

Statistical modelling

#1.c Exploratory Data Analysis

Dr. Léo Belzile
HEC Montréal

Type of data

Data base typically comprise many different *variable types*.

Distinguishing between the later is needed for

- + correct model choice,
- + proper graphical representation,
- + interpretation of effects.

Numerical variables

CONTINUOUS

measured data, can have ∞ values within possible range.



I AM 3.1" TALL
I WEIGH 34.16 grams

DISCRETE

OBSERVATIONS can only exist at LIMITED VALUES, OFTEN COUNTS.

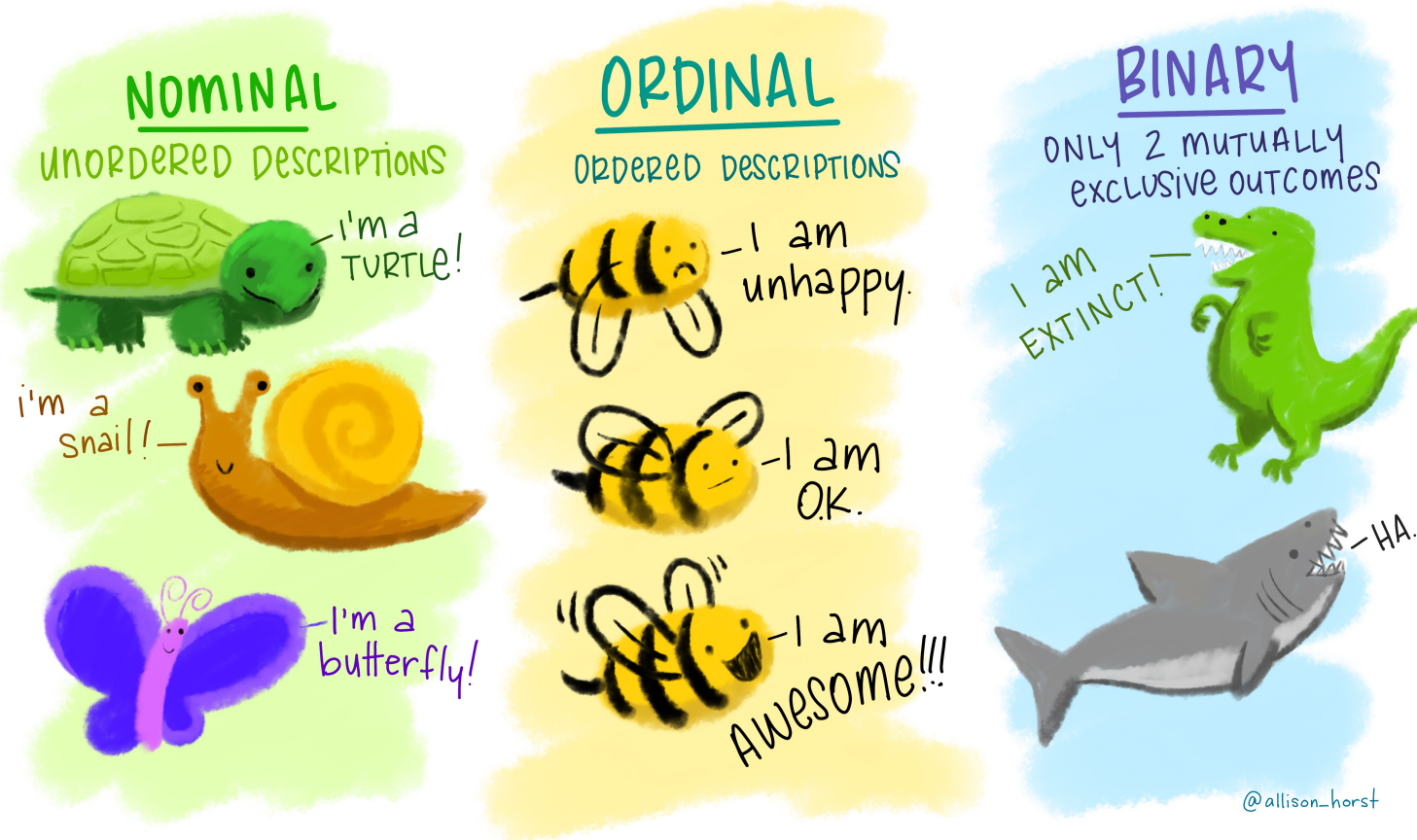


I HAVE 8 LEGS
and
4 SPOTS!

@allison_horst

Drawing by Allison Horst of continuous (left) and discrete (right) numerical variables.

Categorical variables



Drawing by Allison Horst of nominal (left), ordinal (middle) and binary (right) categorical variables.

Graphics and data

The simple graph has brought more information to the data analyst's mind than any other device.

— John Tukey

What is a good graph?

Graphical excellence consists of complex ideas communicated with clarity, precision, and efficiency ... Graphical excellence is that which gives to the viewer the greatest number of ideas in the shortest time with the least ink in the smallest space.

— Tufte, 1983

Grammar of graphics

Wilkinson, L. (2005), *The Grammar of Graphics*(2nd ed.) Statistics and Computing, New York: Springer.

- + Elements (layers):
 - + data
 - + mapping (variables -> aesthetics)
 - + geometric objects
 - + transformations
 - + positioning
- + Scale / guide
- + Coordinates (facets, coordinate system)

Some golden rules for effective visualization

Rule 1: the choice of graphic depends on the variable type

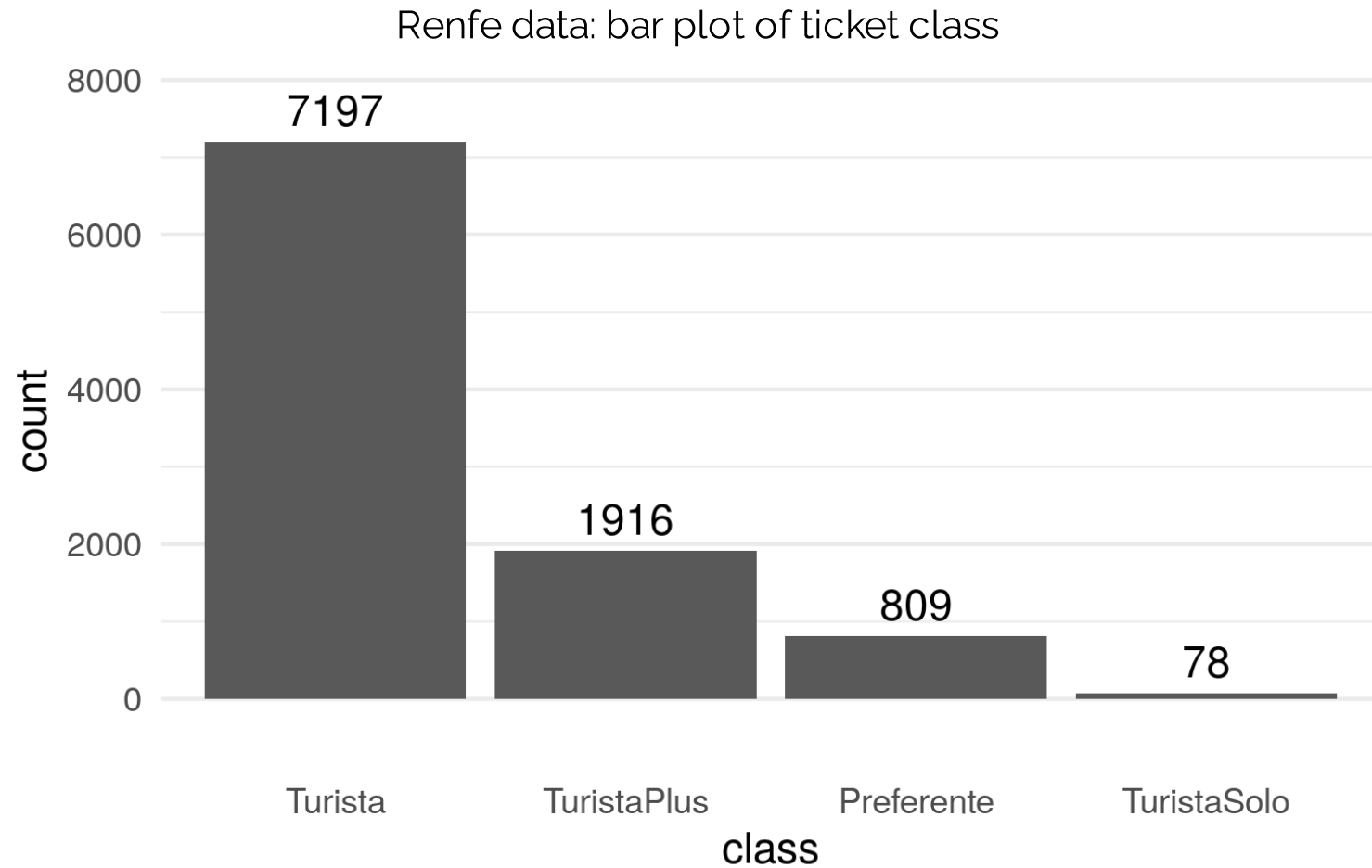
One variable

- + continuous: histogram/density plot
- + discrete: bar plot
- + categorical: bar plot (frequency or percentage)

Two variables

- + continuous: scatterplot
- + categorical: bar plots (one group via color), heatmap
- + continuous \times categorical: box-and-whisker plot, violin plot

- R graph + R code + SAS graph + SAS code



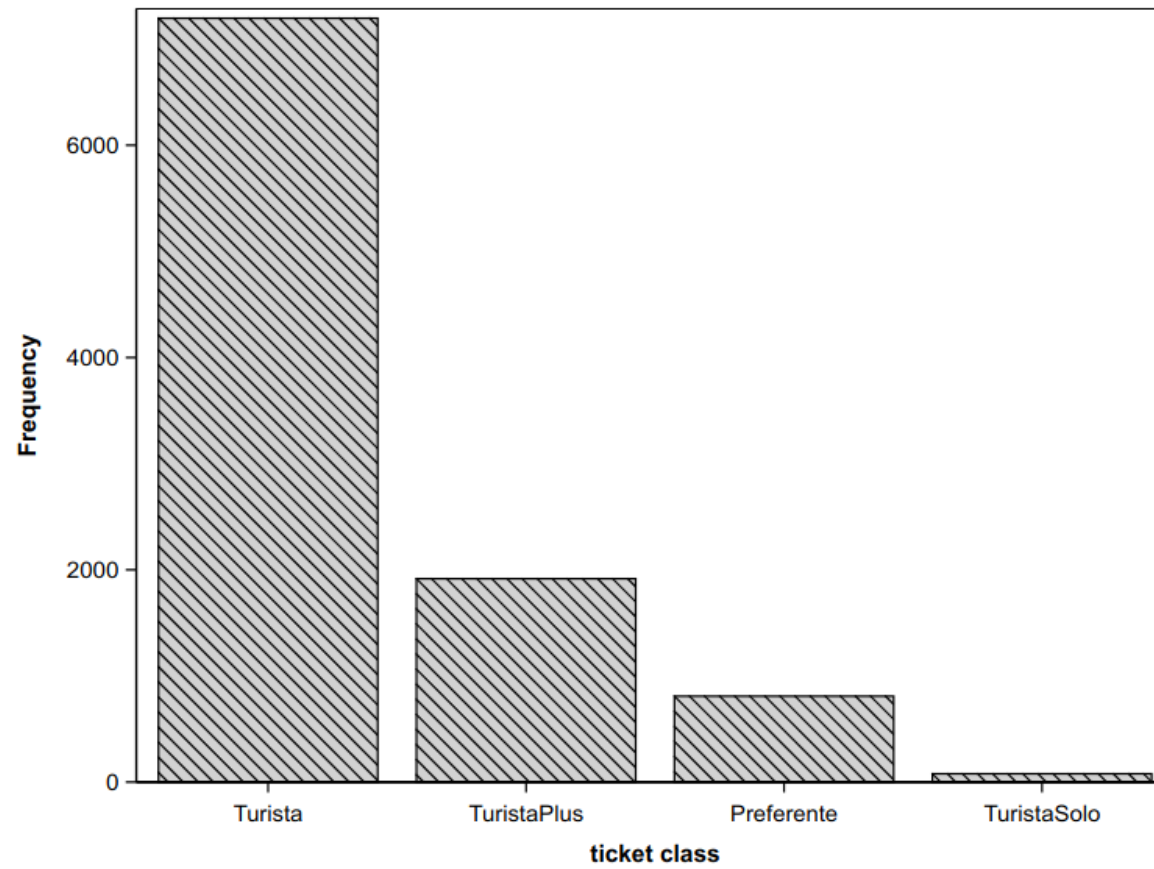
- R graph + R code + SAS graph + SAS code

One categorical variable: bar plot

```
ggplot(data = renfe,  
       aes(x = forcats::fct_infreq(classe))) +  
  geom_bar() +  
  geom_text(stat='count', aes(label=..count..), vjust=-0.5) +  
  labs(x = "class",  
       y = "count") +  
  scale_y_continuous(expand = c(.125, 0)) +  
  theme(panel.grid.major.x = element_blank())
```

- + Order values by frequency.
- + If the labels are too long, rotate the axis (+ `coord_flip()`).

- R graph + R code + SAS graph + SAS code

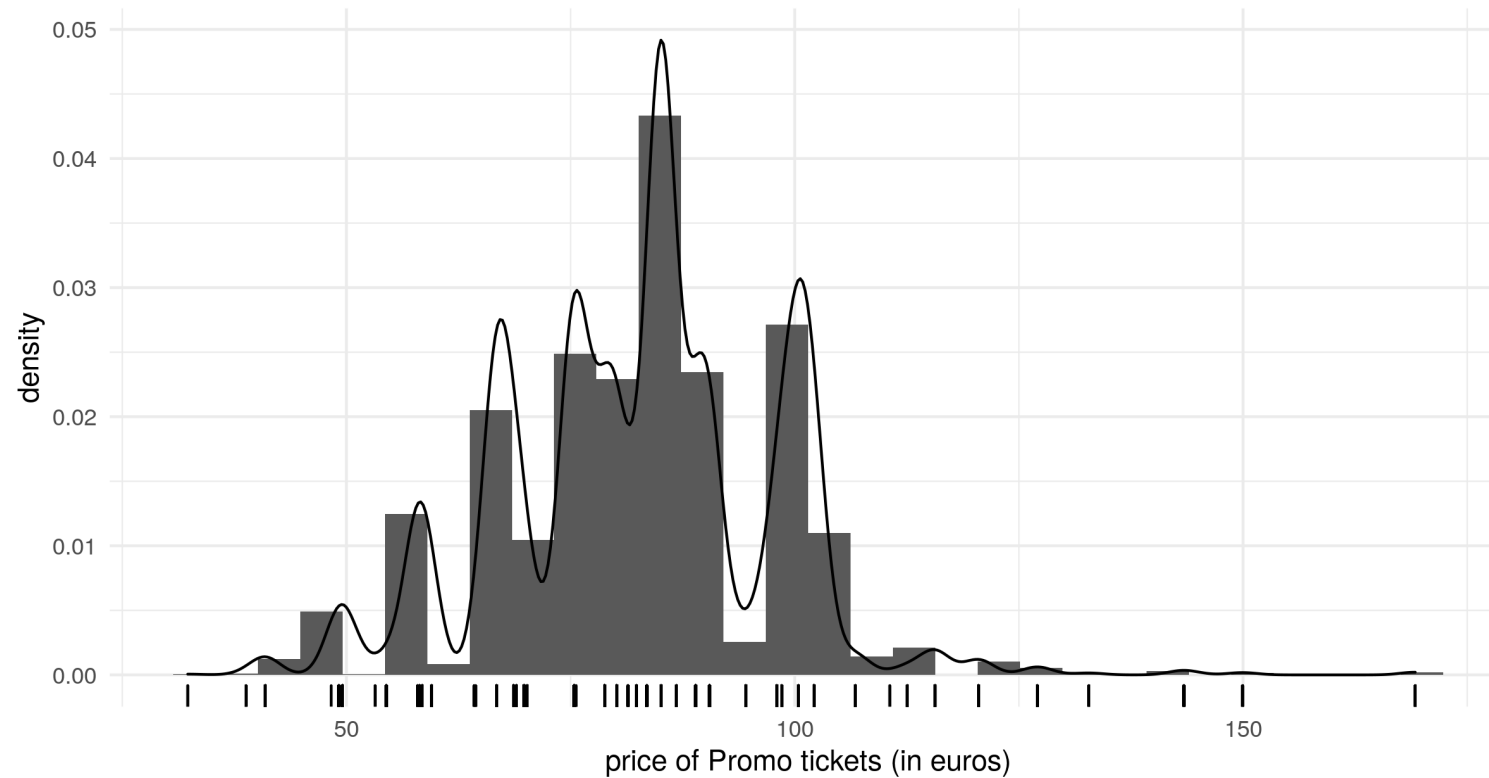


-
- R graph + R code + SAS graph + SAS code
-

```
proc sgplot data=statmod.renfe;  
vbar class / categoryorder=respdesc;  
xaxis label="ticket class";  
run;
```

- R graph + R code + SAS graph + SAS code

Renfe data: histogram of Promo ticket price

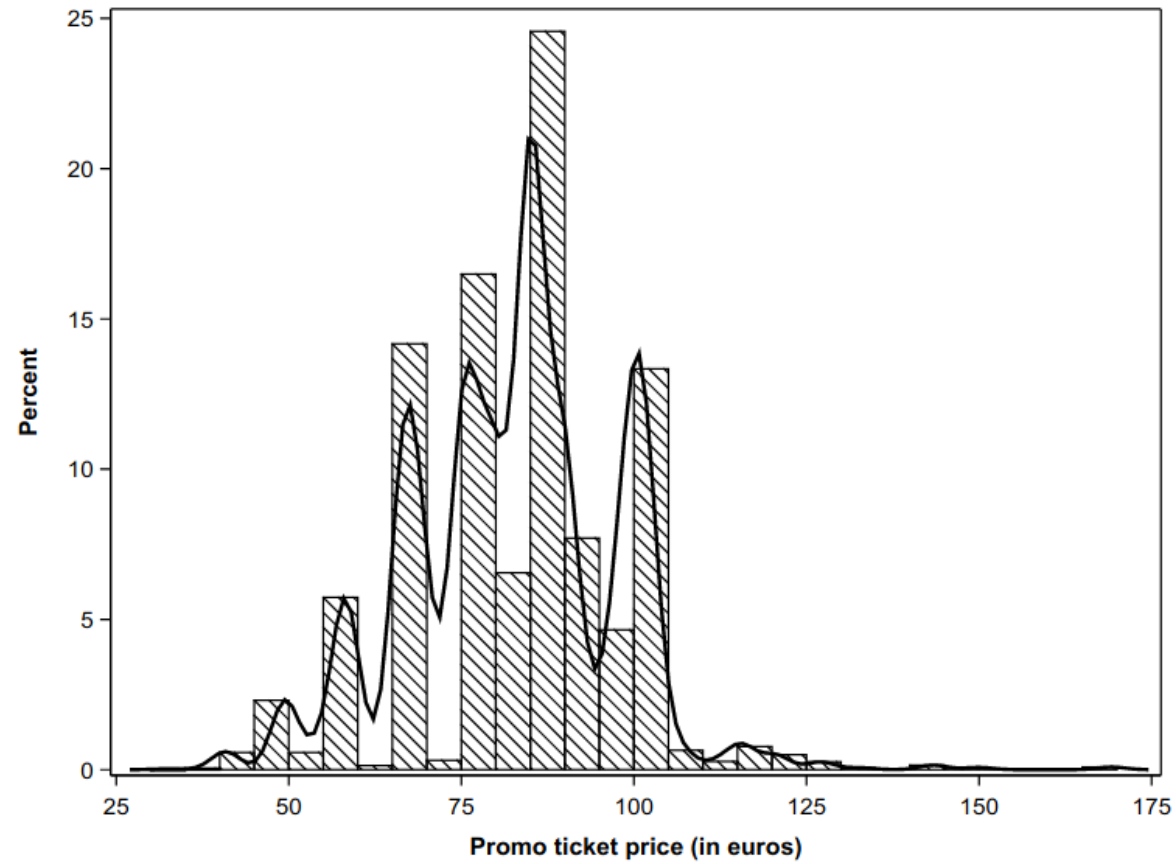


-
- R graph + R code + SAS graph + SAS code
-

One continuous variable: histogram/density

```
renfe %>% subset(fare == "Promo") %>%  
  ggplot(aes(x = price)) +  
    geom_histogram(aes(y = ..density..), bins = 30) +  
    geom_density() +  
    geom_rug(sides = "b") +  
    labs(x = "price of Promo tickets (in euros)",  
         y = "density")
```

- R graph + R code + SAS graph + SAS code



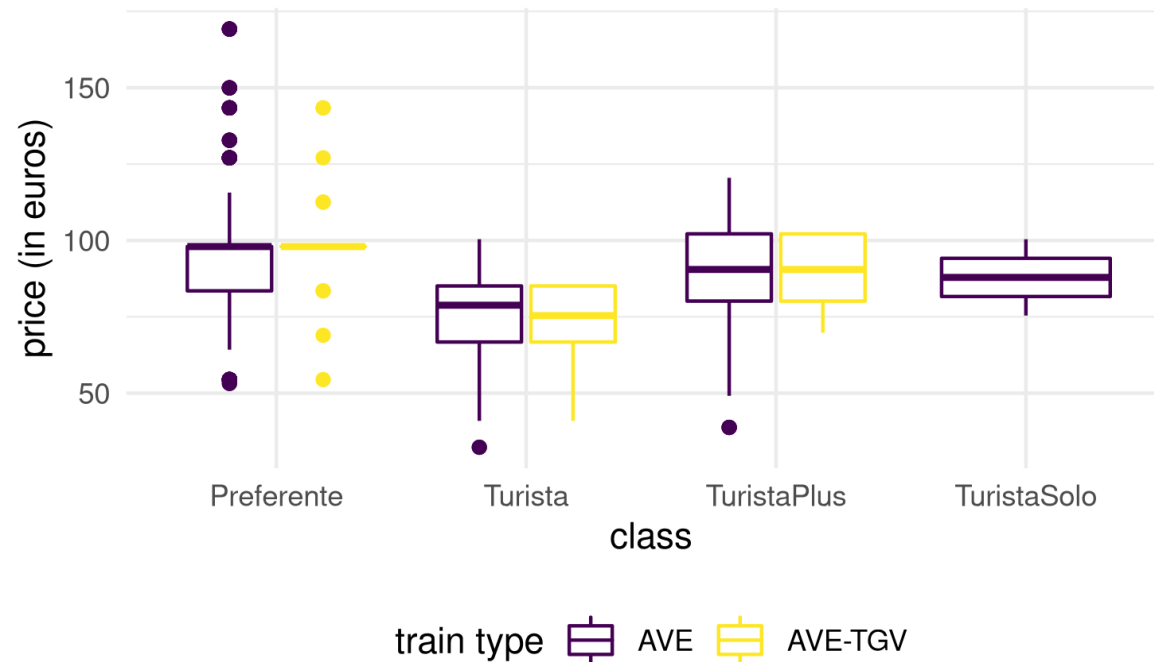
- R graph + R code + SAS graph + SAS code

```
data renfe_promo;
set statmod.renfe;
where fare="Promo";
run;

proc sgplot data=renfe_promo noautolegend;
histogram price;
density price / type=kernel;
xaxis label = "Promo ticket price (in euros)";
run;
```

- R graph + R code + SAS graph + SAS code

Renfe data: box-and-whiskers plot of Promo tickets price as a function of class



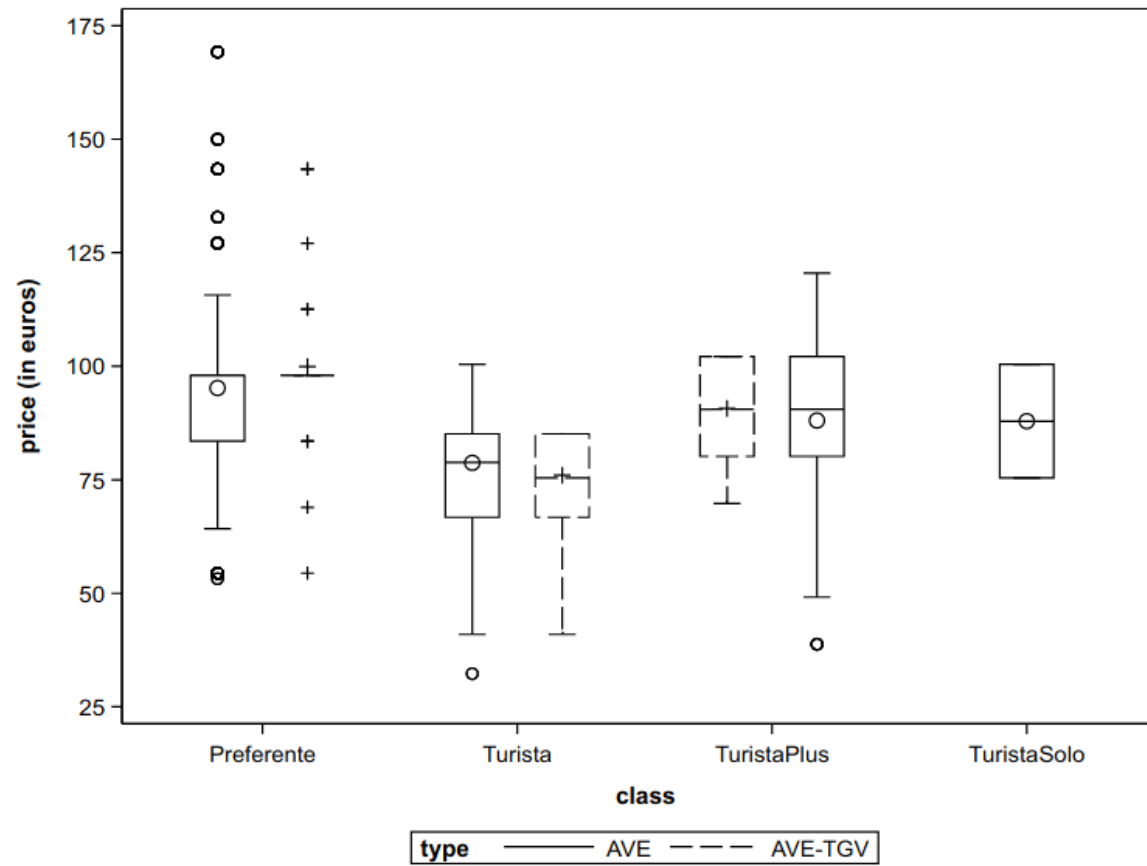
-
- R graph + R code + SAS graph + SAS code
-

Two variables (continuous and categorical): boxplot

```
renfe %>% subset(fare == "Promo") %>%  
  ggplot(aes(y = price, x = class, col = type)) +  
  geom_boxplot() +  
  labs(y = "price (in euros)", col = "train type") +  
  theme(legend.position = "bottom") +  
  scale_colour_viridis_d()
```

- + We added a categorical variable (`type`) through use of colour.
- + Use an appropriate color palette (for color-blind people and for black and white printing).

- R graph + R code + SAS graph + SAS code

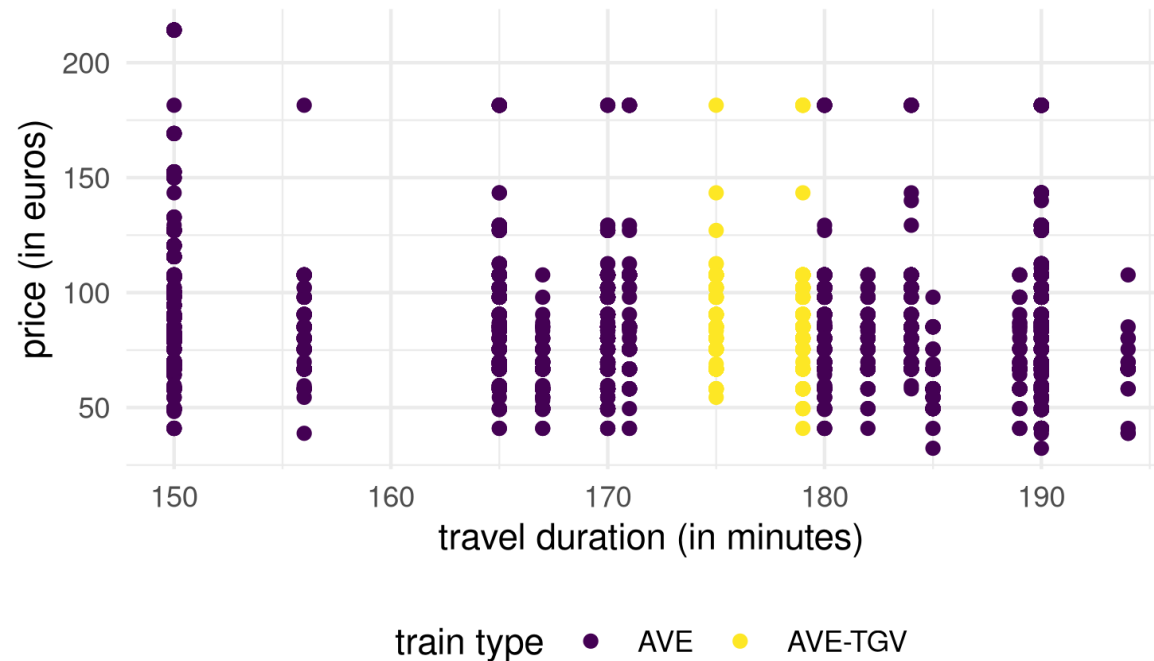


- R graph + R code + SAS graph + SAS code

```
proc sgplot data=renfe_promo;  
vbox price / category=class group=type;  
yaxis label = "price (in euros)";  
run;
```

- R graph + R code + SAS graph + SAS code

Renfe data: scatterplot of ticket price as a function of travel time for high speed trains



- R graph + R code + SAS graph + SAS code

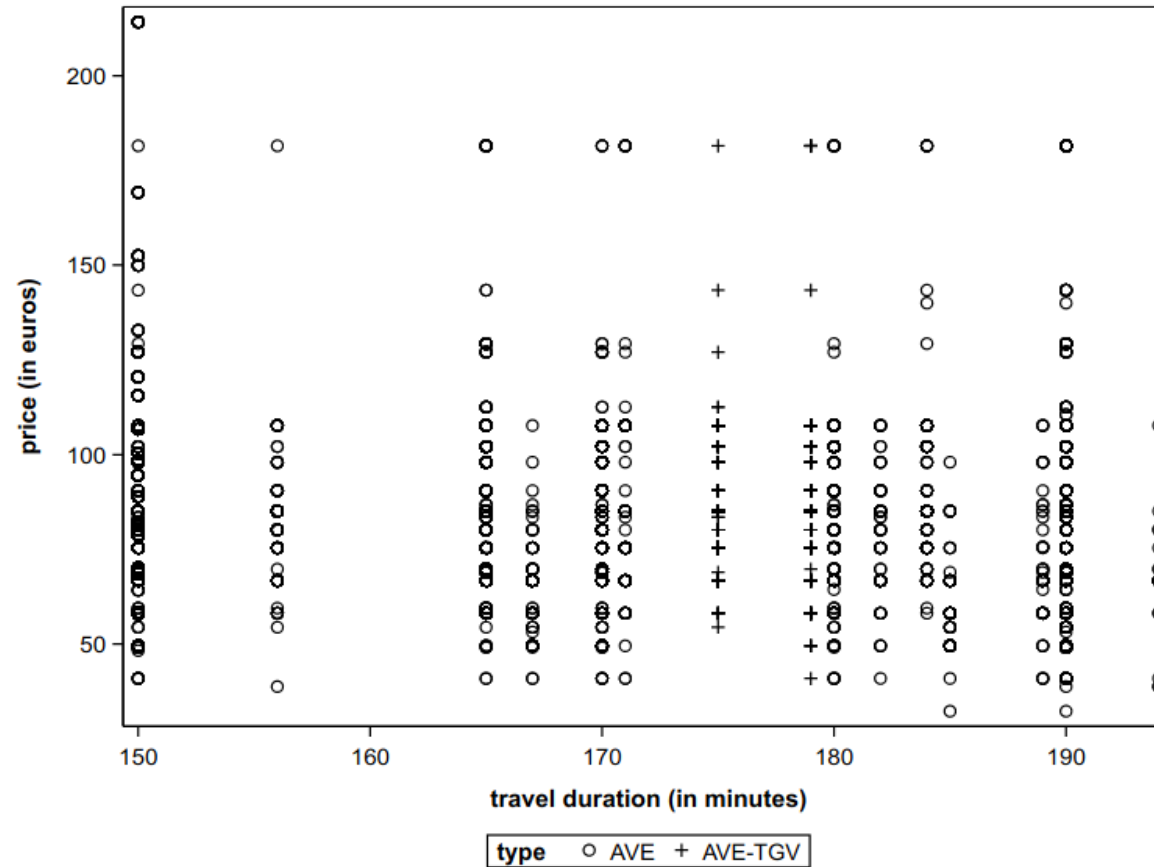
Two continuous variables and one categorical variable: scatterplot

```
renfe %>% subset(type != "REXPRESS") %>%  
  ggplot(aes(x = duration, y = price, col = type)) +  
  geom_point() +  
  labs(y = "price (in euros)",  
       x = "travel duration (in minutes)",  
       col = "train type") +  
  theme(legend.position = "bottom") +  
  scale_colour_viridis_d()
```

What is wrong with the previous display?

How could we fix the graph?

- R graph + R code + SAS graph + SAS code



-
- R graph + R code + SAS graph + SAS code
-

```
data renfe_ave;
set statmod.renfe;
where type NE "REXPRESS";
run;

proc sgplot data=renfe_ave;
scatter y=price x=duration / group=type;
axis label="travel duration (in minutes)";
yaxis label="price (in euros)";
run;
```

Rule 2: A graph tells a story by itself

Your graphic must be stand-alone with the legend.

- + some visualization choices are more effective/adequate than others
- + include both variable name **and** units if ambiguous
- + add a description in the text and cross-reference
- + pay attention to scale (adequate font size for legibility)

Rule 3: pay attention to human visual perception

Avoid junk chart

- + ratio length/width
- + spacing between bands
- + axis limits (with or without zero)
- + choice of color (grayscale, color-blind friendly palettes)
- + comparing areas is difficult
- + avoid 3D graphs / rotation

Graphical exploratory data analysis

Numerical quantities focus on expected values, graphical summaries on unexpected values.

— John Tukey

- + Ask questions related to the data
- + Look for answers using graphs
- + Infirm or confirm your intuitions
- + Refine questions based on preliminary findings
- + Rinse and repeat
- + Write a summary of key findings

References

- + *Fundamentals of Data Visualization* par Claus O. Wilke
- + Chapter 3 of *R for Data Science* by Garrett Grolemund and Hadley Wickham
- + Chapter 1 of *Data Visualization: A practical introduction* by Kieran Healy