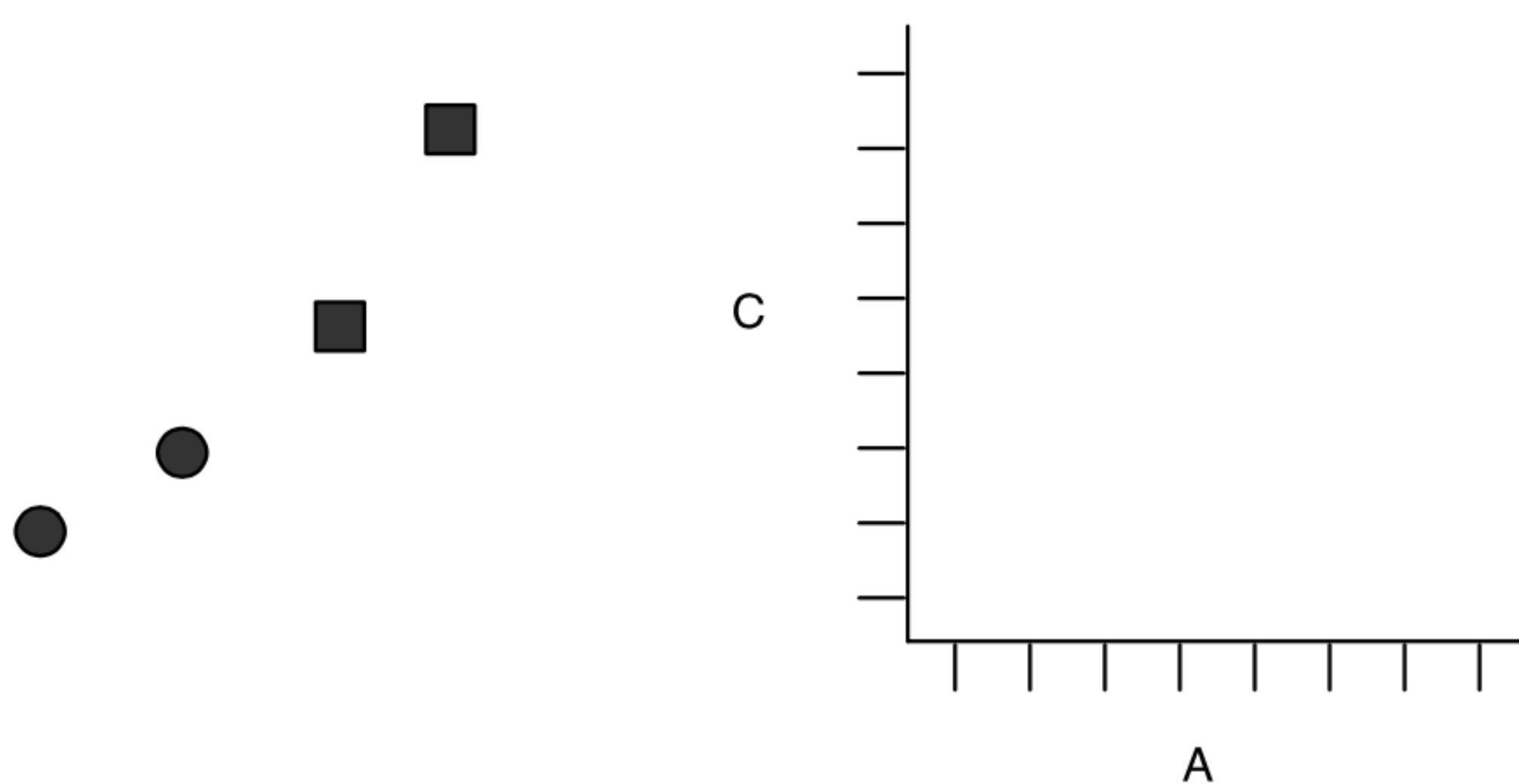
We map the x -axis to **A**, the y -axis to **C**, and shape to **D**

x	y	S
2	4	cin
1	1	cin
4	15	squ
9	80	squ

<http://vita.had.co.nz/papers/layered-plotting.html>

Building a p

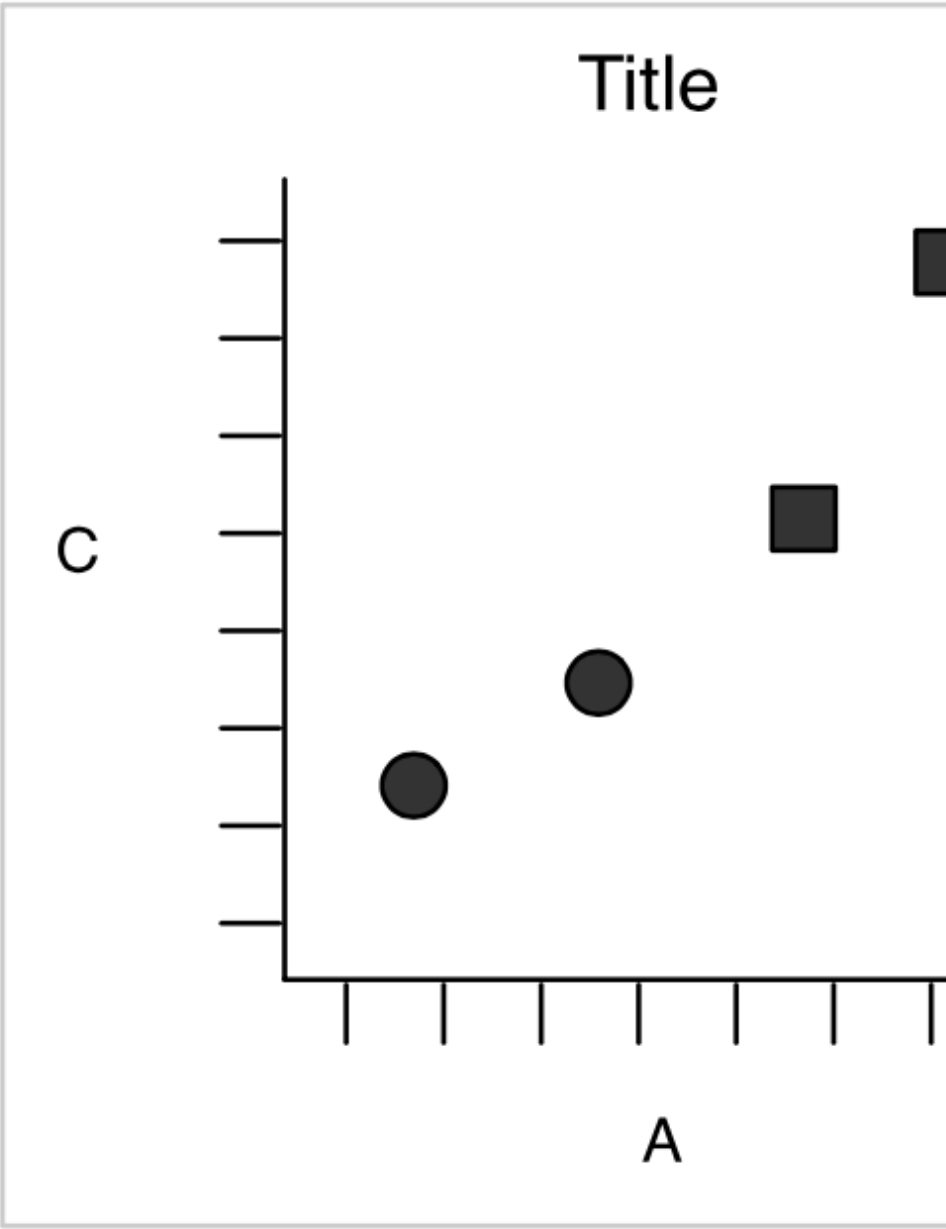
We have (1) marks or **geometric objects**, (2) **scales** and a c



<http://vita.had.co.nz/papers/layered-plotting.html>

Building a plot

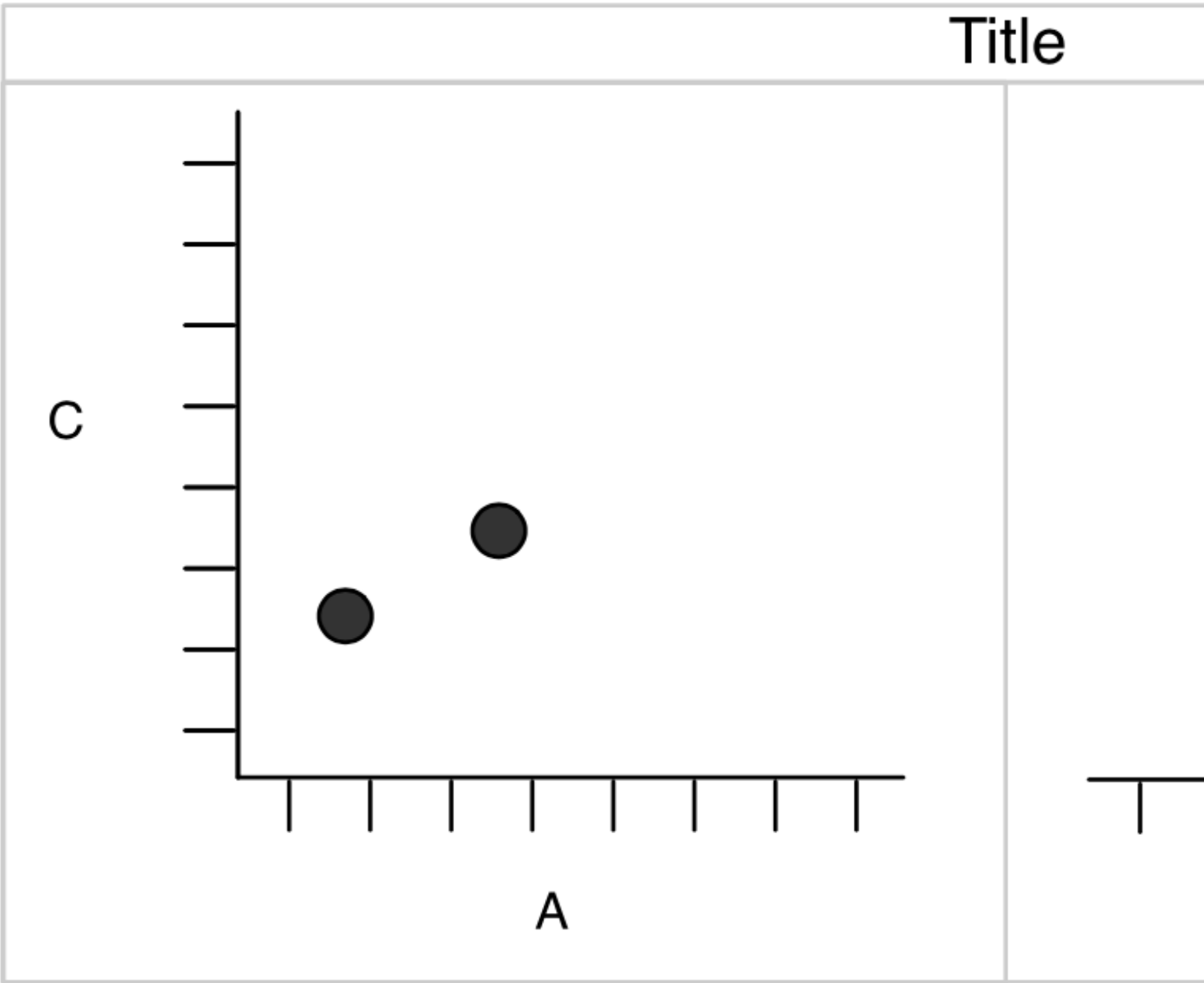
Putting the ingredients together, we have a plot:



<http://vita.had.co.nz/papers/layered-plot.html>

Building a p

If we want to compare the relationship between **A** and



<http://vita.had.co.nz/papers/layered-plotting.html>

Faceting splits the data into subsets and crea

Building a vocabulary

We can build more complicated plots by adding to

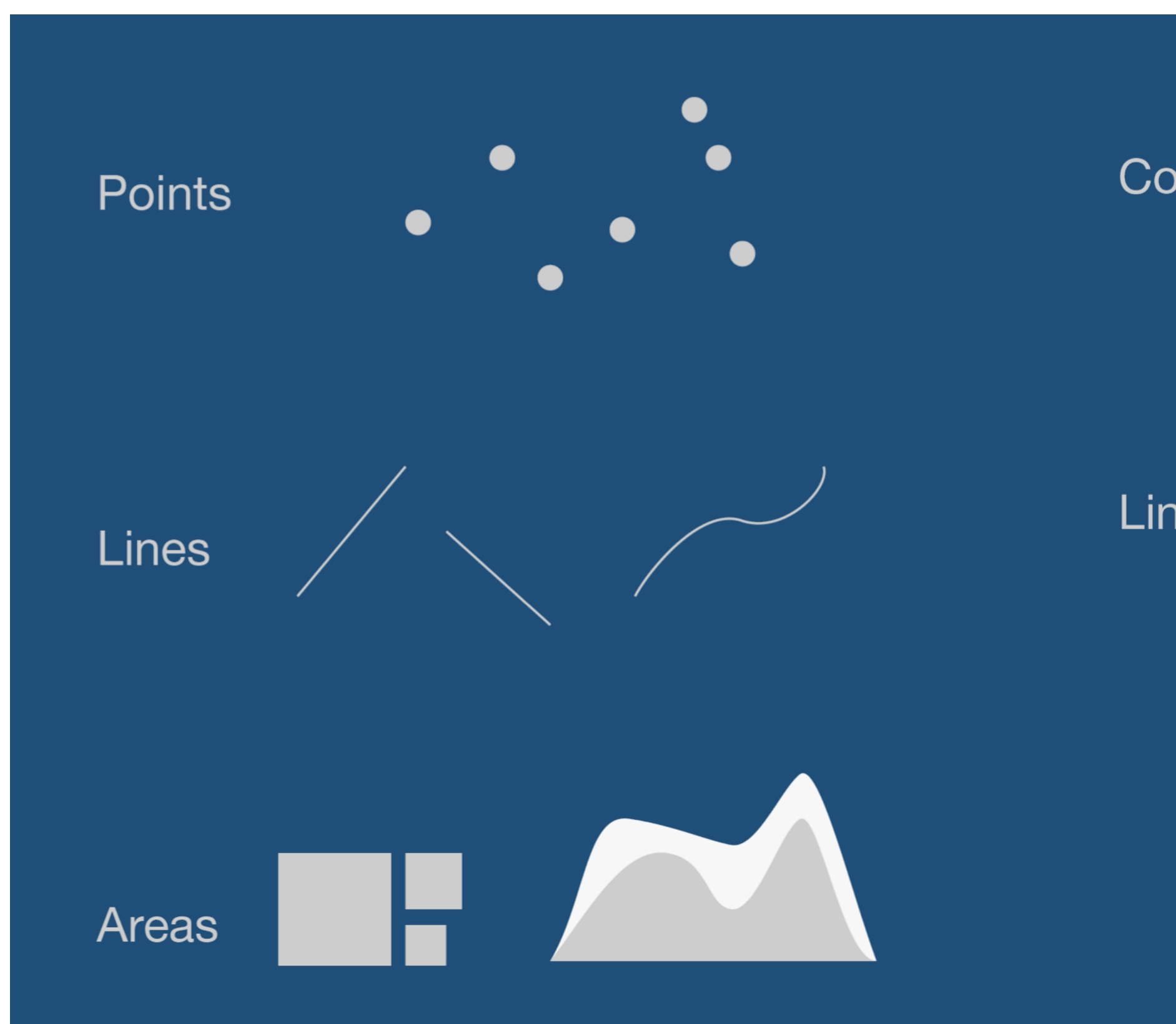
- Layers to overlay plots on top of each other
- Multiple datasets on the same plot
- Apply statistical transformations
- Apply position adjustments (faceting)
- *A way to build such plots programmatically*

A layered grammar

- Default dataset
- Default set of mappings from variables to aesthetics
- One or more layers, each having:
 - ◆ Mark, or geometric object
 - ◆ Statistical transformation
 - ◆ Position adjustment
 - ◆ (Optionally) new dataset
 - ◆ (Optionally) new set of aesthetic mappings
- Scale for each mapped aesthetic
- Facet specification

Visual encodings

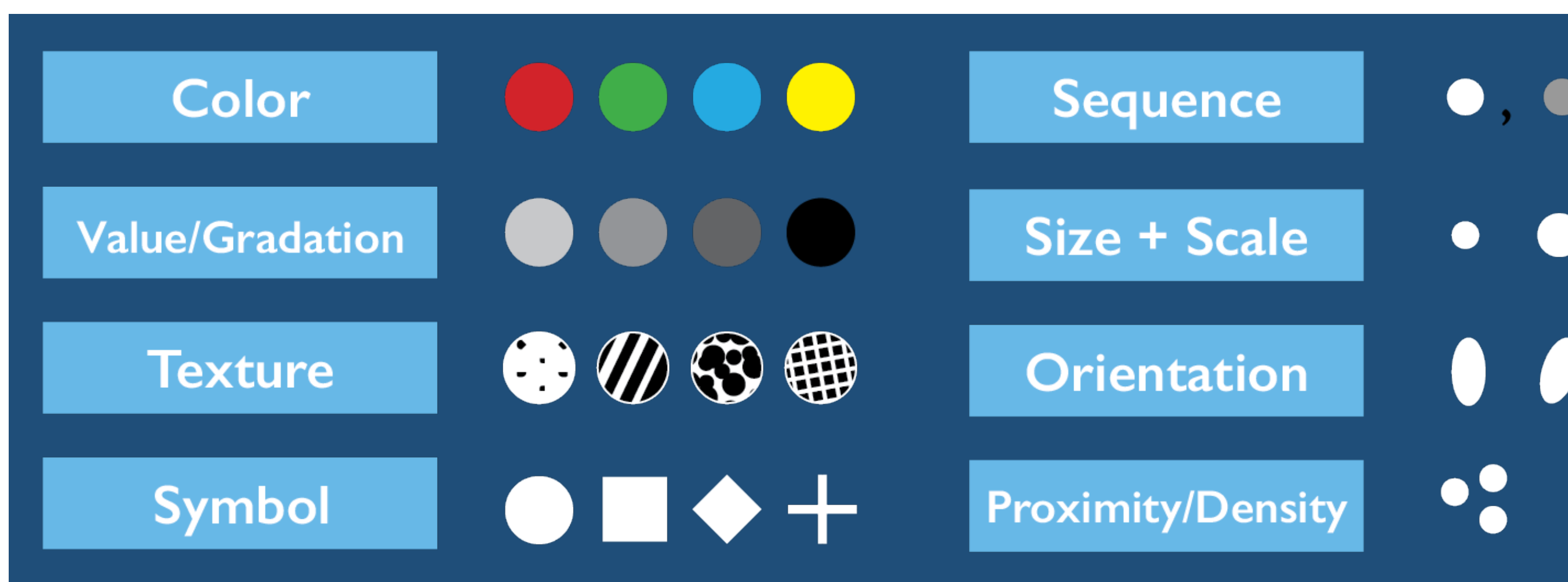
Marks, or **geometric objects**, of



Courtesy of Steven Braun, CAMD

Visual encodings: Channels

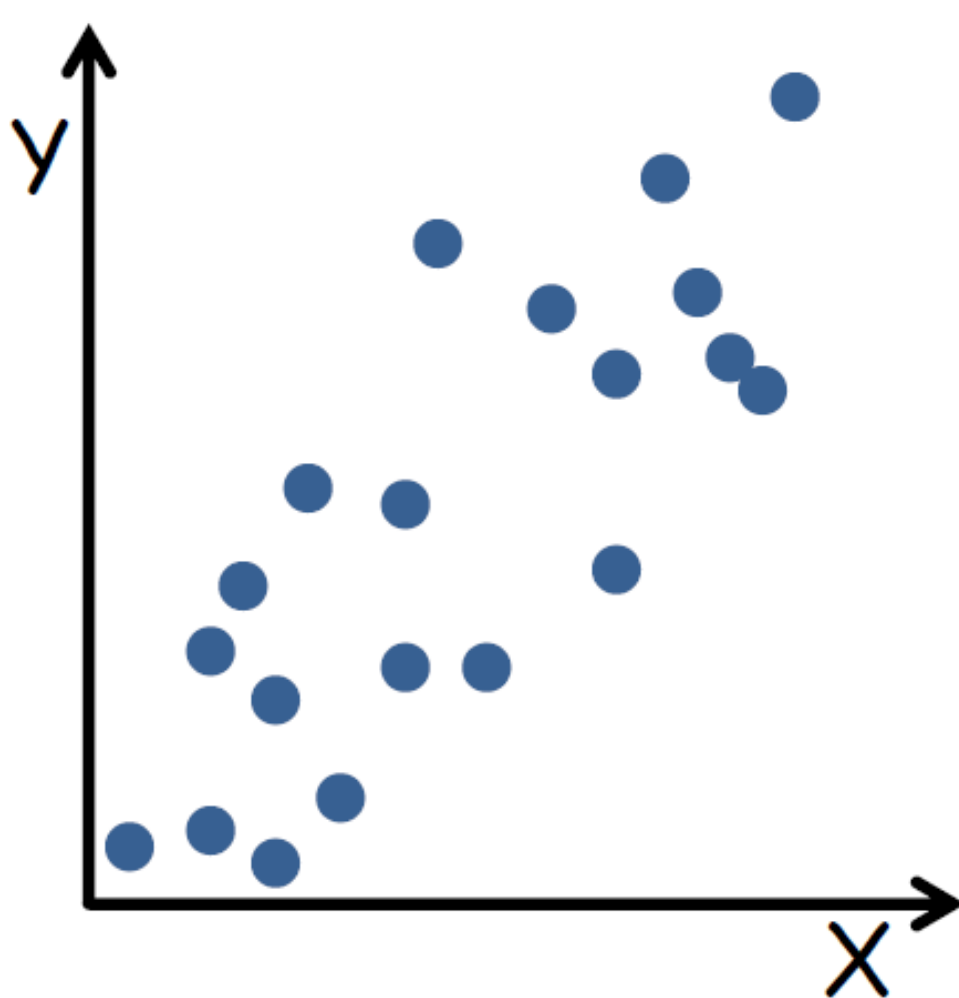
Channels, or **aesthetics** + **scales**



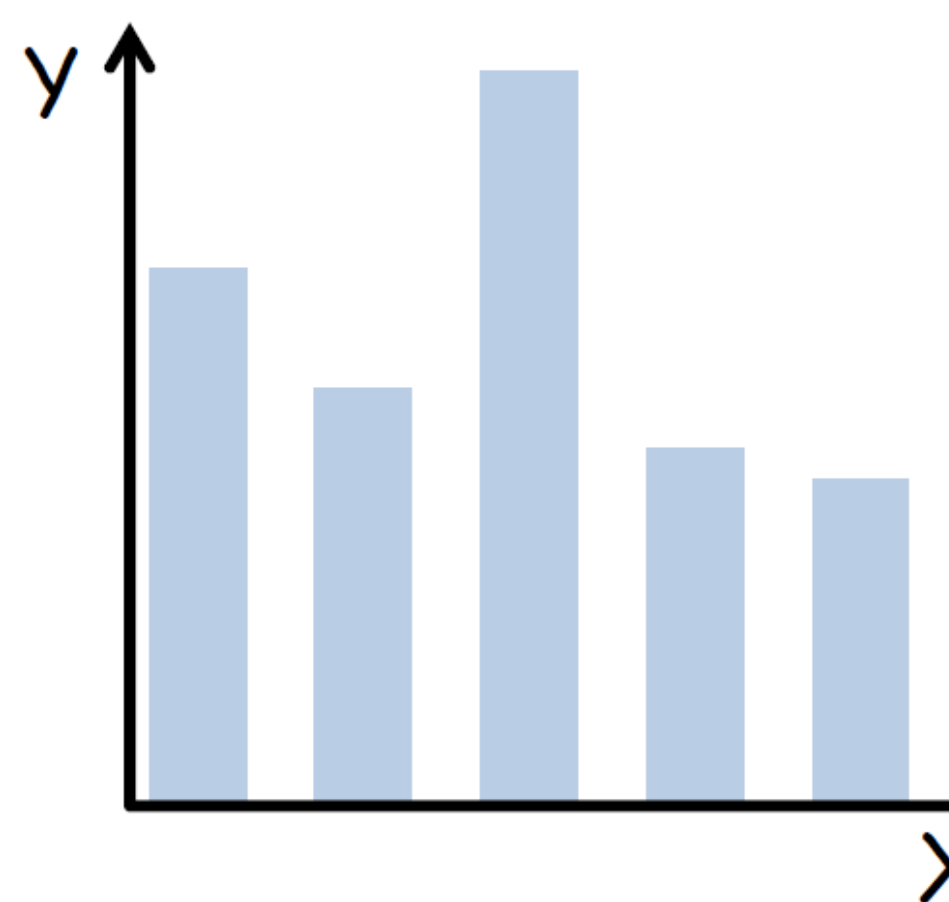
Courtesy of Steven Braun, CAMD

Choosing visual e

Your choice of **geometric**
aesthetic mappings, and **scales** o



Marks: points
Channels: position



Marks: lines
Channels: length, po

Courtesy of Steven Braun, CAMD

Statistical transform

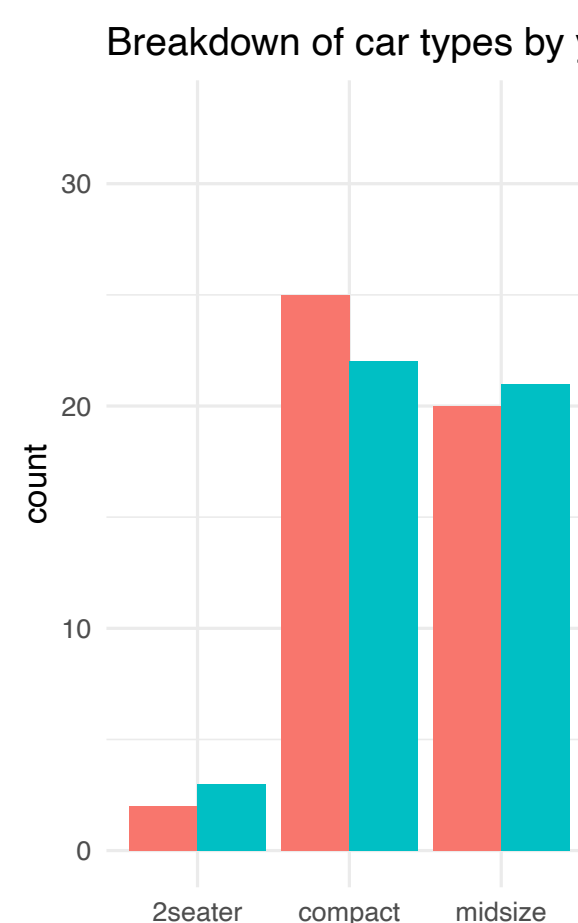
Many statistical graphics utilize **stat**

- Box plot
 - ◆ Five-number summary + outlier
- Histogram
 - ◆ Binning
- Bar plot
 - ◆ Counting

Position adjust

Many statistical graphics require

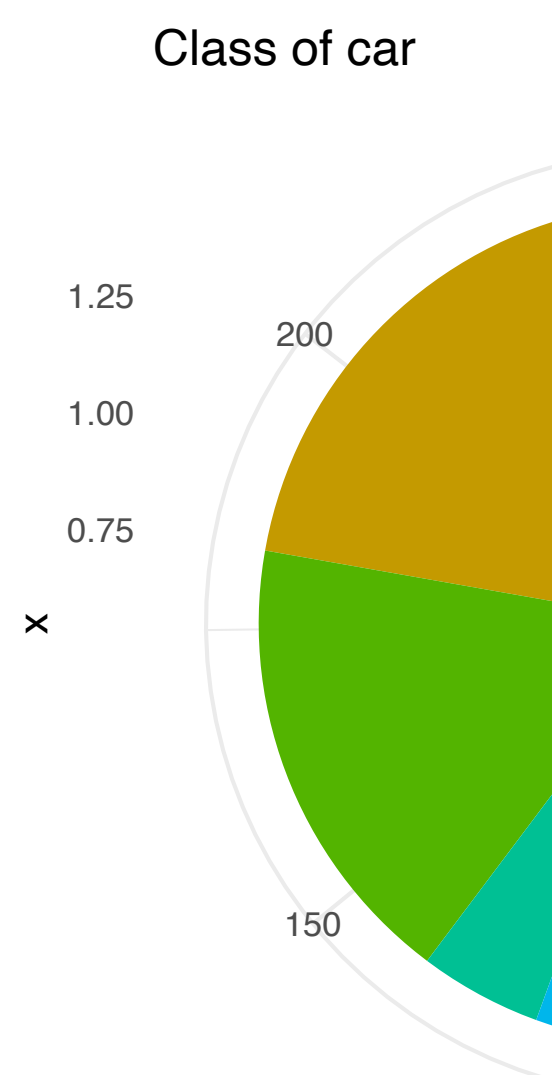
- Scatter plot
 - ◆ Jitter
- Bar plot
 - ◆ Dodge
 - ◆ Stack



Coordinate systems

Some graphics may require different

- Cartesian
- Polar
- Map



Implementing a grammar

A version of the “layered grammar of graphics”



```
ggplot(data = <D T SET>,
       mapping = es(<M PPINGS>
  layer(geom = <GEOM>,
       stat = <ST T>,
       position = <POSITION>)
  <SCALE_FUNCTION>() +
  <COORDINATE_FUNCTION>() +
  <Facet_FUNCTION>()
```

or, more simply

```
ggplot(data = <D T SET>,
       mapping = es(<M PPINGS>
  <GEOM_FUNCTION>()
```

Recipes for common sta

- Scatter plot
 - ◆ Geom = “point”
 - ◆ Stat = “identity”
- Line plot
 - ◆ Geom = “line”
 - ◆ Stat = “identity”
- Box plot
 - ◆ Geom = “boxplot”
 - ◆ Stat = “boxplot”

GGPLOT