

# Introduction to R for data analysis

- plots -

Carl Herrmann & Carlos Ramirez  
R4SC - Freiburg June 2024



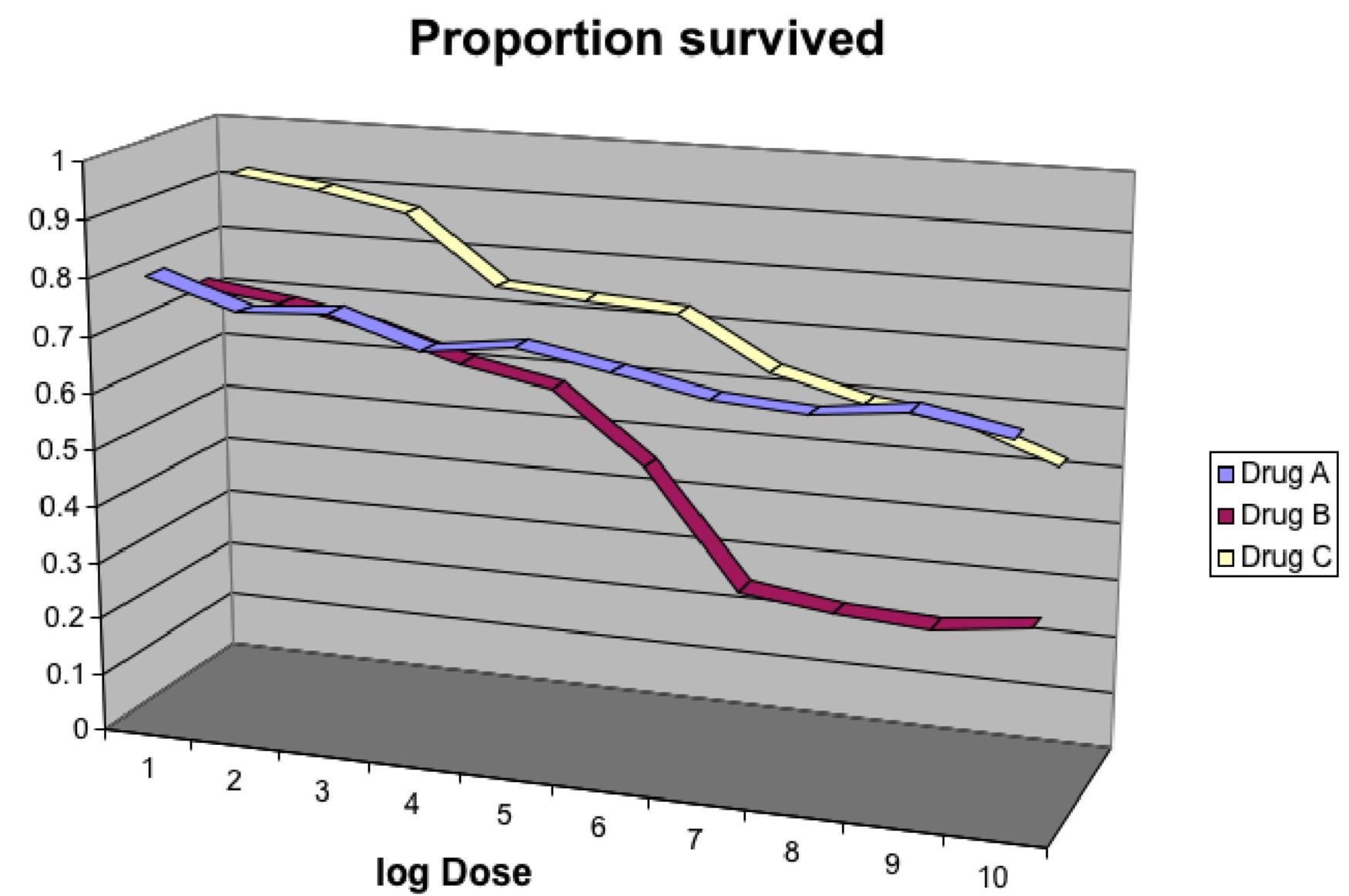
Institut für Pharmazie und  
Molekulare Biotechnologie



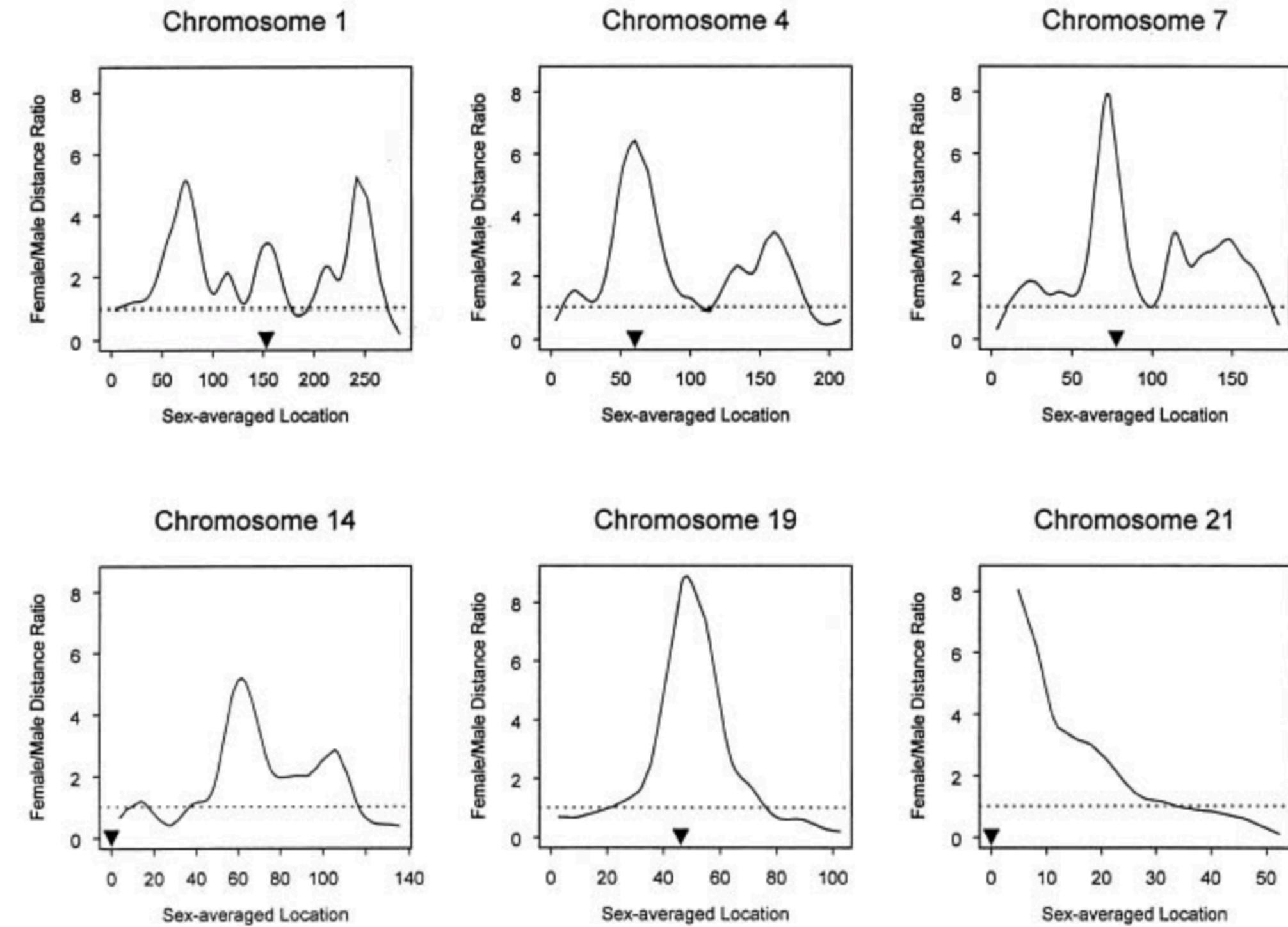
UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

# Graphical representation

- Appropriate graphical representation depends on the type of data
  - categorical
  - counts
  - continuous data
- Aim of good data graphics: **display data accurately and clearly** (Karl Broman  
[https://www.biostat.wisc.edu/~kbroman/topten\\_worstgraphs/](https://www.biostat.wisc.edu/~kbroman/topten_worstgraphs/))
- Bad practice:**
  - as little information as possible
  - make things obscure through inappropriate graphics
  - pseudo 3D
  - poor scales



# Example of bad plot



**Figure 1** Plots of the female:male genetic-distance ratio against sex-averaged genetic location (in cM) along six selected chromosomes. Approximate locations of the centromeres are indicated by the triangles. The dashed lines correspond to equal female and male distances.

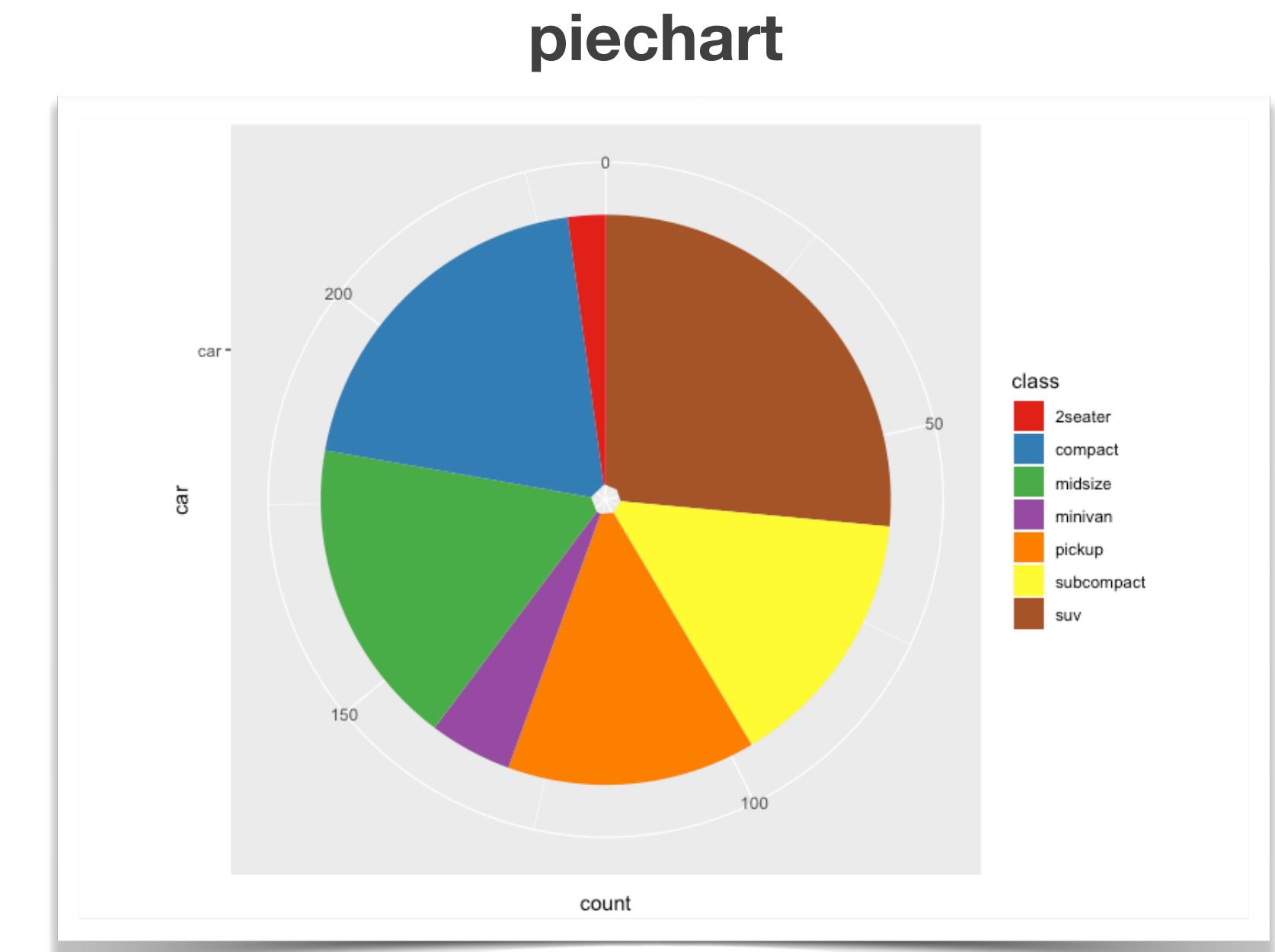
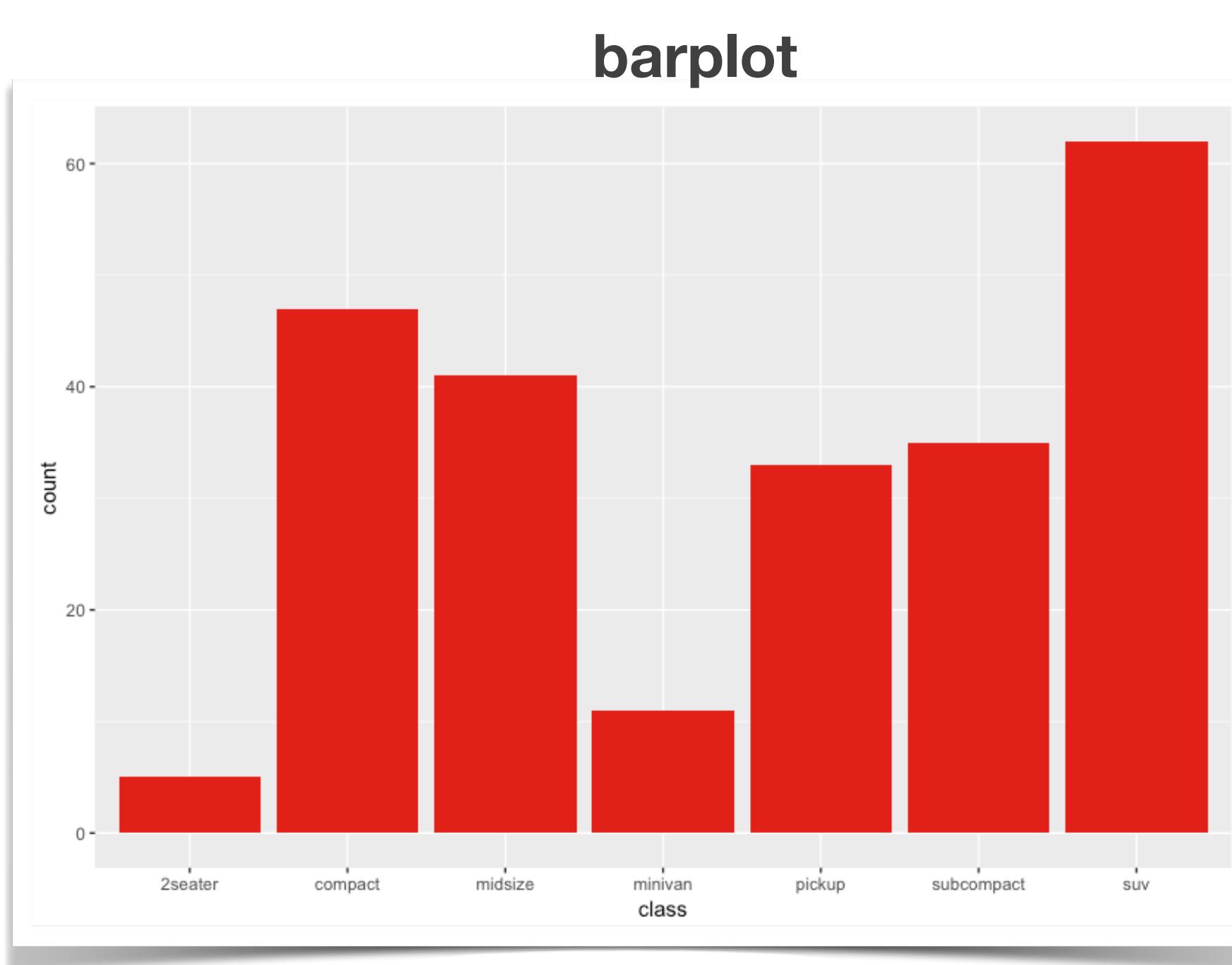
**What's wrong with this plot?**

[https://www.biostat.wisc.edu/~kbroman/topten\\_worstgraphs/](https://www.biostat.wisc.edu/~kbroman/topten_worstgraphs/)

# Categorical Data

## Barplots

- How many instances in each category?
- Only meaningful measure: **MODE** (= category with highest counts)
- Possible plots: **barplots; piecharts**

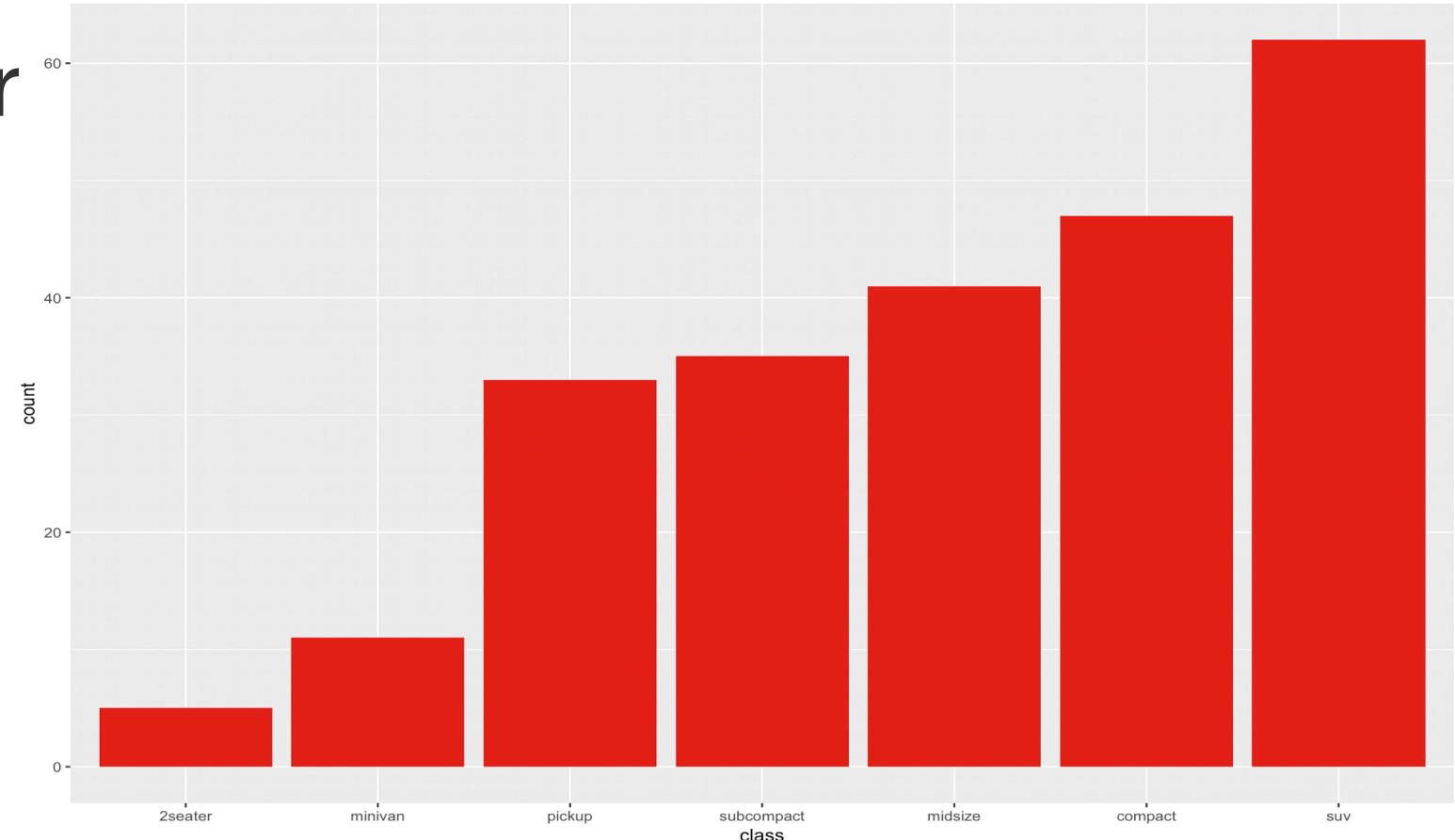


***Avoid piechart : areas are more difficult to judge than length!***

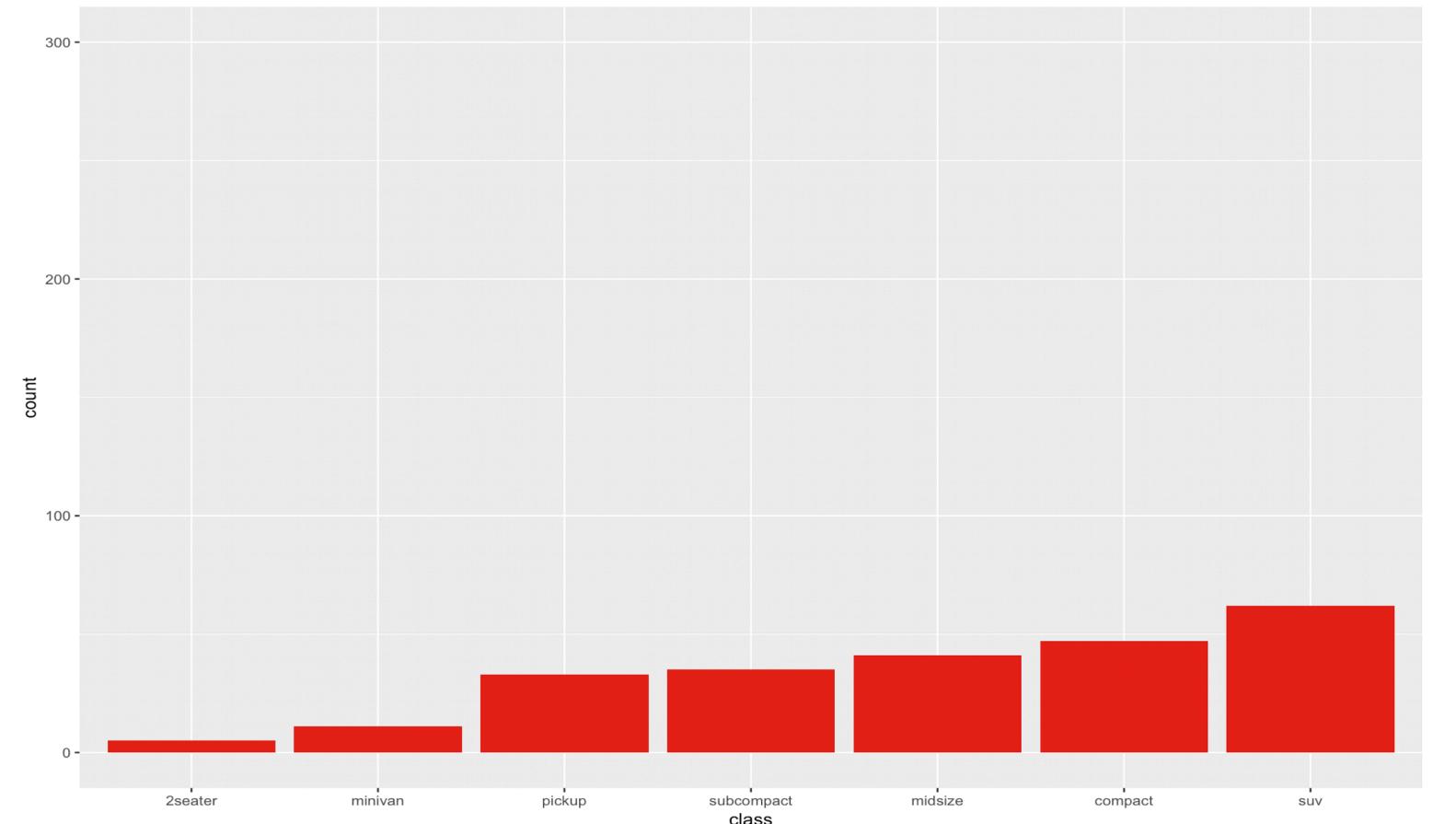
# Categorical Data

## Barplots

- Consider the natural order of categories for nominal data)



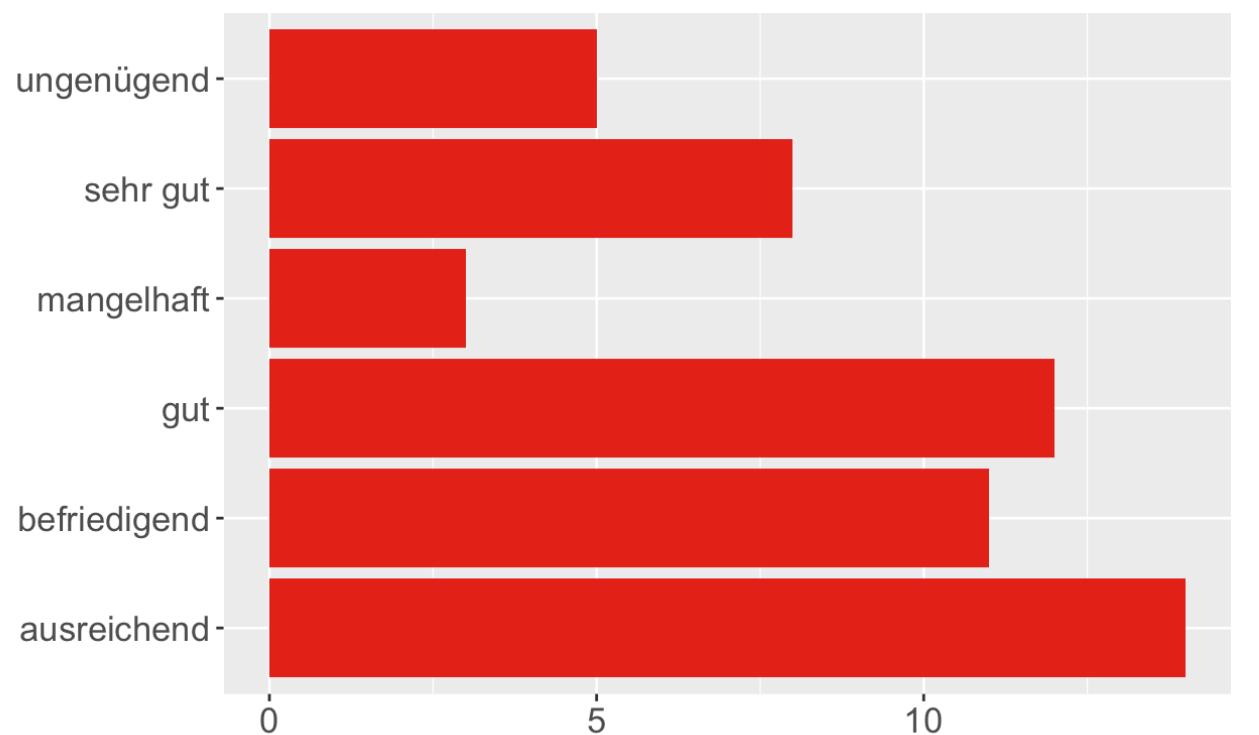
- Beware of selecting the proper scales for plotting!



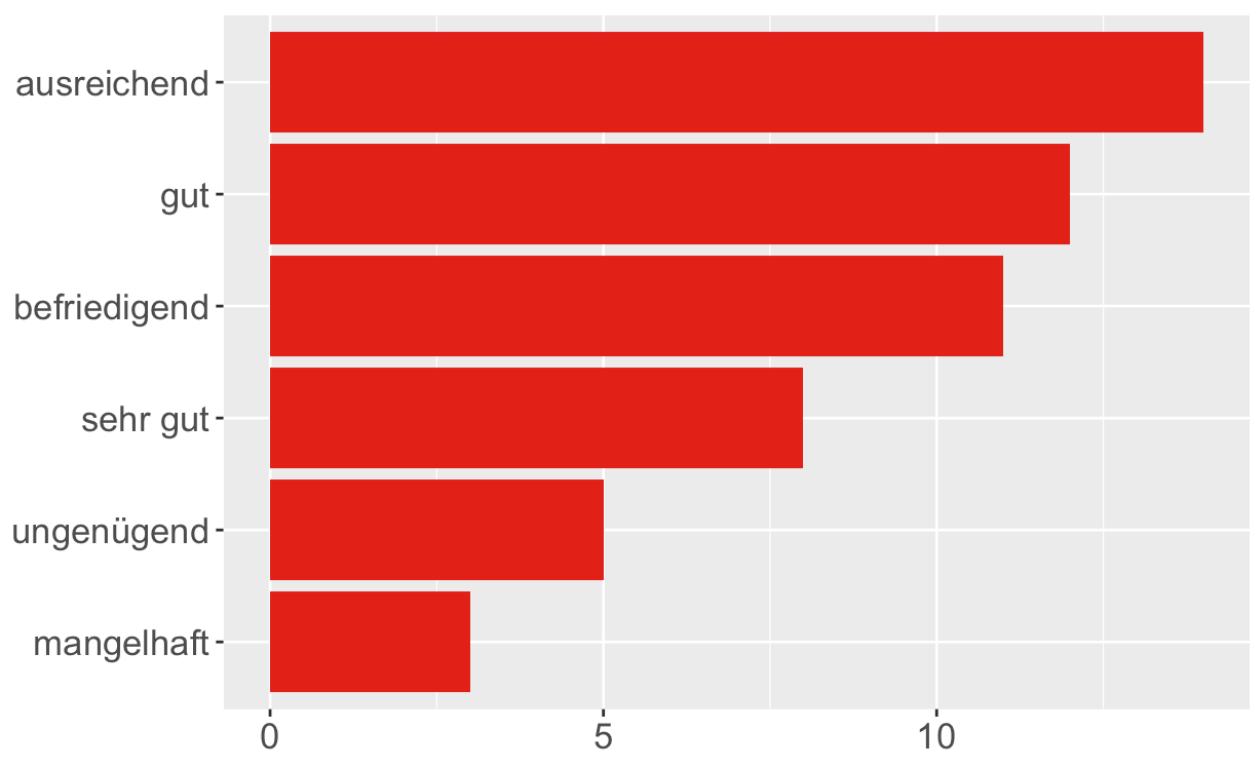
# Categorical Data

## Barplots

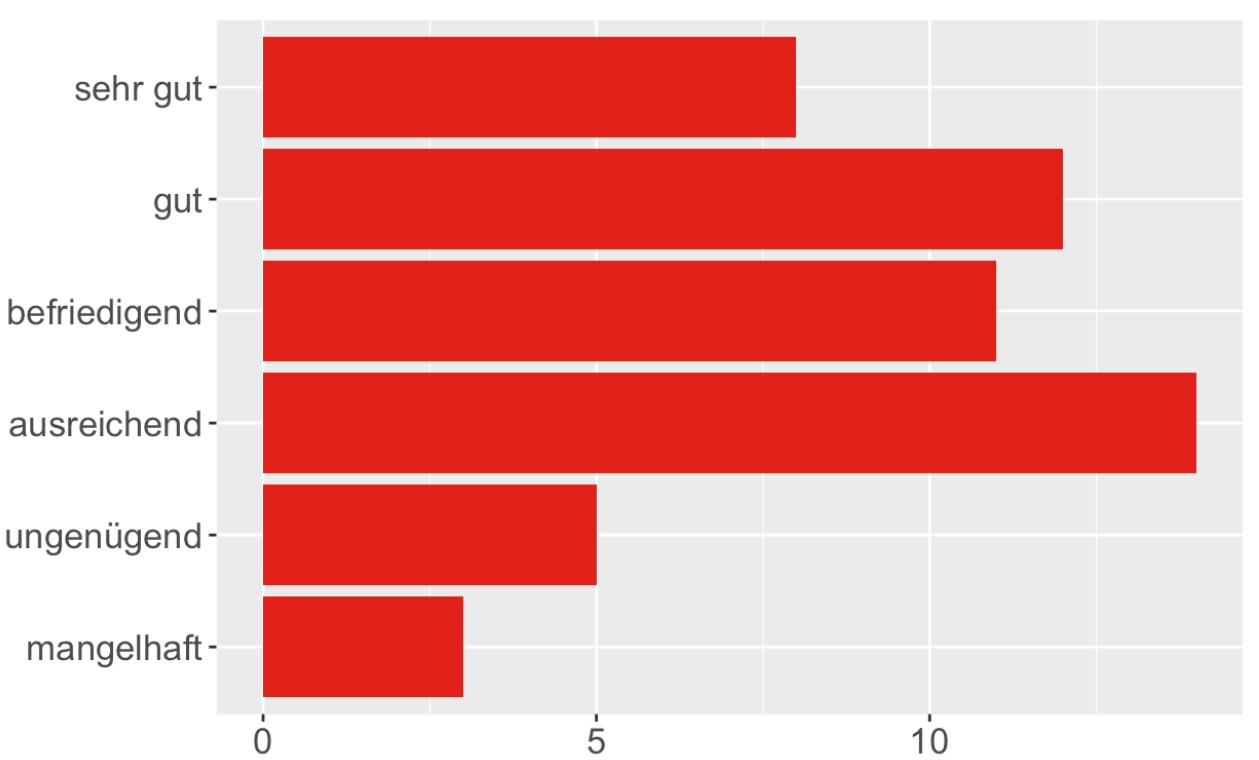
- Random order



- Order by increasing / decreasing counts



- Natural order of ordinal data



Mode

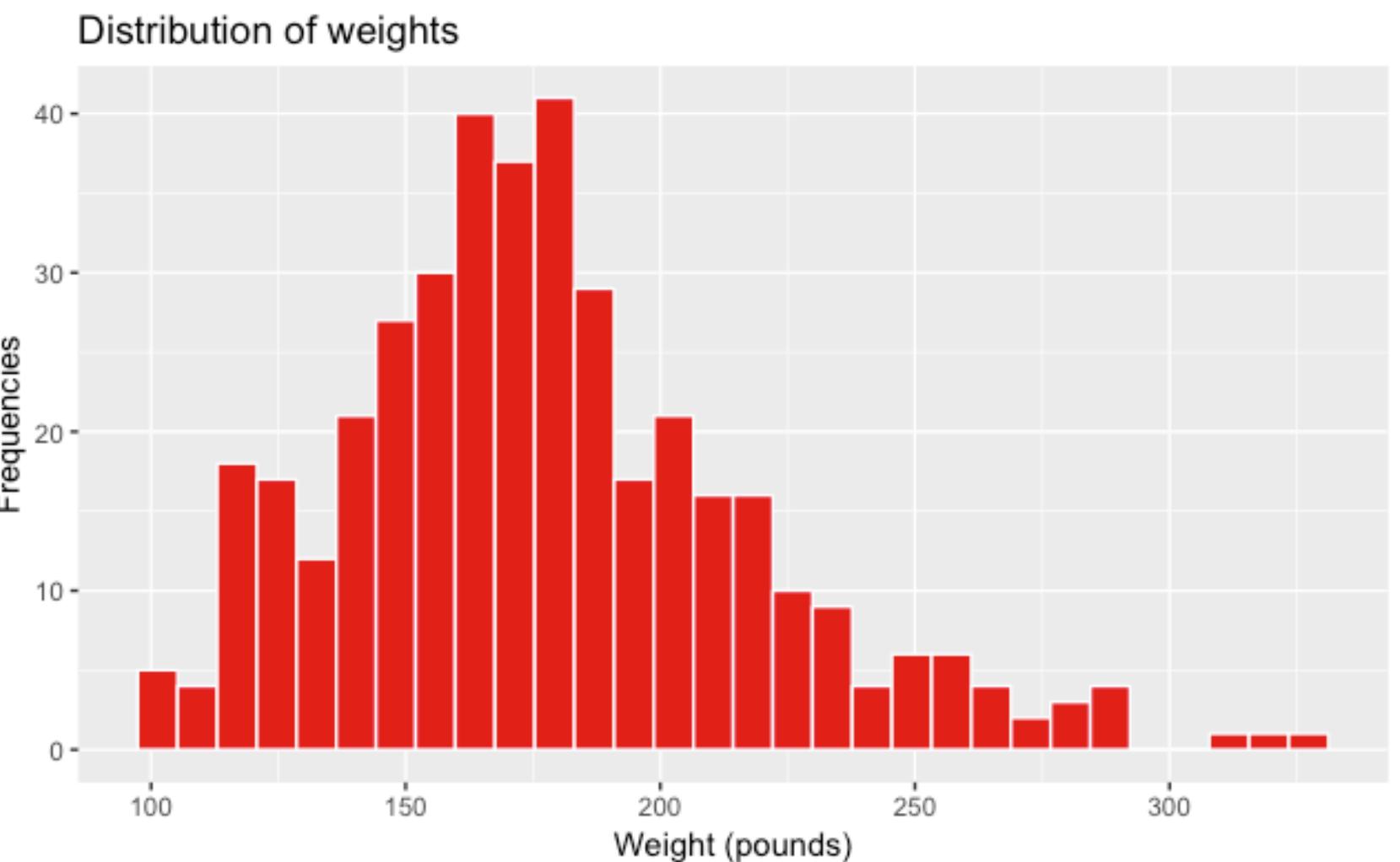
# Numerical variables

- Numerical data are instances of underlying **random variable**
- Random variables  $X$  have
  - **Density distributions**  $p(X)$
  - **Expectation values**  $E(X)$
  - **Variances**  $\text{Var}(X)$

# Numerical data

## Histograms

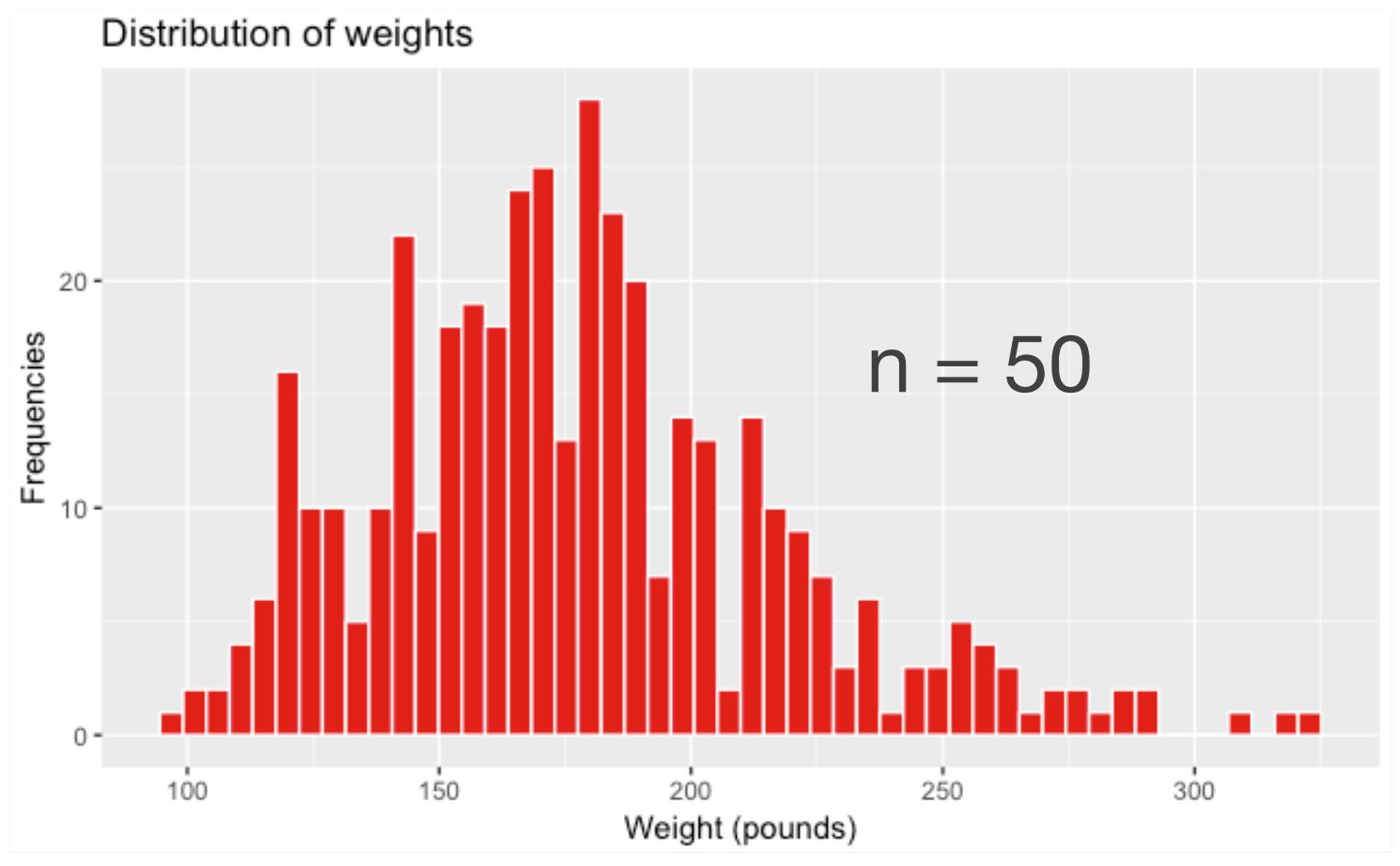
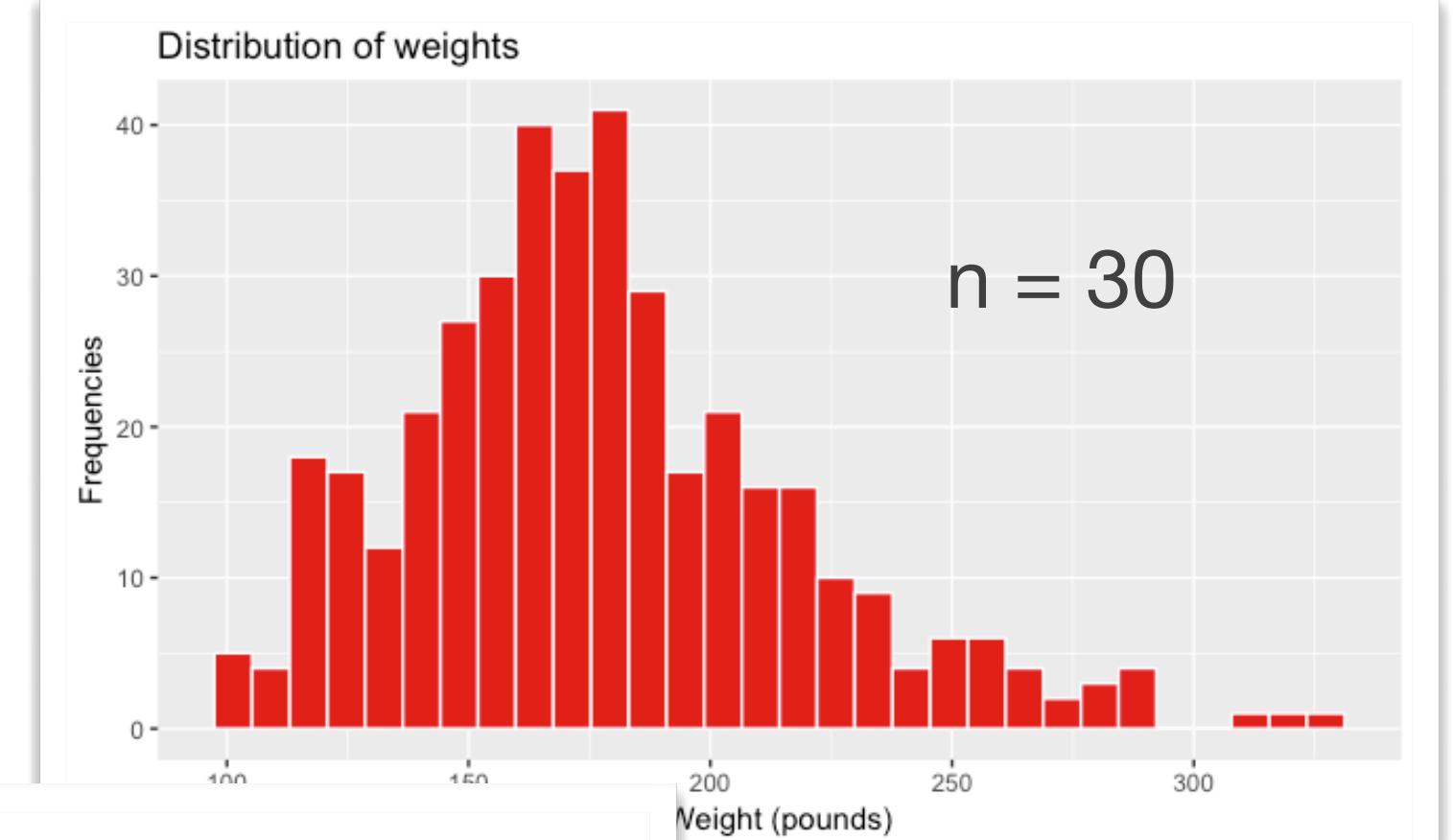
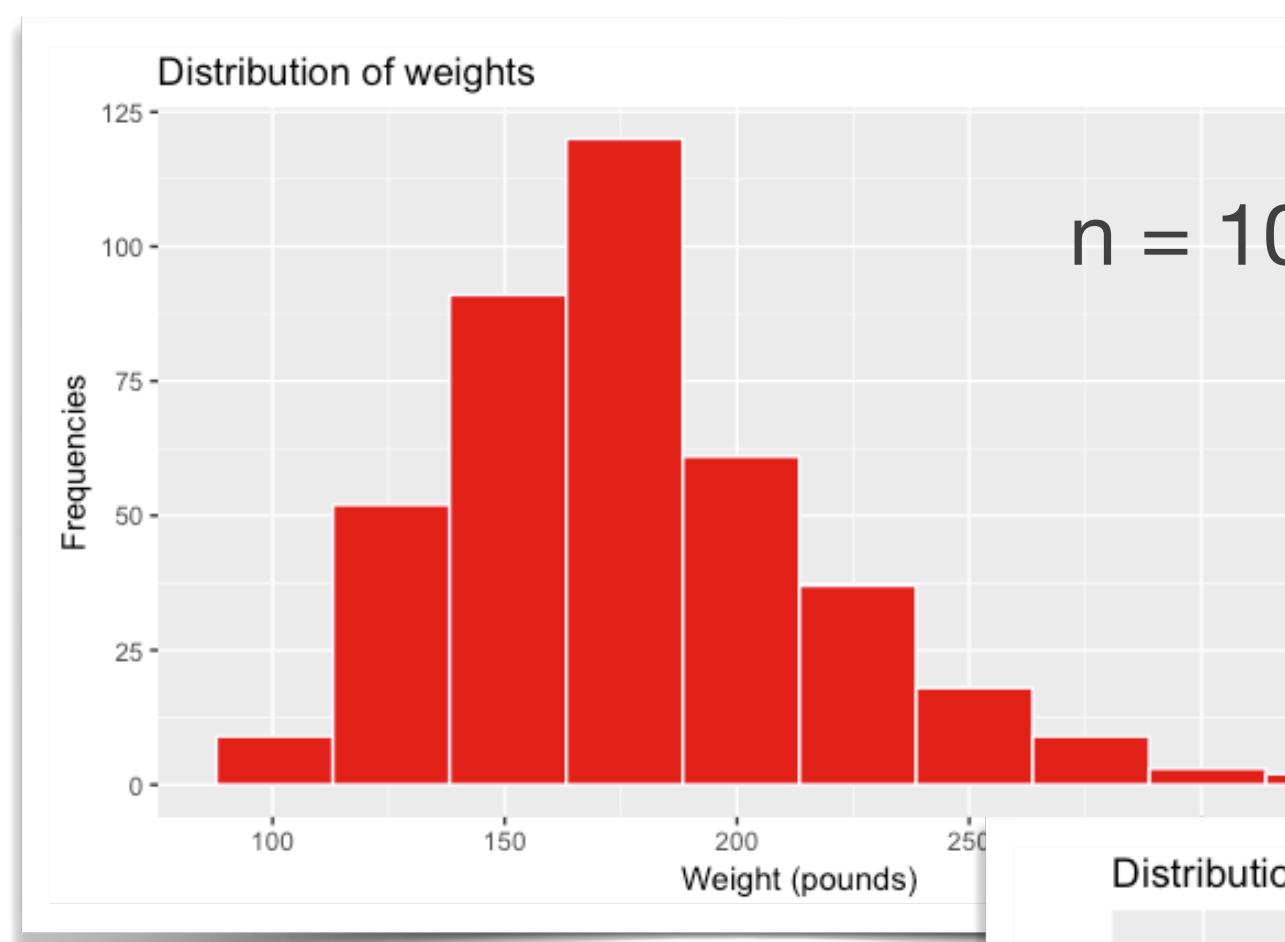
- Categorical data → counts (**barplot**)
- Numerical data → counts within intervals (**histogram**)
  - define **discrete intervals** for numerical variable → **ordinal variable**  
e.g. [0,10), [10,20), [20,30), ...
  - count occurrences within intervals and plot
- histograms represent the **distribution of the variable**



# Numerical data

## Histograms

- Right choice of interval depends on the data type
- Number of bins has a strong impact on appearance of plot!**

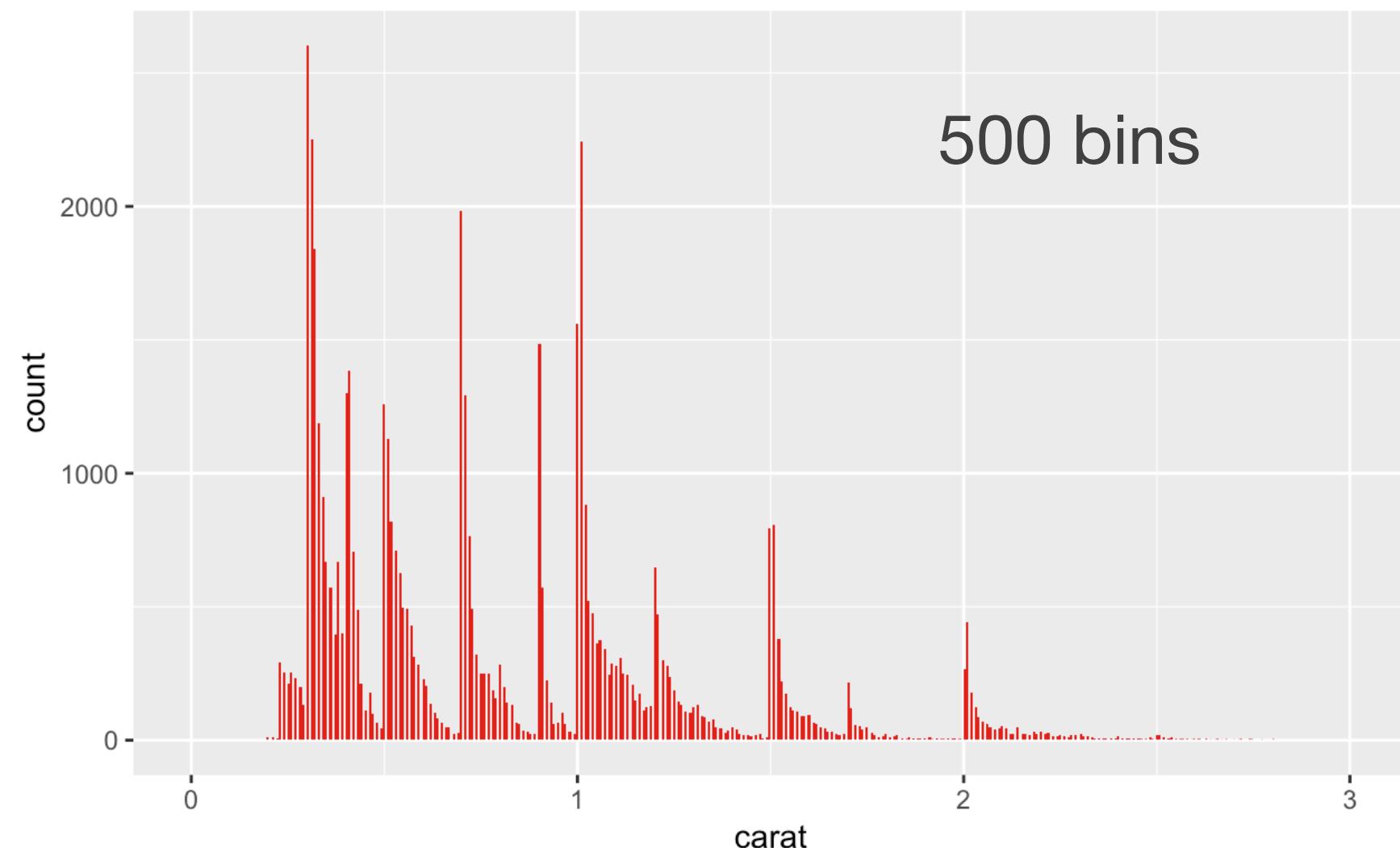
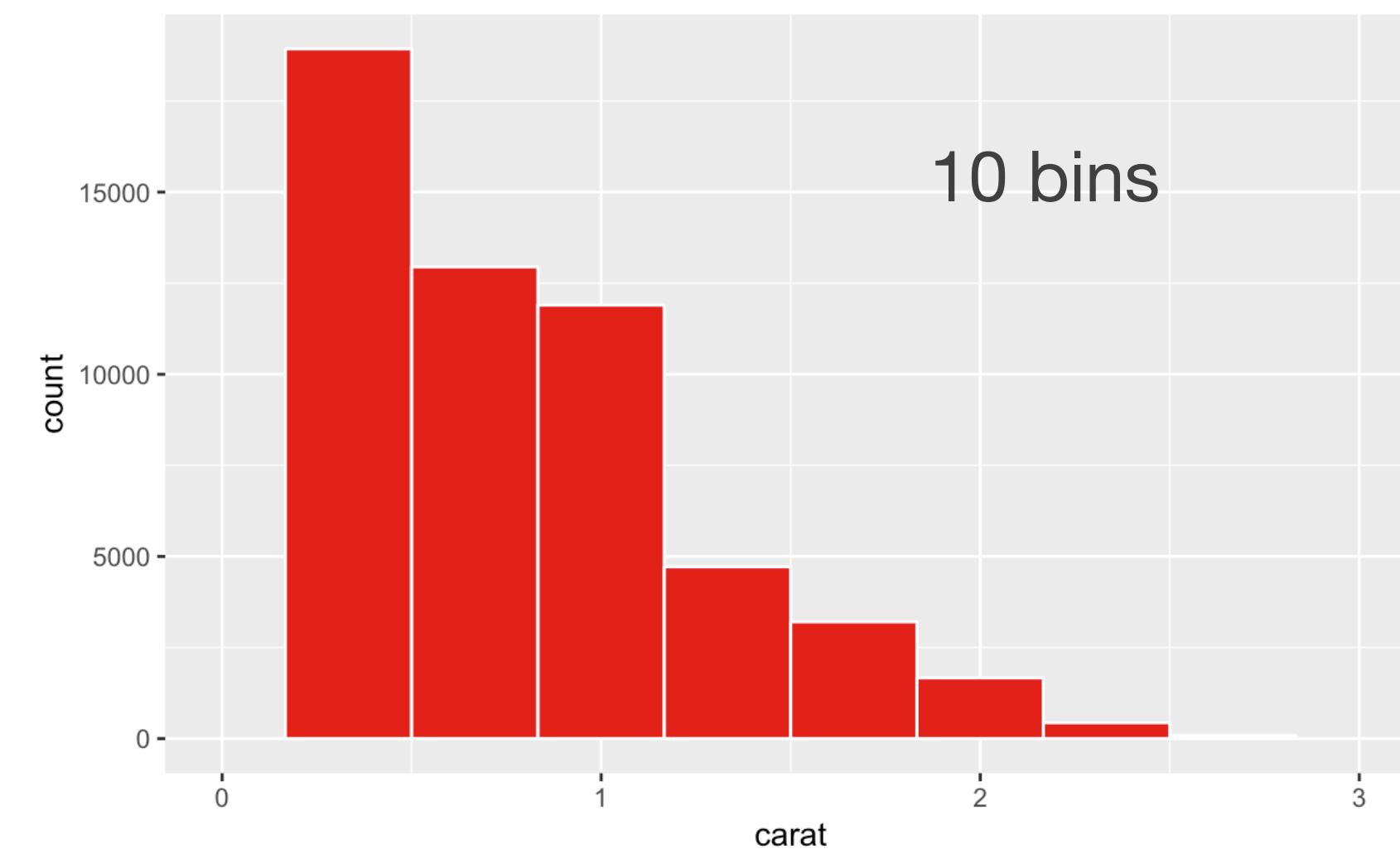


# Numerical data

## Histograms

- pattern becomes visible at high resolution
- peaks around integer values (why?)
- tail on the right of integer values (why?)

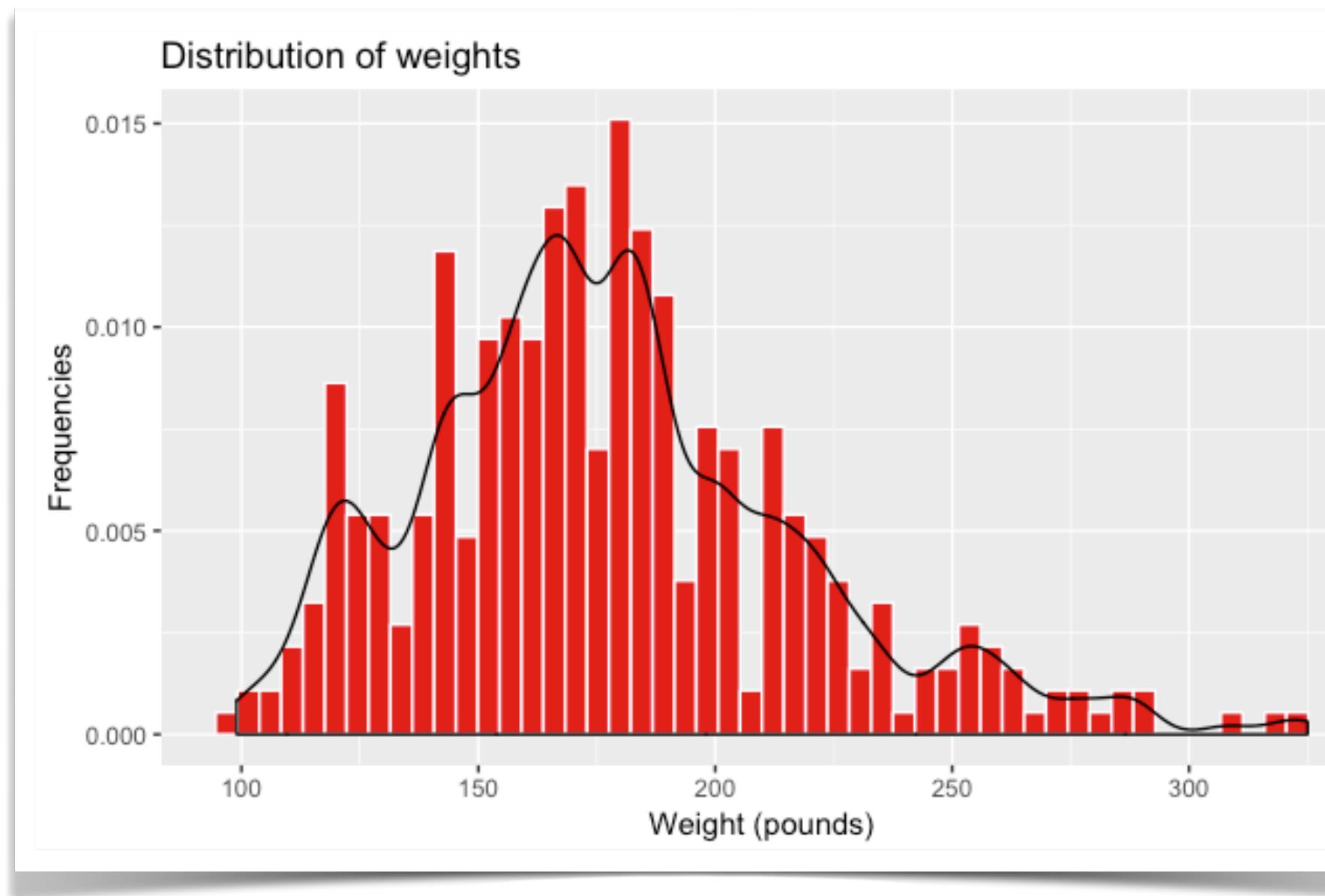
Distribution of carat values for diamonds



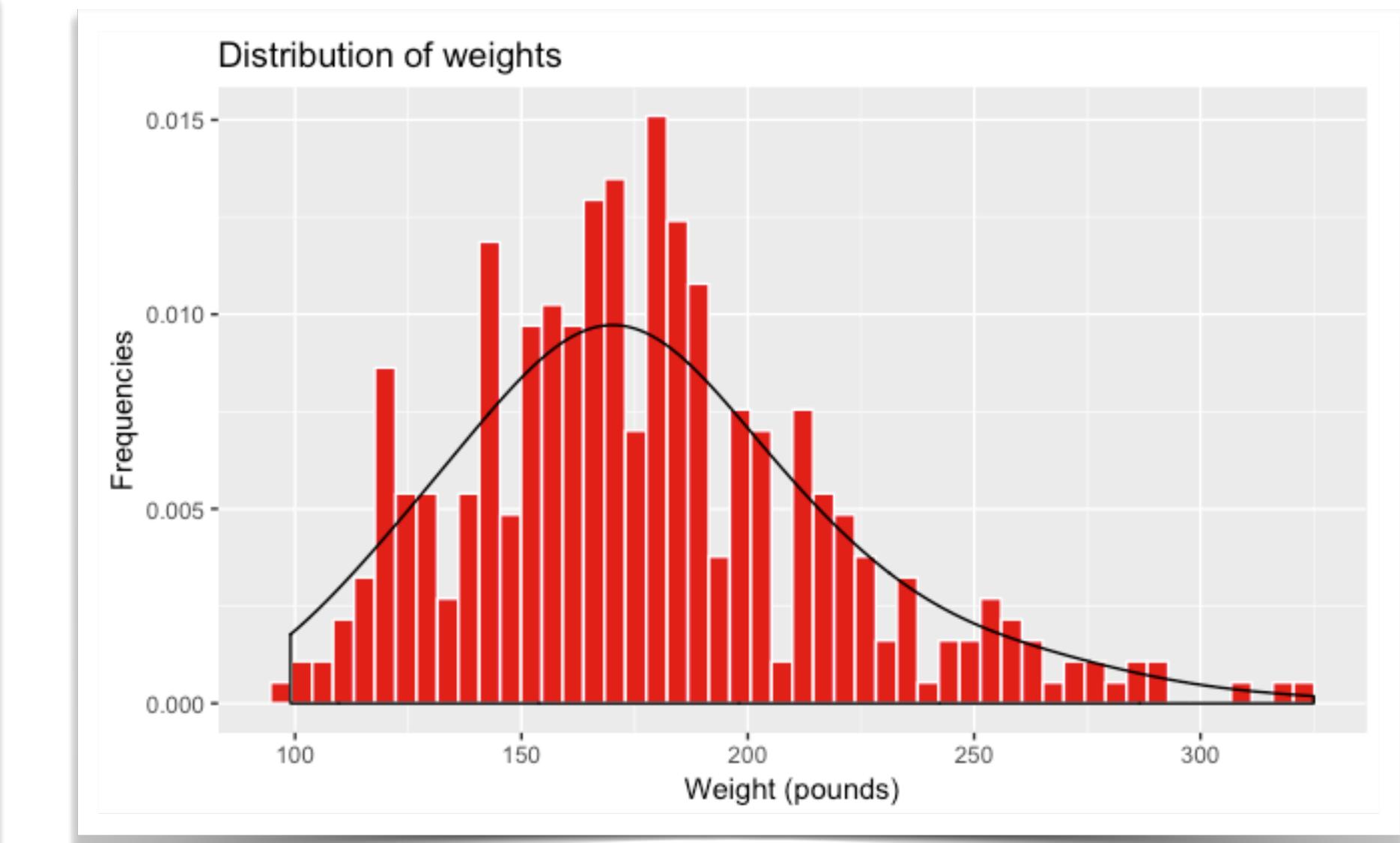
# Numerical data

## Histograms

- Frequency distributions (= histograms) can be shown using a smoothed **density curve**
- Smoothing depends on the bandwidth (~size of the interval over which to smooth)



*bandwidth = 5*

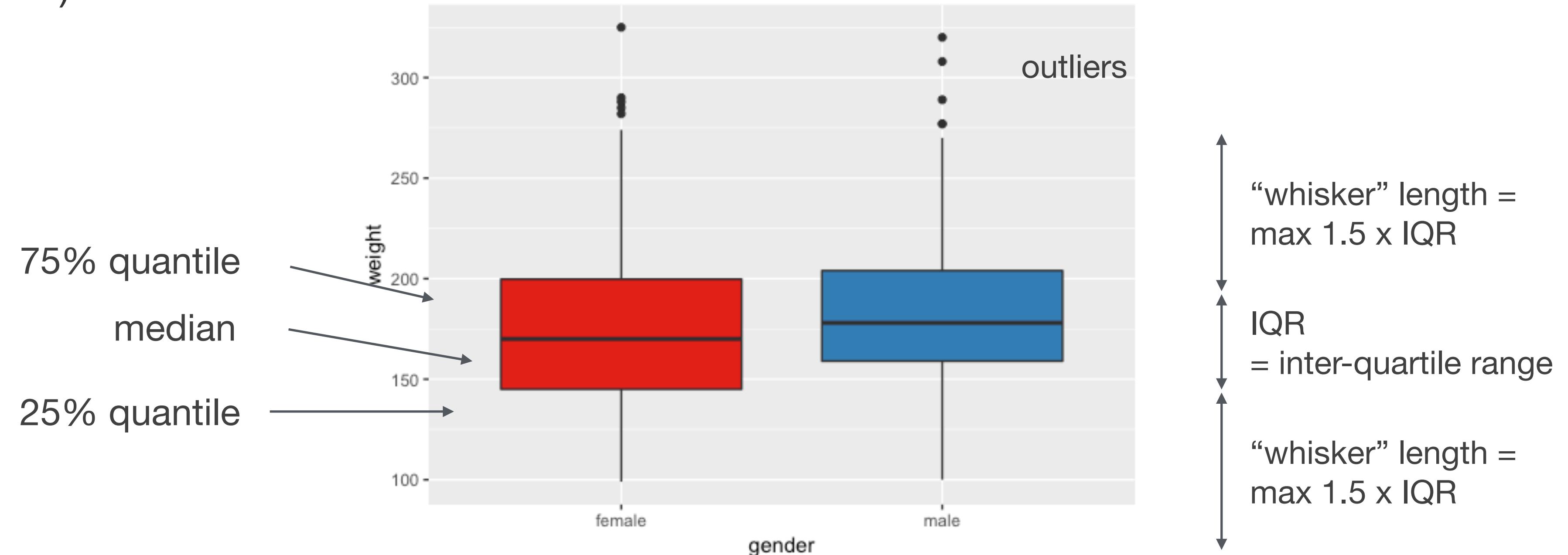


*bandwidth = 20*

# Numerical values

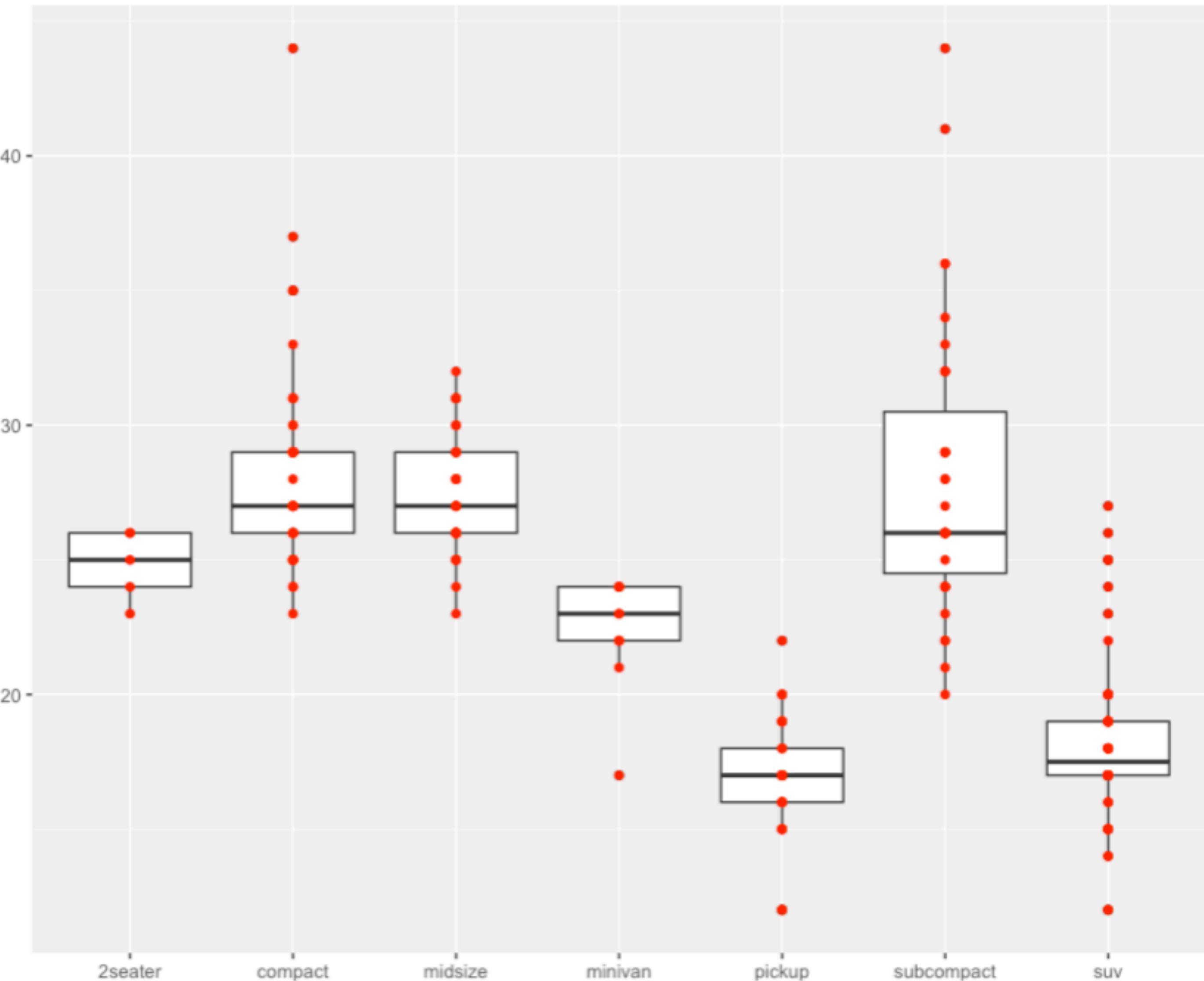
## boxplots

- Boxplot give an indication on the shape of the distribution (median / symmetry / outlier / ...)



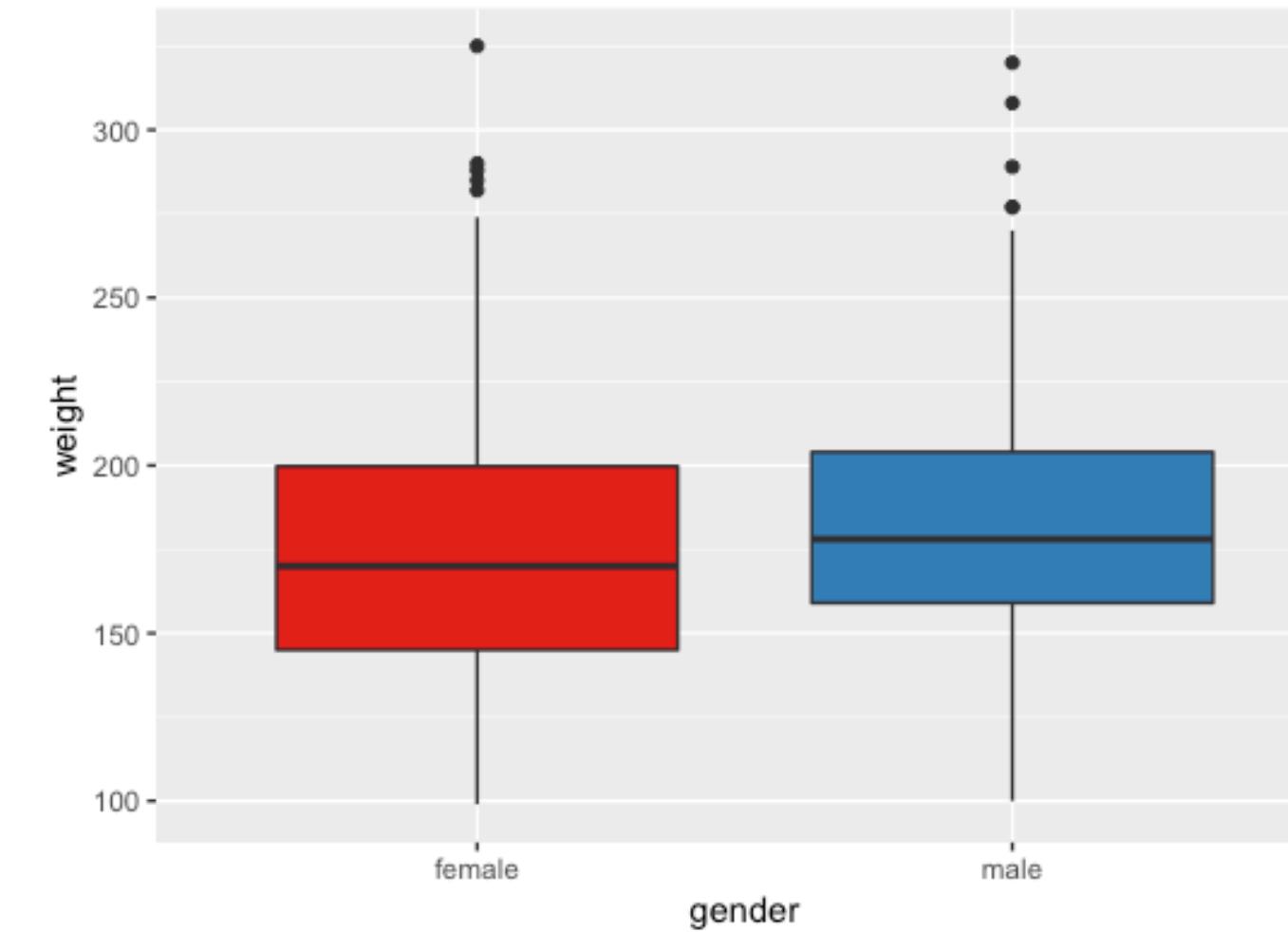
- Upper Whisker extend to the last point that is not larger than  $Q75 + 1.5 \times IQR$
- Lower Whisker extends to the last point that is not smaller than  $Q25 - 1.5 \times IQR$
- Whisker does not go beyond maximum or minimum value ! (Hence both whisker can have different length <  $1.5 \times IQR$ )

# Numerical values boxplots

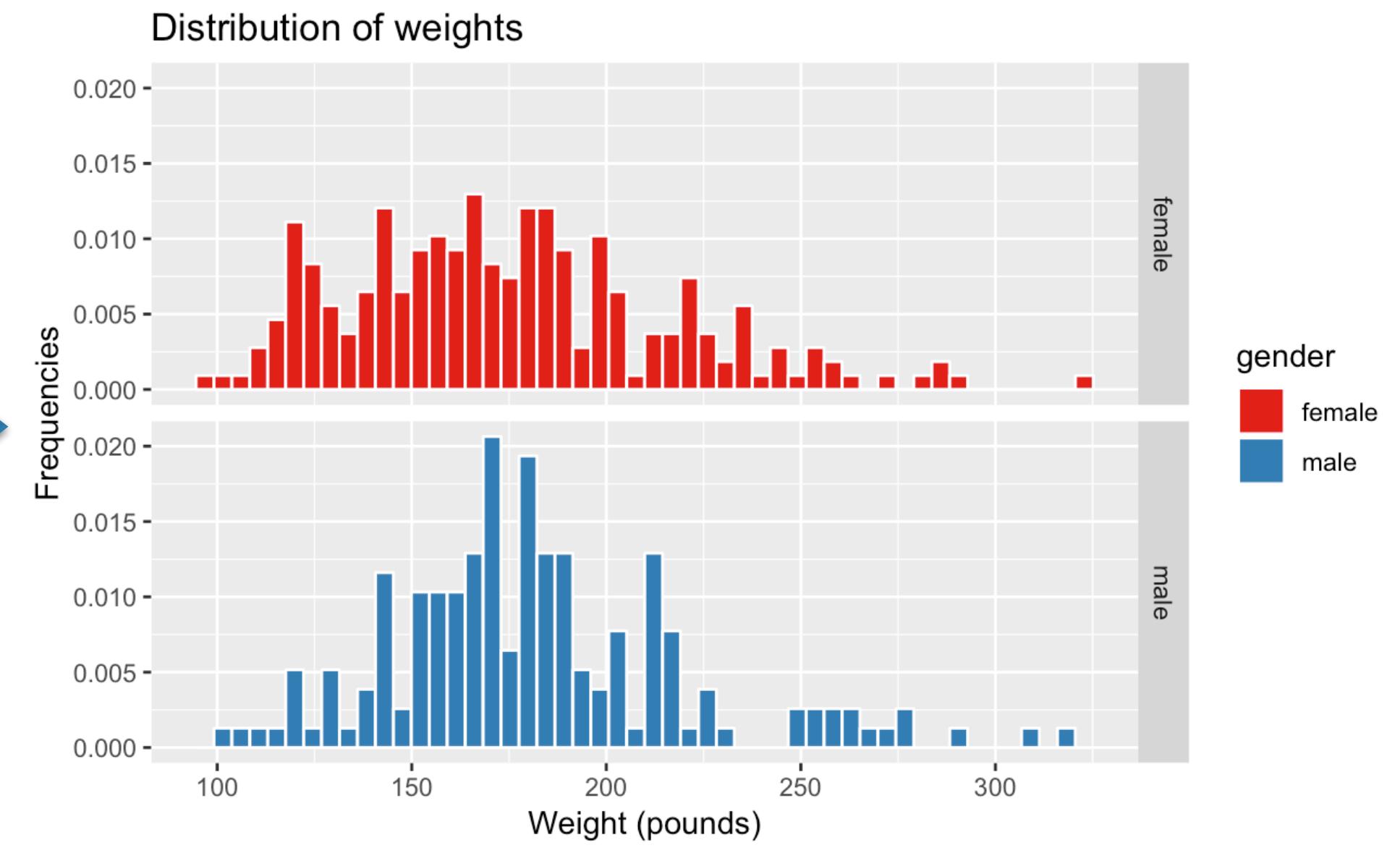


# Numerical values boxplots

Boxplot

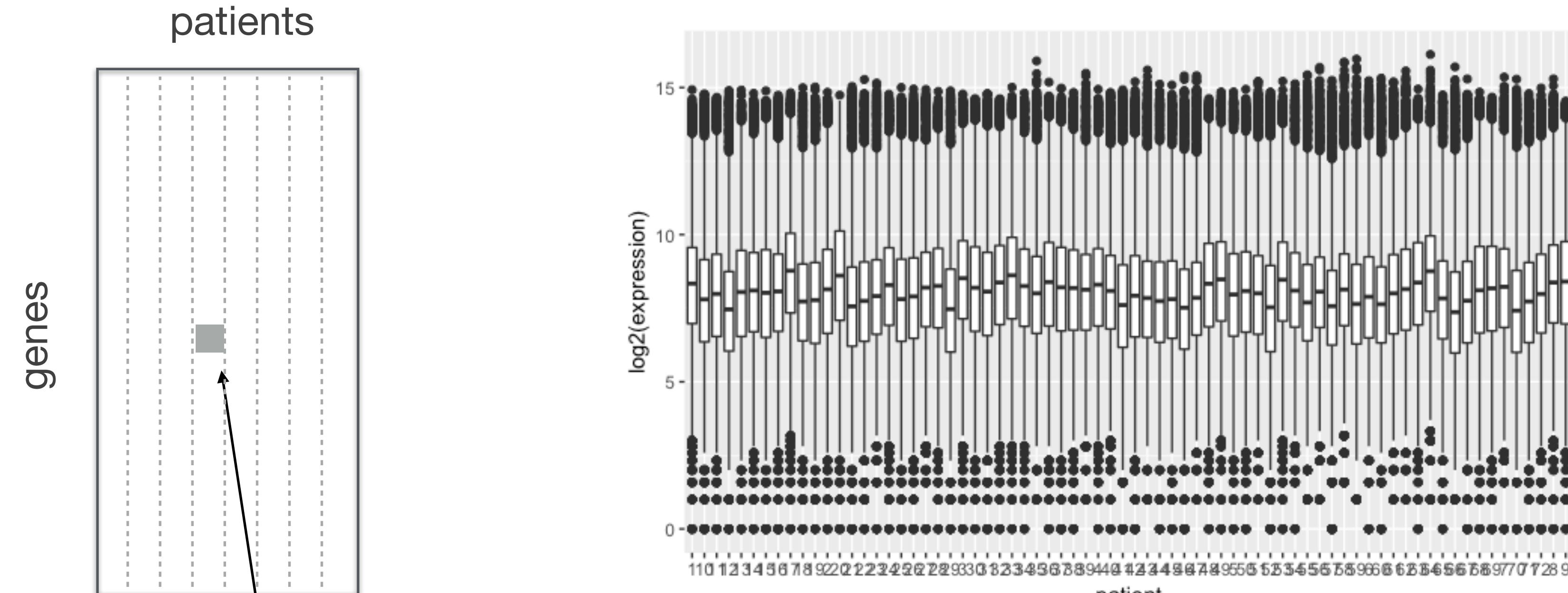


Histogram



Boxplots summarize the properties of the distribution  
Usefull to compare many distributions side-by-side

# Numerical values boxplots



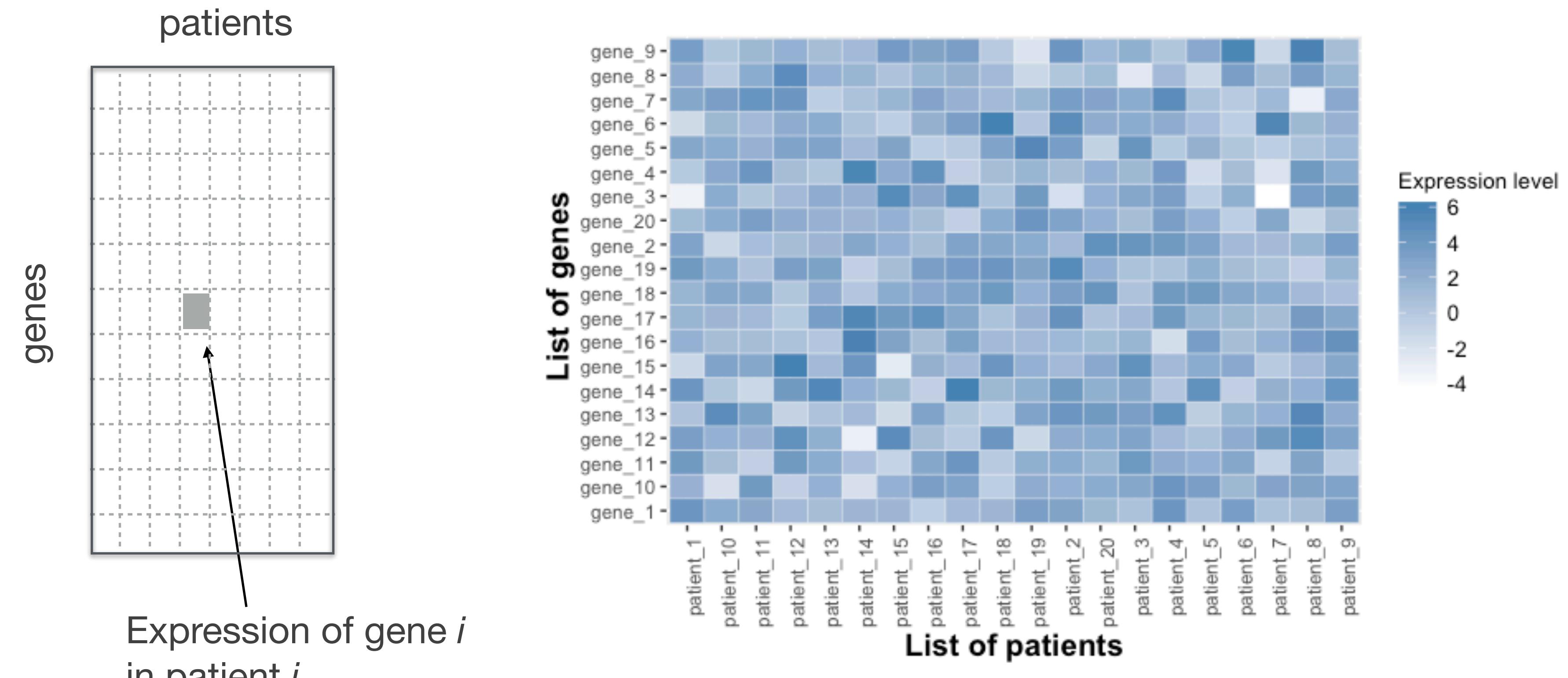
Expression of gene  $i$   
in patient  $j$

Question: do some patients have  
a different median gene expression?

→ *values for individual genes are lost in this type of plot!*

# Numerical Data : heatmaps

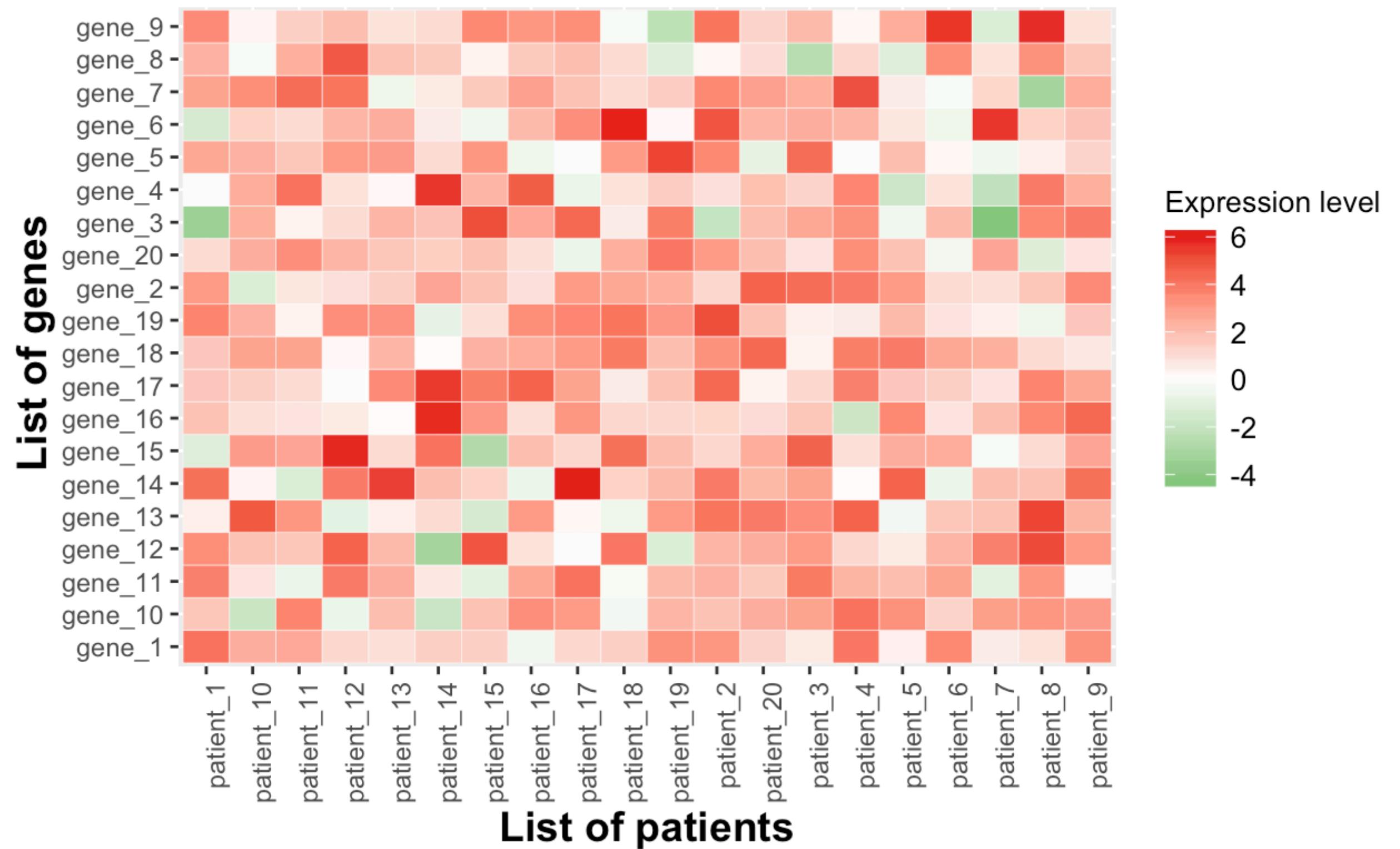
Heatmaps display numerical values in a data matrix using a color scheme



Question: do some genes in some patients have a different gene expression?

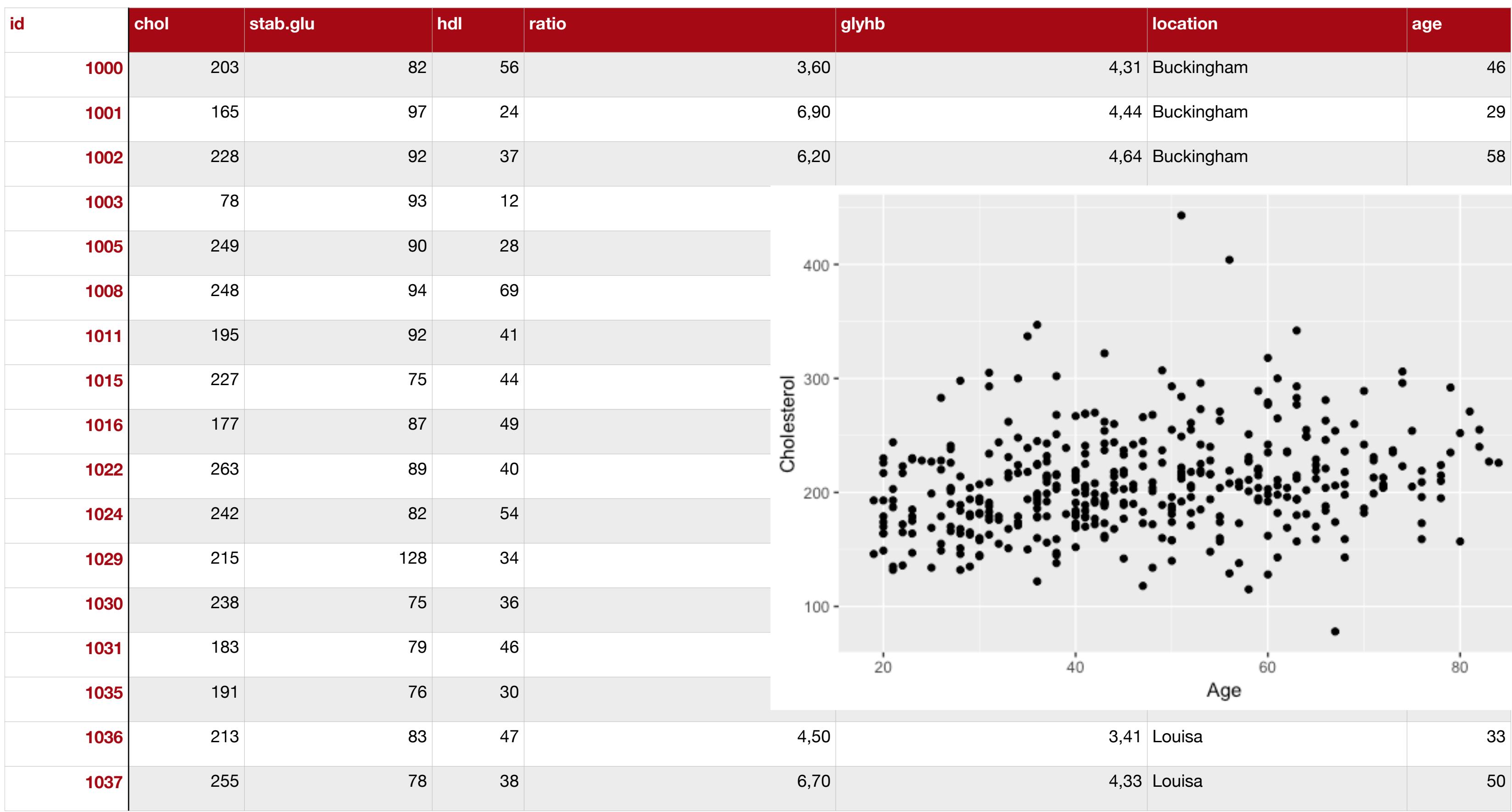
# Numerical Data : heatmaps

use symmetrical color scales for symmetrical ranges



# Numerical data comparing variable with scatter plots

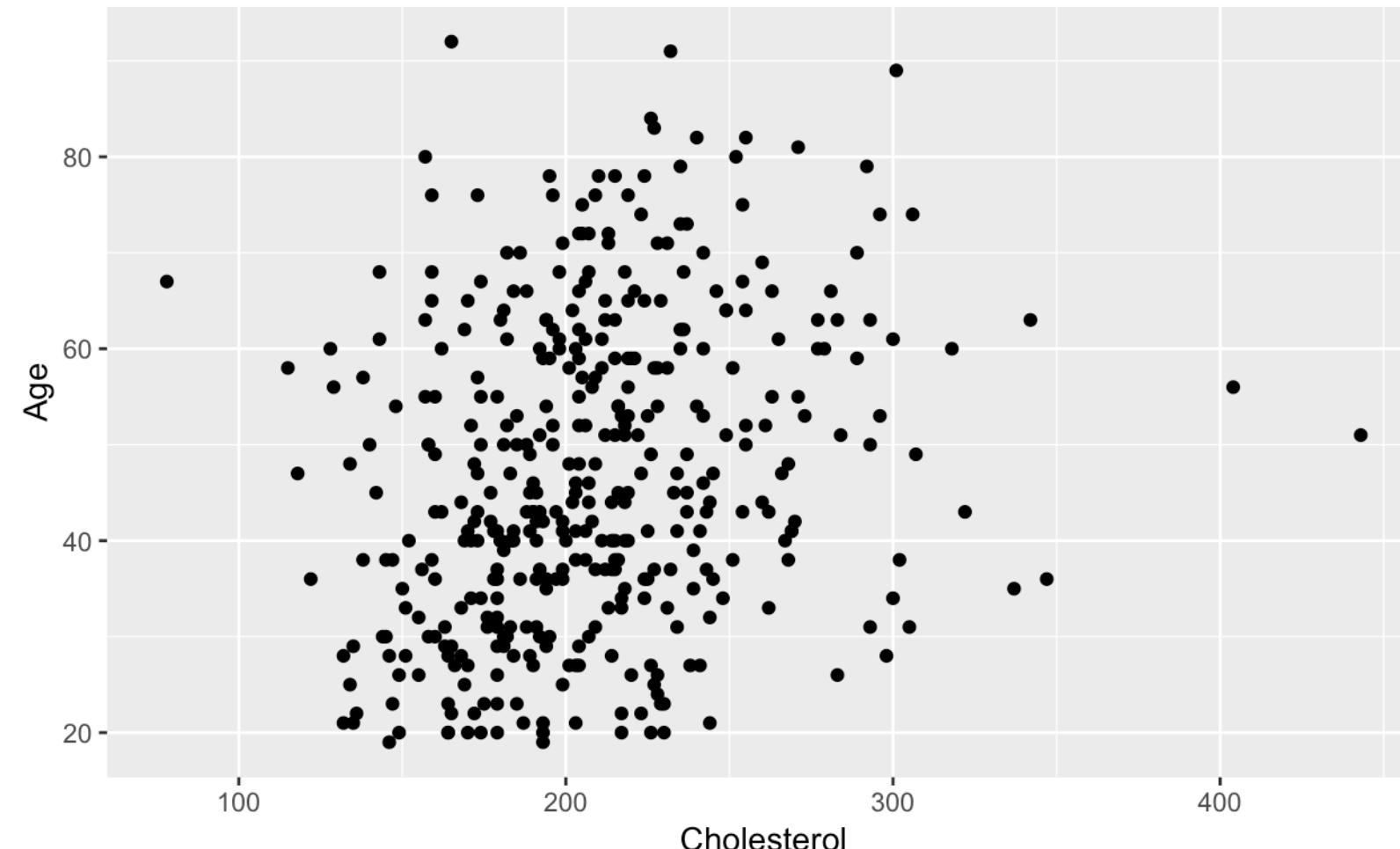
*any relation between age and cholesterol?*



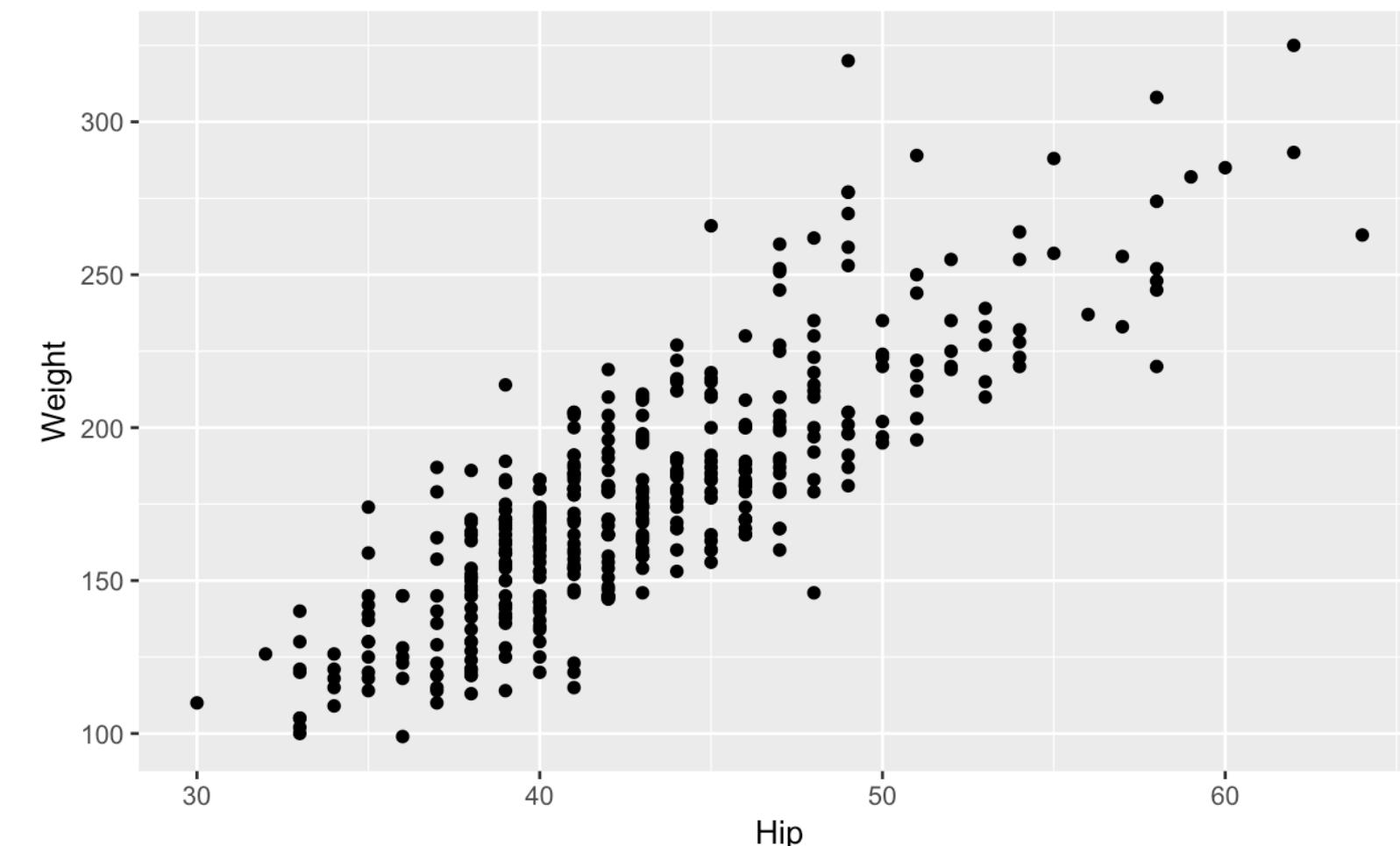
each dot is a patient

# Numerical data comparing variable with scatter plots

- we will later **quantify** this relationship in terms of **covariance / correlation** and determine how **significant** this relationship is!



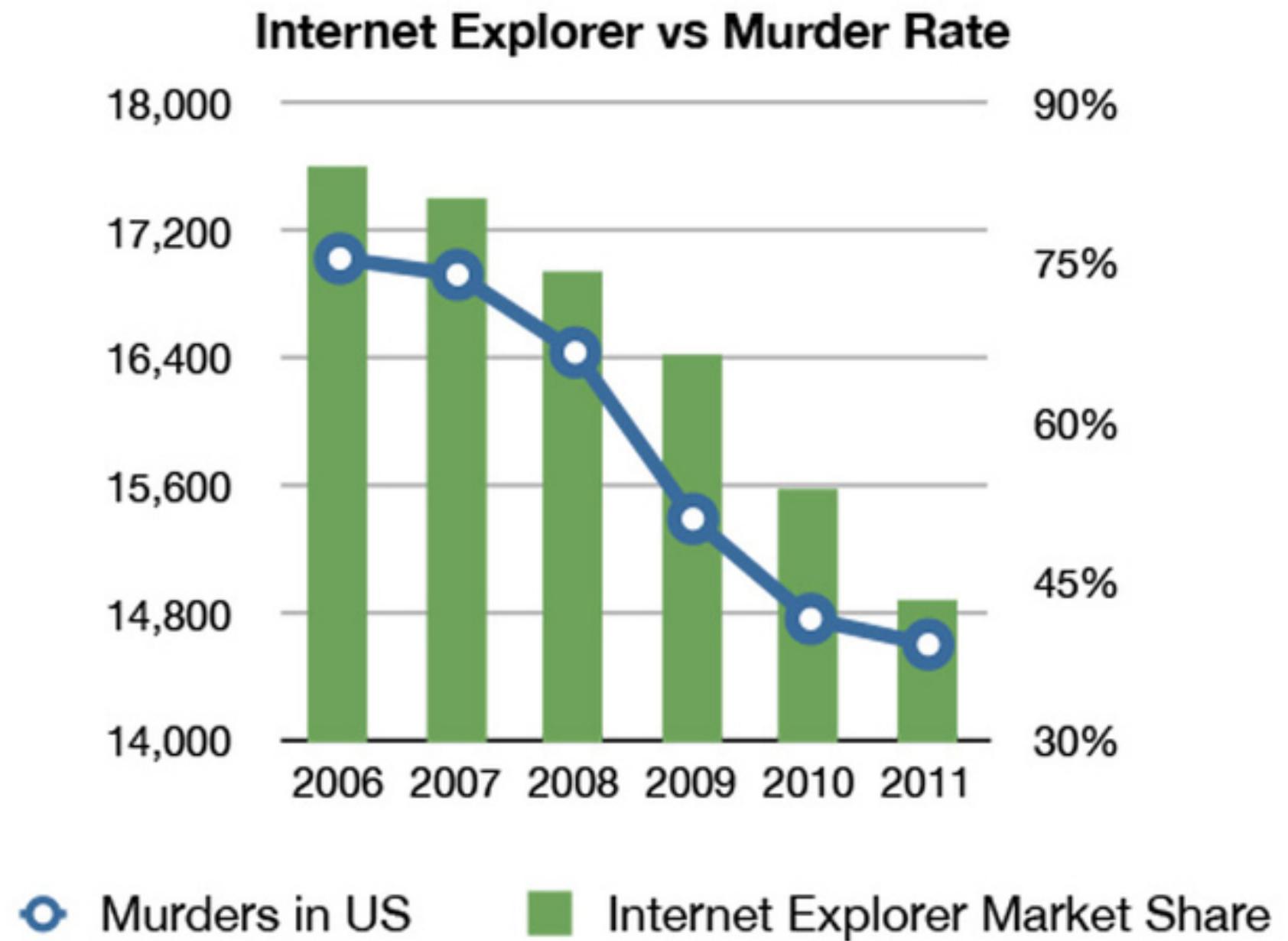
Weak relationship  
between Cholesterol and age



Strong relationship  
between Hip and Weight

# Numerical data comparing variable with scatter plots

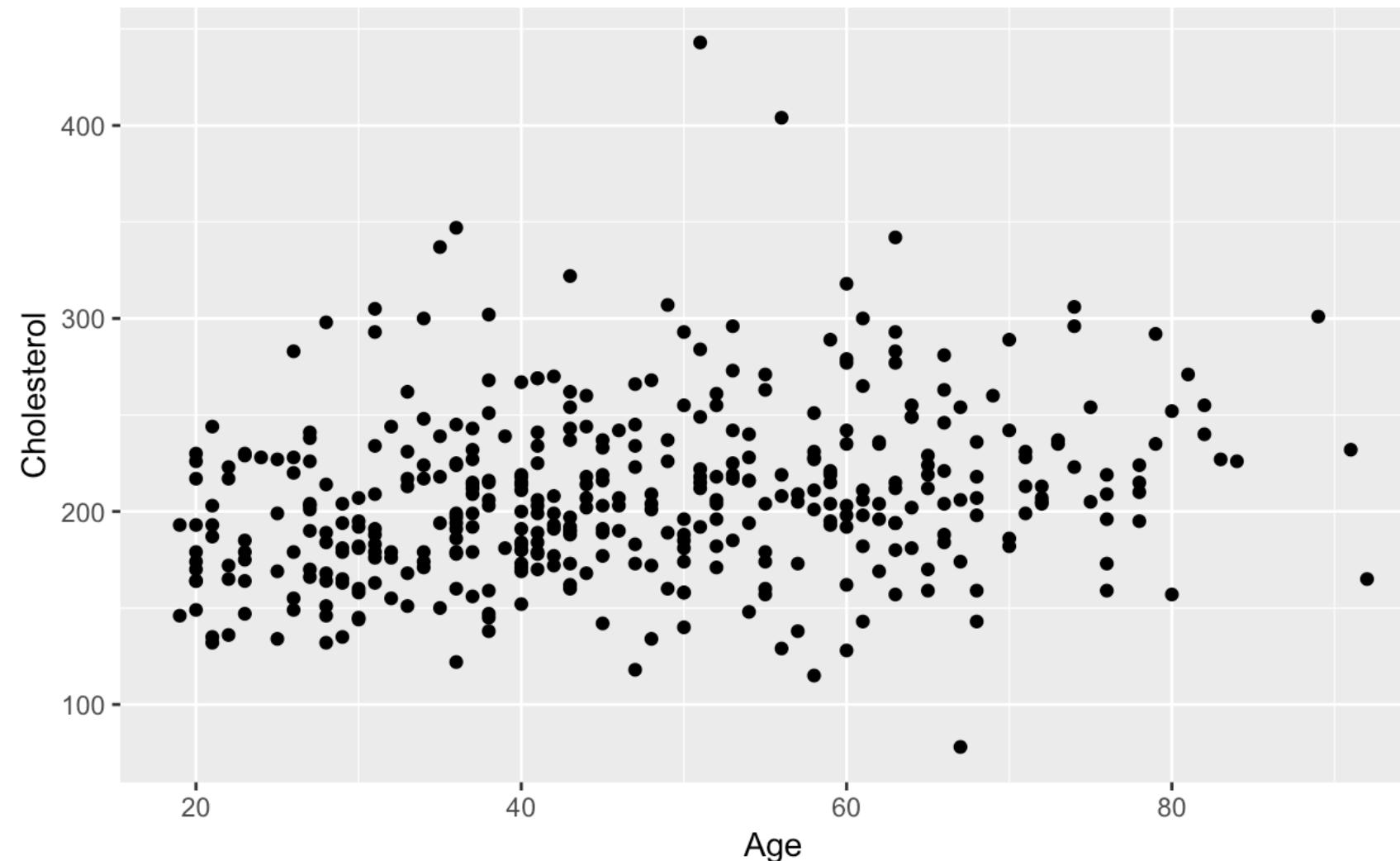
- Do not over-interpret scatter plots!
- Existence of relation between variables does not mean that there is a causal relationship between them!
- Correlation is NOT causality!!



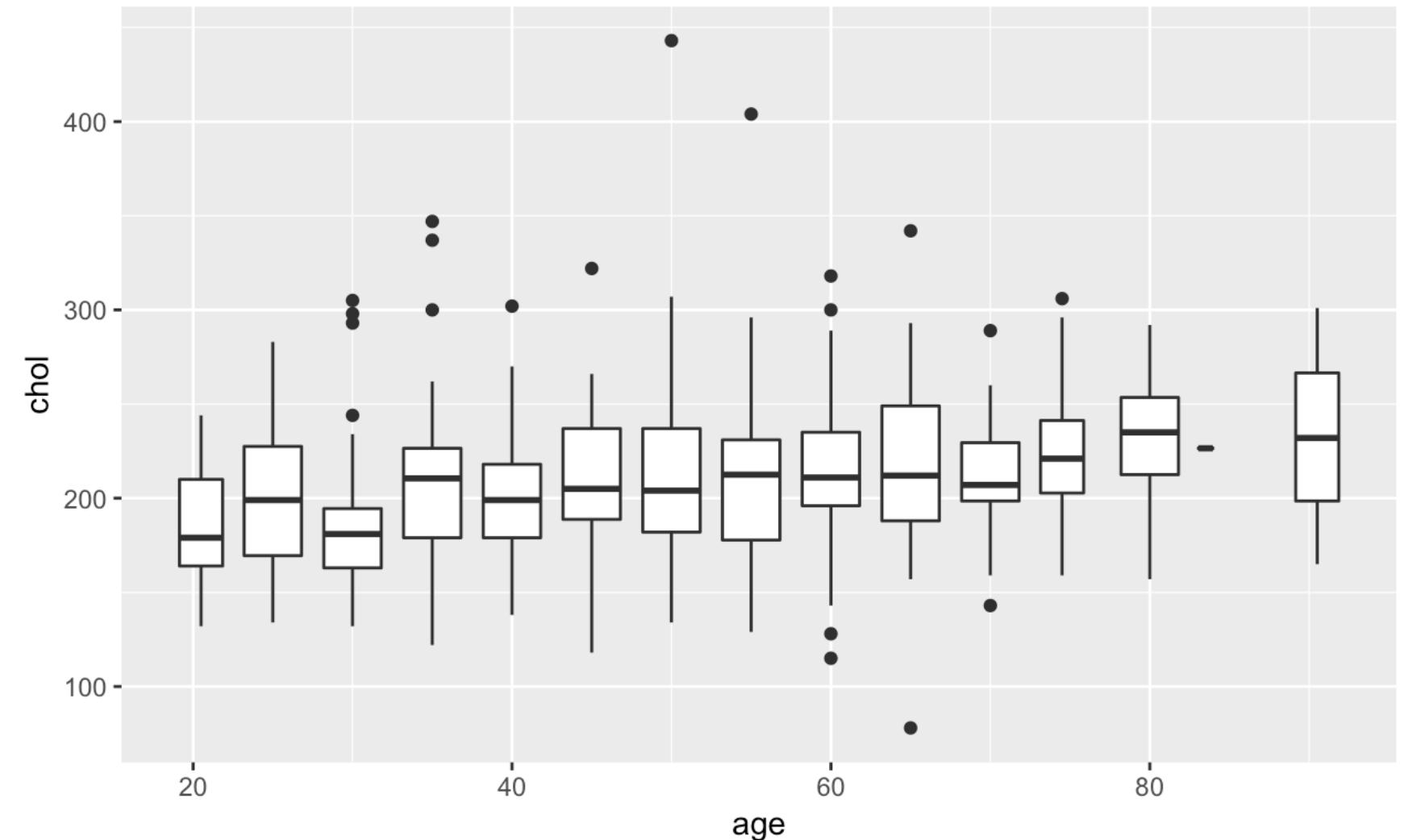
<http://www.tylervigen.com/spurious-correlations>

# Numerical data comparing variable with scatter plots

- A **continuous numerical** variable can always be transformed into an **ordinal categorical** variable through binning



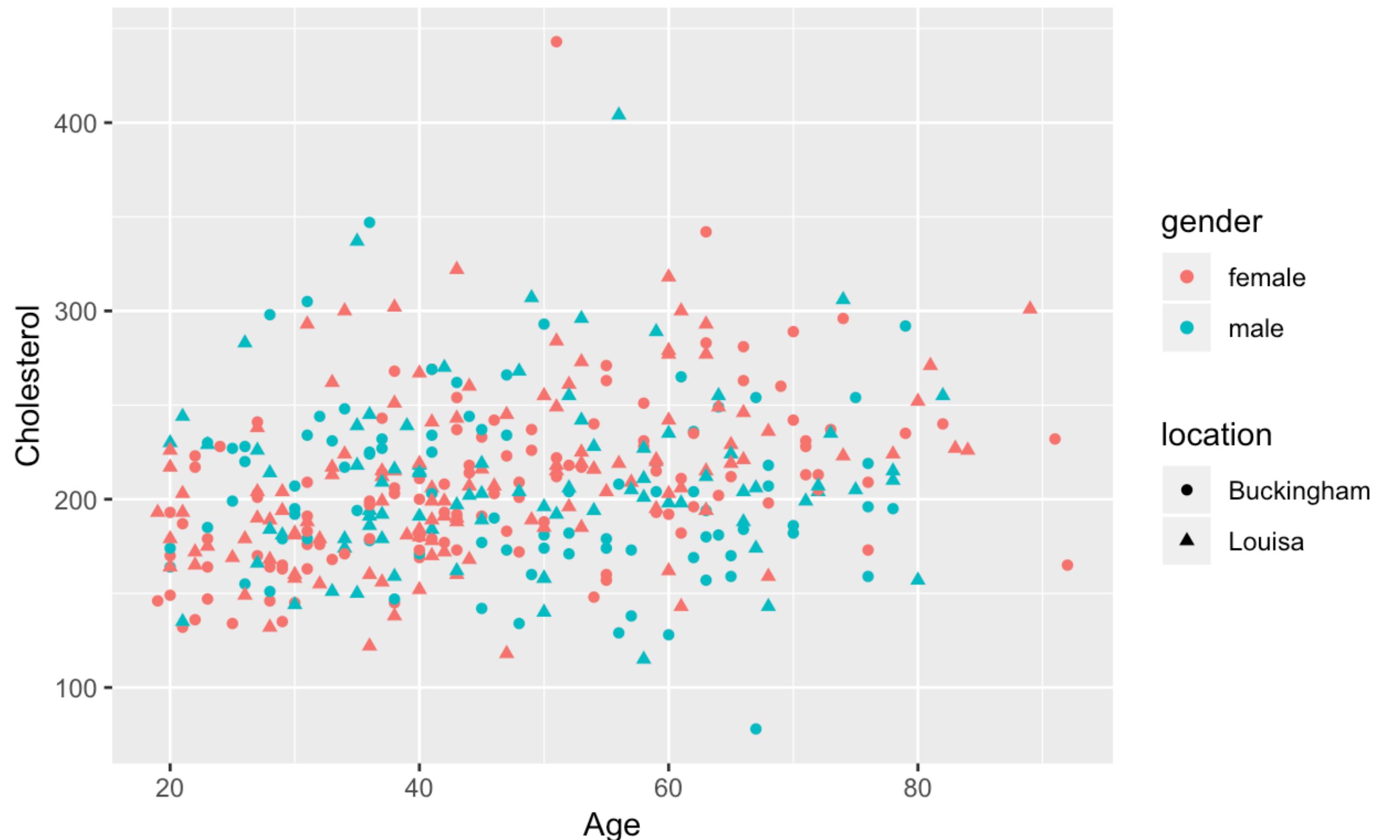
Age = continuous variable



Age = ordinal variable (bins)

# Numerical data comparing variable with scatter plots

- Additional categorical / numerical variables can be added using color, shape, size of dot ,...



# Summary on visualization

## Single variable plot

plot counts

type of plot	
continuous variable	histogram
categorical variable	barplot

## Two variable plot

plot relationship

continuous variable	continuous variable	categorial data
categorical variable	scatter plot	boxplot
categorical variable		heatmap