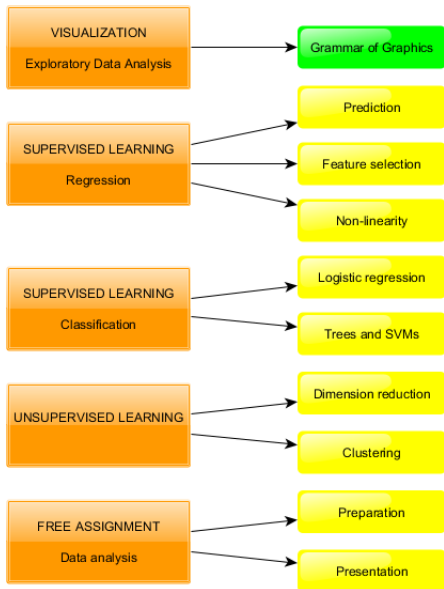


Visualization

Grammar of Graphics

Maarten Cruyff

Program



1. Data visualization
2. Grammar of Graphics
3. Grammar of Data Manipulation
4. Tufte's Principles of Graphical Excellence

What's data visualization?

Communication of data by encoding it as visual objects, i.e.

- dots, lines, bars, etc.

to make data more accessible, understandable and usable.

. . . an integral part of data science process

Data Science Process

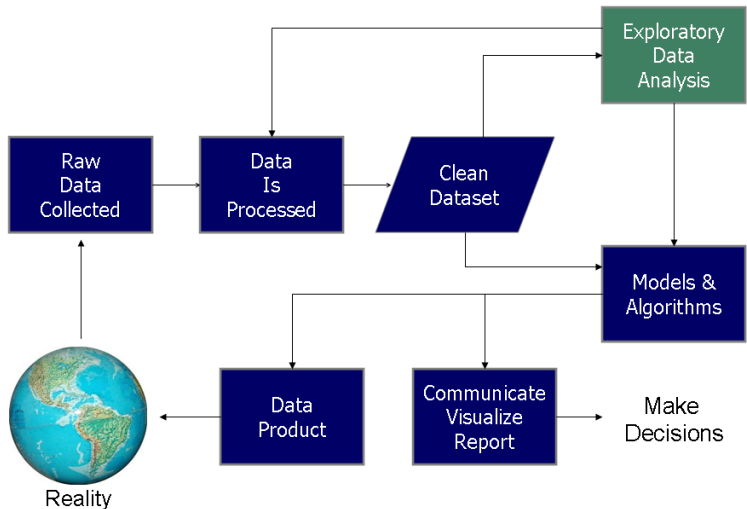
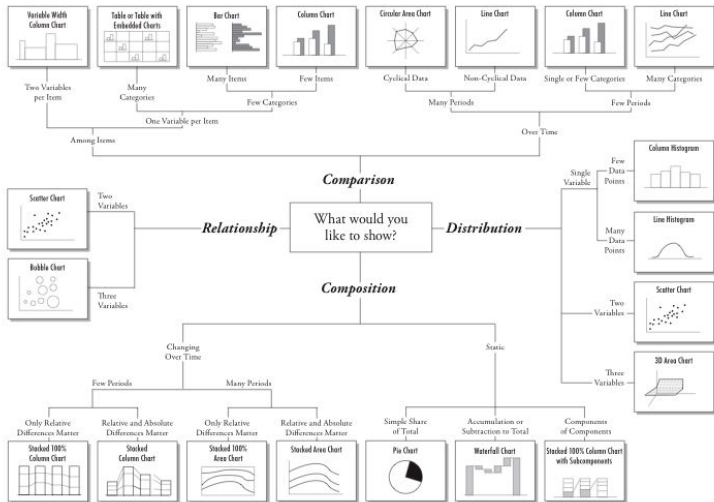


Chart Suggestions—A Thought-Starter



Why data visualization?

Human brain excels in distinguishing differences in

- length
- shape
- orientation
- hue

Brain processes graphical information better than text:

- Retina: 10^6 or 10^7 bits/sec
- Reading: 10^2 to 10^3 bits/sec

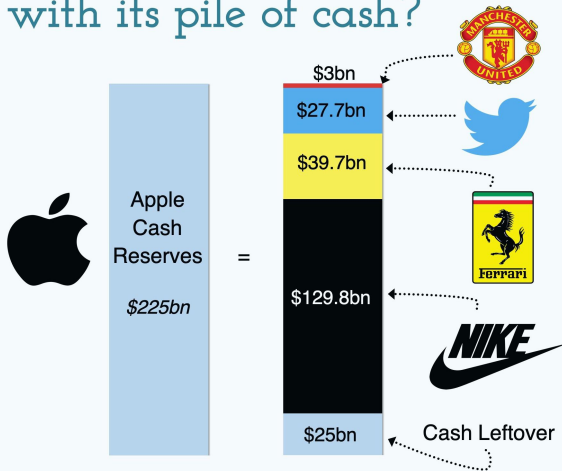
Early example

John Snow finds source of 1854 Soho cholera outbreak (not miasma)



Recent example

What could Apple buy with its pile of cash?



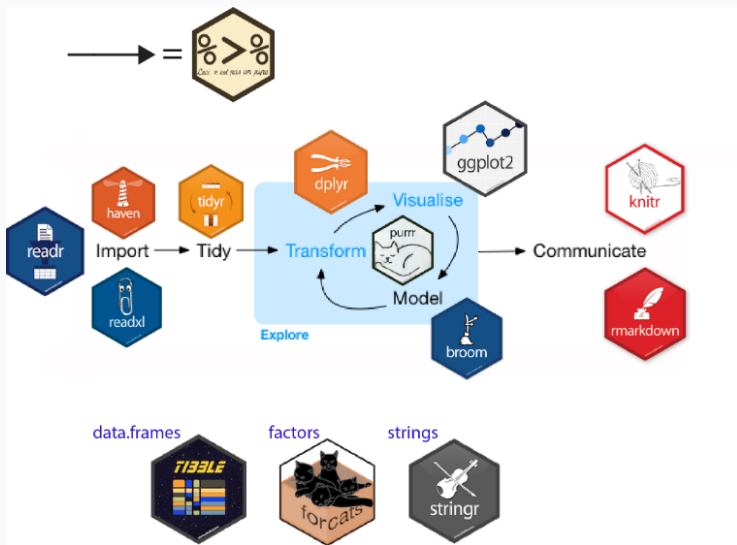
Instagram: @chartdaily

Source: Apple 10-Q, Yahoo Finance
(Market Values as of June 26th, 2019)

Grammar of Graphics

Data science in R

Packages and pipes



Packages for data visualization

ggplot2

- package for making plots
- based on the *grammar of graphics*

dplyr

- package for data preparation
- based on the *grammar of data manipulation*

Both packages work beautifully together

- with pipe operator `%>%`

Grammar of Graphics

Building plots layer-by-layer

Describes all the non-data ink

Plotting space for the data

Statistical models & summaries

Rows and columns of sub-plots

Shapes used to represent the data

Scales onto which data is mapped

The actual variables to be plotted

Theme

Coordinates

Statistics

Facets

Geometries

Aesthetics

Data

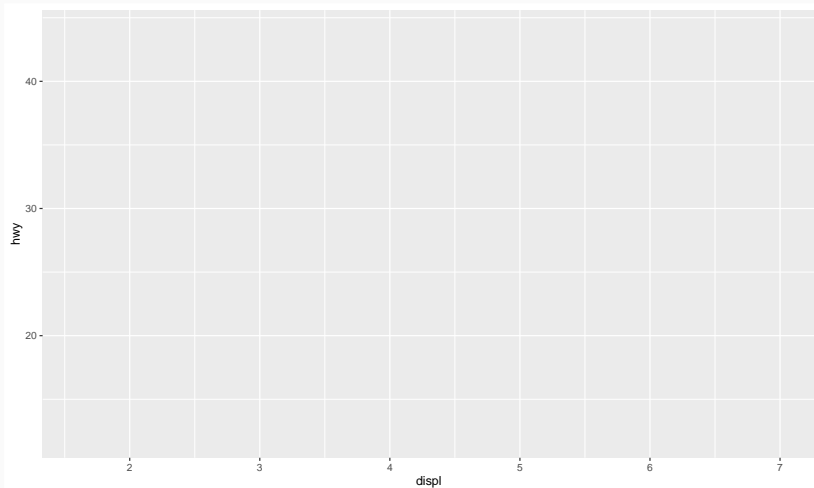


Example of ggplot2 code

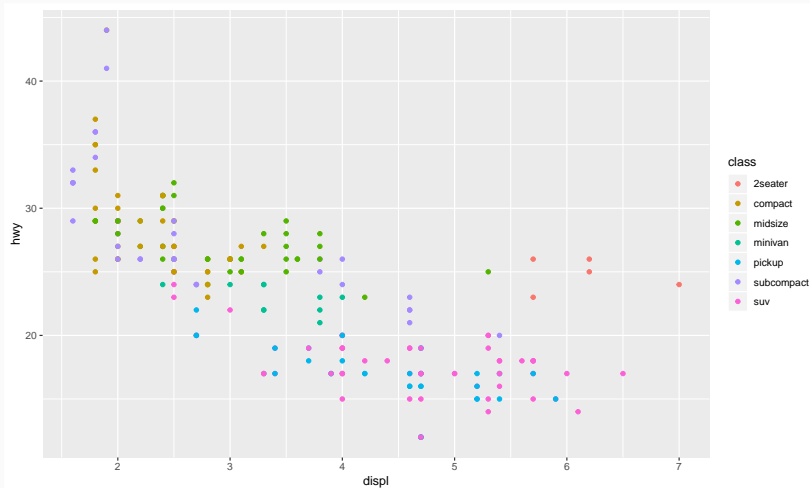
```
ggplot(data = <data> ,  
       mapping = aes(x = <variable>, y = <variable>)) +  
  geom_point(aes(col = <variable>)) +  
  geom_smooth() +  
  facet_grid(rows = vars(<variable>)) +  
  theme_minimal()
```

- aes variable to x- and y-axis
- geom... type of plot (point, density, bar, histogram, etc.)
- aes(col = <variable>) color points by scores/levels of
- facet_grid produce plot for each level vars(<variable>)
- theme... changes appearance of plot

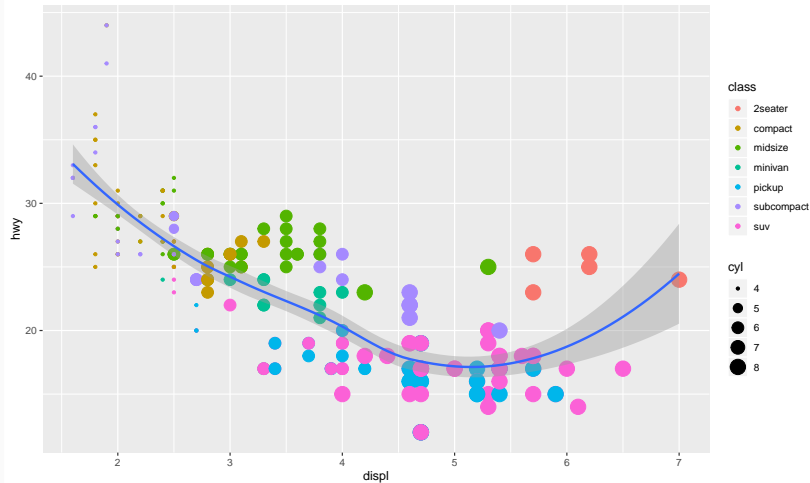
```
ggplot(data = mpg, aes(x = displ, y = hwy))
```



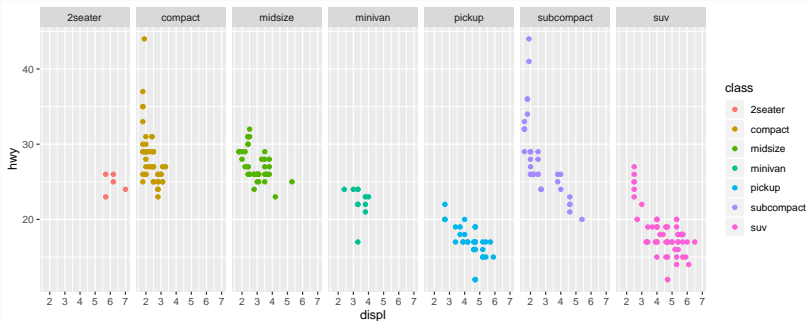
```
ggplot(data = mpg, aes(x = displ, y = hwy)) +  
  geom_point(aes(color = class))
```



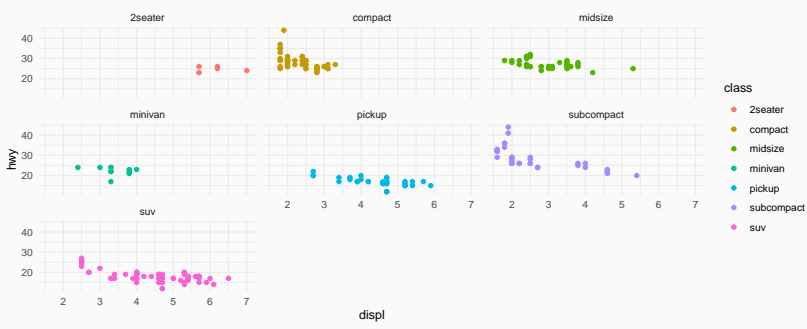

```
ggplot(data = mpg, aes(x = displ, y = hwy)) +  
  geom_point(aes(color = class, size = cyl)) +  
  geom_smooth()
```



```
ggplot(data = mpg, aes(x = displ, y = hwy)) +  
  geom_point(aes(color = class)) +  
  facet_grid(cols = vars(class))
```



```
ggplot(mpg, aes(x = displ, y = hwy)) +  
  geom_point(aes(color = class)) +  
  facet_wrap(vars(class)) +  
  theme_minimal()
```



Grammar of data manipulation

Prepare data for plotting with dplyr functions

- `filter()` select subgroup(s) of cases
- `select()` select variables
- `mutate()` make new variables
- `summarise()` derive summary statistics
- `group_by()` stratify on the levels of one or more variables
- `arrange()` sort the outcome

Example

Display mean of log of Sepal.Width for selected plant species excluding setosa of the data set iris.

1. Deselect setosa
2. Compute log of Sepal.Width
3. Compute group means
4. Display means ordered from low to high

This can be done in many ways, e.g.

- a. with traditional R functions
- b. with dplyr functions without the pipe operator
- c. with dplyr and pipe operator

Traditional R functions

```
newiris <- subset(iris, subset = Species != "setosa")
lg      <- log(newiris$Sepal.Width)
mean_lg <- aggregate(lg,
                      by = list(newiris$Species),
                      FUN = mean)
mean_lg[order(mean_lg$x), ]
```

```
##      Group.1      x
## 1 versicolor 1.012282
## 2  virginica 1.084180
```

- Not easy to oversee what is going on
- Pollution workspace with byproducts `newiris` and `sq`

With dplyr

```
arrange(  
  summarize(  
    mutate(  
      group_by(  
        filter(iris, Species != "setosa"),  
        Species),  
      lg = log(Sepal.Width)),  
    mean_lg = mean(lg)),  
  mean_lg)
```

```
## # A tibble: 2 x 2  
##   Species    mean_lg  
##   <fct>      <dbl>  
## 1 versicolor  1.01  
## 2 virginica   1.08
```

- Code has to be read from the inside out, but no workspace pollution

Using %>%

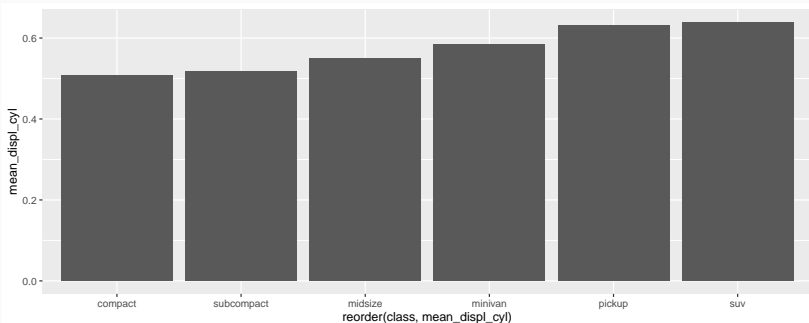
```
iris %>%  
  filter(Species != "setosa") %>%  
  mutate(lg = log(Sepal.Width)) %>%  
  group_by(Species) %>%  
  summarise(mean_lg = mean(lg)) %>%  
  arrange(mean_lg)
```

```
## # A tibble: 2 x 2  
##   Species    mean_lg  
##   <fct>      <dbl>  
## 1 versicolor  1.01  
## 2 virginica   1.08
```

- easy to follow and no workspace pollution

Combination ggplot and dplyr

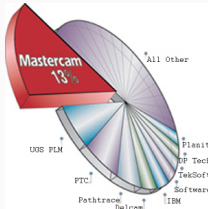
```
mpg %>%  
  filter(class != "2seater") %>%  
  group_by(class) %>%  
  mutate(displ_cyl = displ/cyl) %>%  
  summarise(mean_displ_cyl = mean(displ_cyl)) %>%  
  ggplot() +  
  geom_col(aes(x = reorder(class, mean_displ_cyl), y = mean_displ_cyl))
```



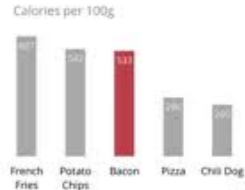
Tufte's Principles of Graphical Excellence

Proportionality, data-to-ink, chartjunk

- proportionality data/shape



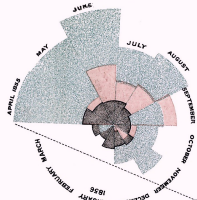
- maximize data-to-ink ratio, omit 'chart junk' (themes in ggplot)



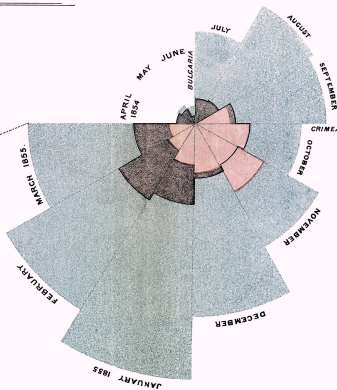
Proportionality principle

DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST.

2.
APRIL 1855 to MARCH 1856



1.
APRIL 1854 to MARCH 1855.



The Areas of the blue, red, & black wedges are each measured from the centre as the common vertex.

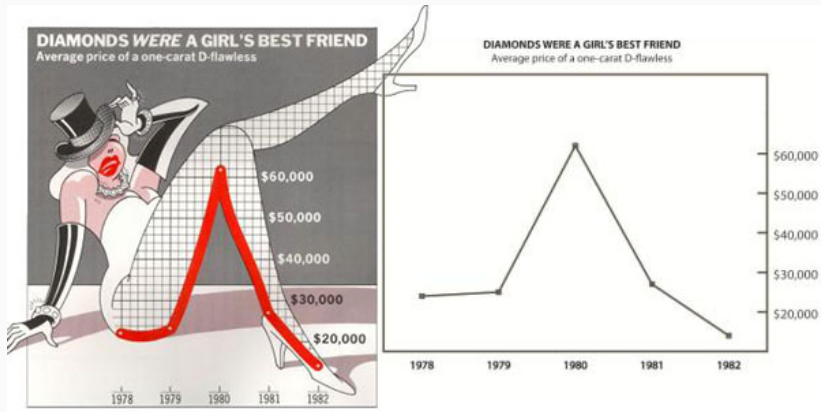
The blue wedges measured from the centre of the circle represent area for area the deaths from Preventible or Mitigable Zymotic diseases, the red wedges measured from the centre the deaths from wounds, & the black wedges measured from the centre the deaths from all other causes.

The black line across the red triangle in Nov^r 1854 marks the boundary of the deaths from all other causes during the month.

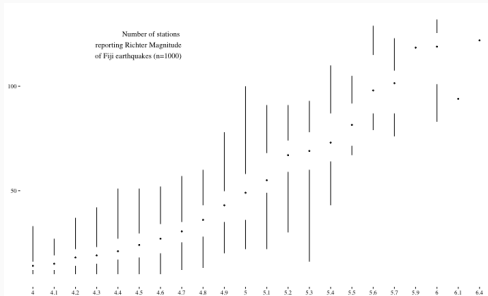
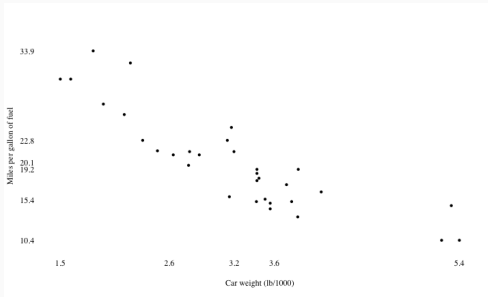
In October 1854, & April 1855, the black area coincides with the red; in January & February 1855, the blue coincides with the black.

The entire areas may be compared by following the blue, the red & the black lines enclosing them.

Chartjunk



Tufte in practice



Data visualization is a huge field

Stick to basic principles:

- Map data to aesthetics, geoms, scales, facets
- Which comparisons do I want?
- Maximize data-ink (within reason).

Lab 1A

Get acquainted with `dplyr` and `ggplot2`

1. Make plots with `ggplot2`
2. Combine plotting with `dplyr` functions

An R Markdown template is available

- Open the template in RStudio
- Save the file in an appropriate folder
- Read the instructions
- Fill the R chunks with the appropriate R code
- Test the R code in each chunk, and if it is error free
- Render an HTML file