Today: Critiquing Statistical Graphics
Introduction to Data
Graphics Principles
Friday: Introduction to R and Reproducibility
Monday: No class (Martin Luther King, Jr. Day)
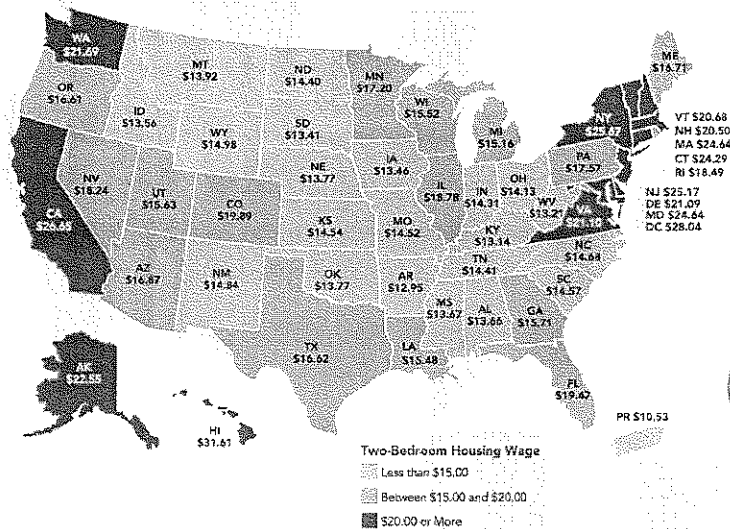
Sam Ventura
36-315

Department of Statistics
Carnegie Mellon University

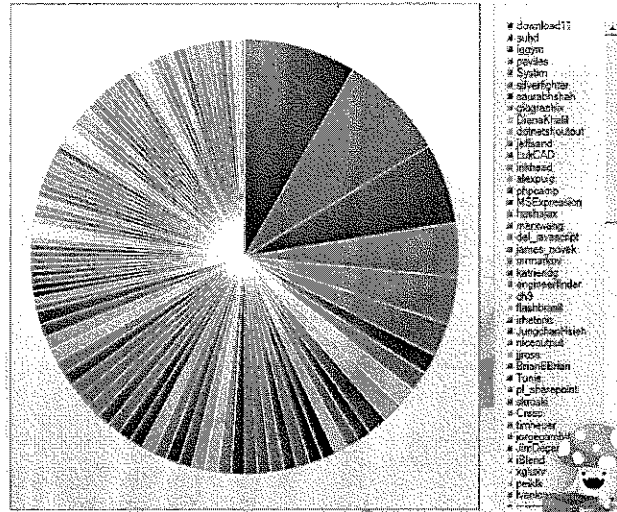January 13, 2016

## Hourly Wages to Afford Two-Bedroom Apartment



VT $20.68
NH $20.50
MA $24.64
CT $24.29
RI $18.49

NJ $25.17
DE $21.09
MD $24.64
DC $28.04

PR $10.53

Two-Bedroom Housing Wage
Less than $15.00
Between $15.00 and $20.00
$20.00 or More

*improve w/ color scale or just more options*

# Top 100 Tweeters

Pie = bad

## 100 Most Active Tweeters

# Never Make 3-D Pie Charts

## Microsoft Word Features By Version Added



- Word 1.0
- Word 1.1
- Word 2.0
- Word 6.0
- Word 95
- Word 97
- Word 2000
- Word 2002
- Word 2003

distortion

## What Is Data?

Information organized in some fixed easy-to-understand way

EX) Tweets

field goal attempts $\left\{ \begin{array}{l} made \\ missed \end{array} \right. \rightarrow$ 3pt sh %

Temperature

Censuses / Surveys → collect
info on population demographics, etc

## How Do We Describe Data?    mean, median, mode

→ rows

Two measurements used to describe datasets:

$n$ = # of obs, people, subjects, objects, etc

$p/d$ = # of covariates, variables, questions, etc
↳ columns

Data is usually in matrix form:

rows → observations

|  | $V_1$ | $V_2$ | ... | $V_d$ |
|------|------|------|------|------|
| $O_1$ | $X_1$ | | | |
| $O_2$ | $X_2$ | | | |
| $O_3$ | $X_3$ | | | |
| : | : | | | |
| : | | | | |
| $O_n$ | $X_n$ | | | |

single row has all answers
to all ?s from a single person

columns → variables

single column has all answers
from a single ? for all
people

Types of Data → *String, integer*

Categorical → *qualitative, describes qualities of obs*

*Ordered : strongly disagree, disagree, neutral, agree, SA*
  *educ. level*

*unordered : "nominal" → race, colors (sort of)*
  *names / general text, gender*

Continuous:
  *↳ real-valued, quantitative, numerical data*

  *Notation* $X = \{X_1, X_2, \ldots, X_d\},$

  *↓*
  *data / variable*

  $X_i \in \mathbb{R},$

  *→ double, int, float*

  $X \in \mathbb{R}^d$

7/13

# Graphics and Their Goals (from Tufte) → *feather of graphics*

Graphics: visually display measured quantities by combining points, lines, coordinate system, numbers, symbols, words, shading, color

Goals: show data!

- ➤ induce viewer to think about substance, not graphical methodology
- ➤ avoid **distorting** the data
- ➤ present numbers in small space
- ➤ make large, complicated datasets more coherent
- ➤ encourage comparison of different pieces of data
- ➤ reveal data at several levels of detail
- ➤ describe, explore, tabulate, or decorate
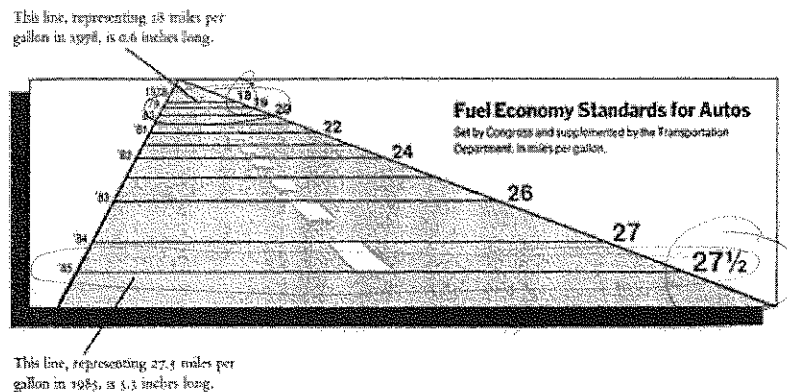- ➤ be closely integrated with statistical/verbal descriptions of dataset

Graphs that do not meet these goals are not successful

Graphs leading viewers to make misleading conclusions should be avoided

8/13

# Distortion

Visual representation of data is inconsistent with numerical representation

In other words: **The graph doesn't match the data**

This line, representing 18 miles per
gallon in 1978, is 0.6 inches long.

**Fuel Economy Standards for Autos**
Set by Congress and supplemented by the Transportation
Department, in miles per gallon.

22

24

26

27

27½

This line, representing 27.5 miles per
gallon in 1985, is 5.3 inches long.

Optimal: $LF \approx 1$

$LF > 1 \rightarrow$ enhance the effect
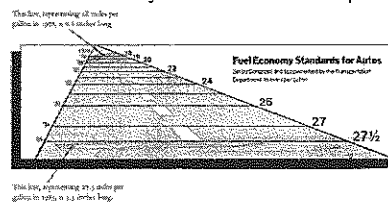
$LF < 1 \rightarrow$ decrease the effect

# Lie Factor

Tufte suggests optimizing the Lie Factor:

$$LF = \frac{\text{size of "effect" in graphic}}{\text{------------ in data}}$$

"effect" =
change in amount
of some feature
or variable

Fuel Economy Standards Example:

$\rightarrow$ Actual % increase (in data)

$$\frac{|27.5 - 18|}{18} \approx 0.528$$

graphical increase (in graph)

$$\frac{|5.3\,in - 0.6\,in|}{0.6\,in} = 7.83$$

$$LF = \frac{7.83}{0.528} = 14.83$$

Graphics should not draw the viewer's attention away from the data.
Extras get in the way.

**Note: Decoration does not refer to appropriate graph labeling.**
Labels should always be clear, detailed, and thorough.
Label key parts of the data. Add text explanations if necessary.

**Data Ink should primarily present information about the data:**
the non-erasable, non-redundant core of a graphic

Tufte suggests using the *data-ink ratio*:

$$DI = \frac{data\ ink}{total\ ink\ on\ graphic}$$

% of ink devoted to non-redundant / useful information.   11/13

Ideally → Maximize DI (max = 1)
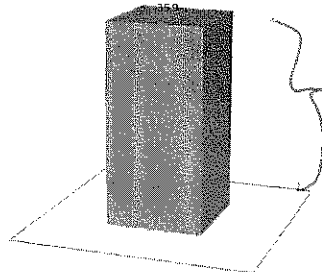won't quite get to 1, because of
axes, grid lines etc

Two ways to increase the proportion of data-ink:

**Remove non-data-ink:**
  ↳ Ink that does not depict statistical info
In class wednesday 1/20, hands on map graphic

**Remove redundant data-ink:**
  → Ink that is unnecessarily redundant/repetitive.

Indications of height:
1) height of front-left line on bar
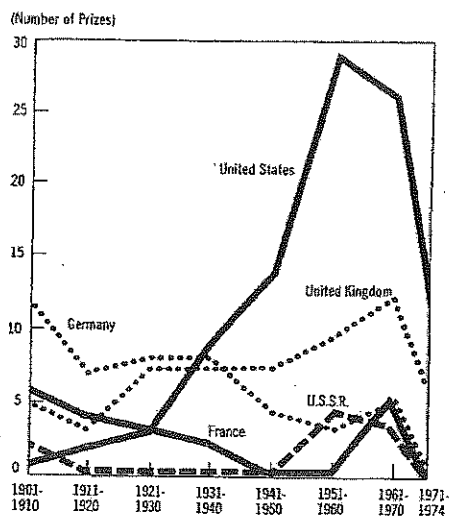2) height of front-right line on bar
3) — — — back-right — — —
4) position of number
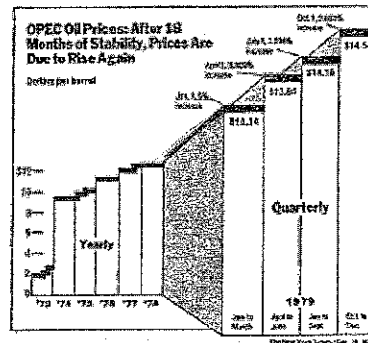5) value of number
6) position of top line (front)
7) position of top line (back)

8) height of shading
9) etc. —

Nobel Prizes Awarded in Science, for Selected Countries, 1901-1974

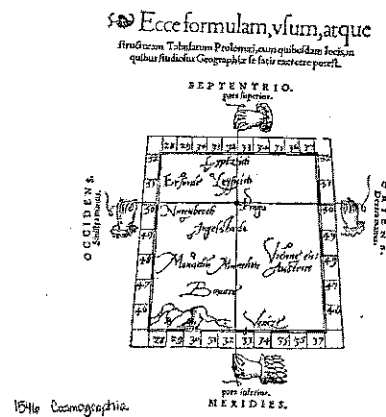# Graph Principles (Tufte)

Minimize **Design Variation**:   Changes in Design may imply Changes in Data



New York Times, December 29, 1978, p. D-3.

# Graph Principles (Tufte)

Maximize **Data-Ink Ratio:** Ink in graphic should primarily show data

# Size and Types of Data

Two measurements used to describe data:

$n$ # of obs    (people, subjects, objects)

$p/d$   # dim or var  (variables or questions)

Data usually in matrix form:

$$X_1 O_1 \quad X_2 O_2 \quad \cdots \quad X_n O_n$$

$$\begin{array}{cccc} V_1 & V_2 & \cdots & V_d \\ P_1 & P_2 & & P_d \end{array}$$

rows → observations
    answers to all the questions for one obs

cols → variables/questions
    answers for all obs to one question

In R: $dim(data) = n \times d$

Often the # of colums is the important piece of the dim; how we graph/visualize depends on $d$ (also to some extent on the $n$; more later)

Two major types we'll be working with: *categorical* and *continuous*

\* <u>Categorical:</u> Qualitative, describes qualities of the obs          string, int (factor?)

  <u>Non-ordered categ</u>    ex?
    - Favorite Ice Cream                 can be numerical or text
    - Nationality

  <u>Ordinal categ</u>: like having levels/factors    ex?
            categ have an inherent order

  e.g. Likert Scale: Strongly Agree, Agree, Neutral, Disagree, Strongly Disagree

\* <u>Continuous:</u>                                     int, double
    <u>quantitative, numerical data</u>
    often see notation like

$$X = \{ \underline{x_1}, x_2, \ldots, x_n \} \in \mathbb{R}^d \leftarrow \text{real-valued number}$$
            ↓                                    $d$-dim space
          $n \times 1$    ∟ vector of $d$ values, one for each var

# Distortion

Visual representation of data is NOT consistent with the Numerical representation; *i.e. Graph doesn't match Data*

Graphics strongly dependent on visual perception of viewer.

Experiments have shown relationships between numerical measures and perceived measures. *people look at different shapes, areas, lines, etc*
*& guess the length, area → big range of responses*

Area Example:

*Perceived Area (may) grow more slowly than Actual area*

$$PA = AA^x \qquad X = 0.8 \pm 0.3$$

Can't design graph for each viewer; What should we do?
*Some say Tables for 20 it's or less → big debate, post discussion*
*well-labeled graph is fine as well; graphs for large sets & higher dim*

## Lie Factor

Tufte suggests optimizing the Lie Factor:

*(effect - change/feature of interest)*

$$Lie\ Factor = \frac{size\ of\ effect\ in\ graphic}{size\ of\ effect\ in\ data}$$

*Optimal = 1* $\qquad\qquad$ *log LF = 0*

$LF < 1$ *graphic diminishes the effect* $\qquad < 0$

$LF > 1$ *graphic enhances the effect* $\qquad > 0$
$\qquad$ *(more common)*

*Fuel Economy Standards Example:* U.S. Congress and Department of Transportation set a series of fuel economy standards to be met by automobile manufacturers.

*Actual % Increase* $\qquad$ *Graphical Increase* $\qquad LF = \frac{7.83}{0.528} = 14.83$

$\frac{27.5 - 18}{18} = 0.528 \qquad\qquad \frac{5.3 - 0.6}{0.6} = 7.83$

# "Decorating"/ Data Ink

Graphics should not draw the viewer's attention away from the data.
Extras get in the way.

**Note: Decoration does not refer to appropriate graph labeling.
Labels should always be clear, detailed, and thorough. Label
key parts of the data. Add text explanations if necessary.
More later.**

Data Ink should be primarily present information about the data:
the Non-Erasable core of a graphic, Non-redundant Ink.

Tufte suggests using the *data-ink ratio*:

$$DI: \frac{\text{data ink}}{\text{total ink in graphic}} = \% \text{ ink devoted to non-redundant display of info about data}$$
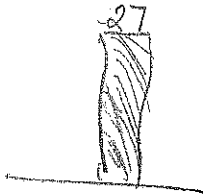
1 - % of graphic that could be erased

We want to maximize this ratio (within reason).

max of 1; 1 not realistic; why?
grid, border, etc

Two ways we can increase the proportion of data-ink:

*Remove non-data-ink:* ink that doesn't depict statistical info
see Playfair example; re disembodied hands

*Remove redundant data-ink:* six indications of height:
1) ht of left line
2) ht of right line
3) height of shading
4) position of top line
5) position of # 6) actual #

<u>Removals should be done within reason; some redundancy will remain.</u>