

STA291

Fall 2008

1

LECTURE 5
5 FEBRUARY 2009

Itinerary

2

- **2.3 Graphical Techniques for Interval Data (mostly review)**
- **2.4 Describing the Relationship Between Two Variables**
- **3 Art and Science of Graphical Presentations**

Administrative Notes and Homework

3

- Use the Study Tools at Cengage Now, click on the “Personalized Study Book” with the same title page as our textbook, and work through “Chapter 2 – Graphical and Tabular Descriptive Techniques”. This involves taking a pre-test, working through a personalized study plan, and then taking a post-test.
- Please read Chapter 3 about the Art & Science of graphical presentations.
- Suggested problems from the textbook (not graded, but good as exam preparation): 2.74, 2.76, 3.12

Review: Graphical/Tabular Descriptive Statistics

4

- Summarize data
- Condense the information from the dataset
- Always useful: Frequency distribution
- Interval data: Histogram (Stem-and-Leaf?)
- Nominal/Ordinal data: Bar chart, Pie chart

Stem and Leaf Plot

5

- Write the observations ordered from smallest to largest (stems, certainly)
- Each observation is represented by a stem (leading digit(s)) and a leaf (final digit)
- Looks like a histogram sideways
- Contains more information than a histogram, because every single measurement can be recovered

Stem and Leaf Plot

6

- Useful for small data sets (<100 observations)
 - Example of an *EDA*
- Practical problem:
 - What if the variable is measured on a continuous scale, with measurements like 1267.298, 1987.208, 2098.089, 1199.082, 1328.208, 1299.365, 1480.731, etc.
 - Use common sense when choosing “stem” and “leaf”

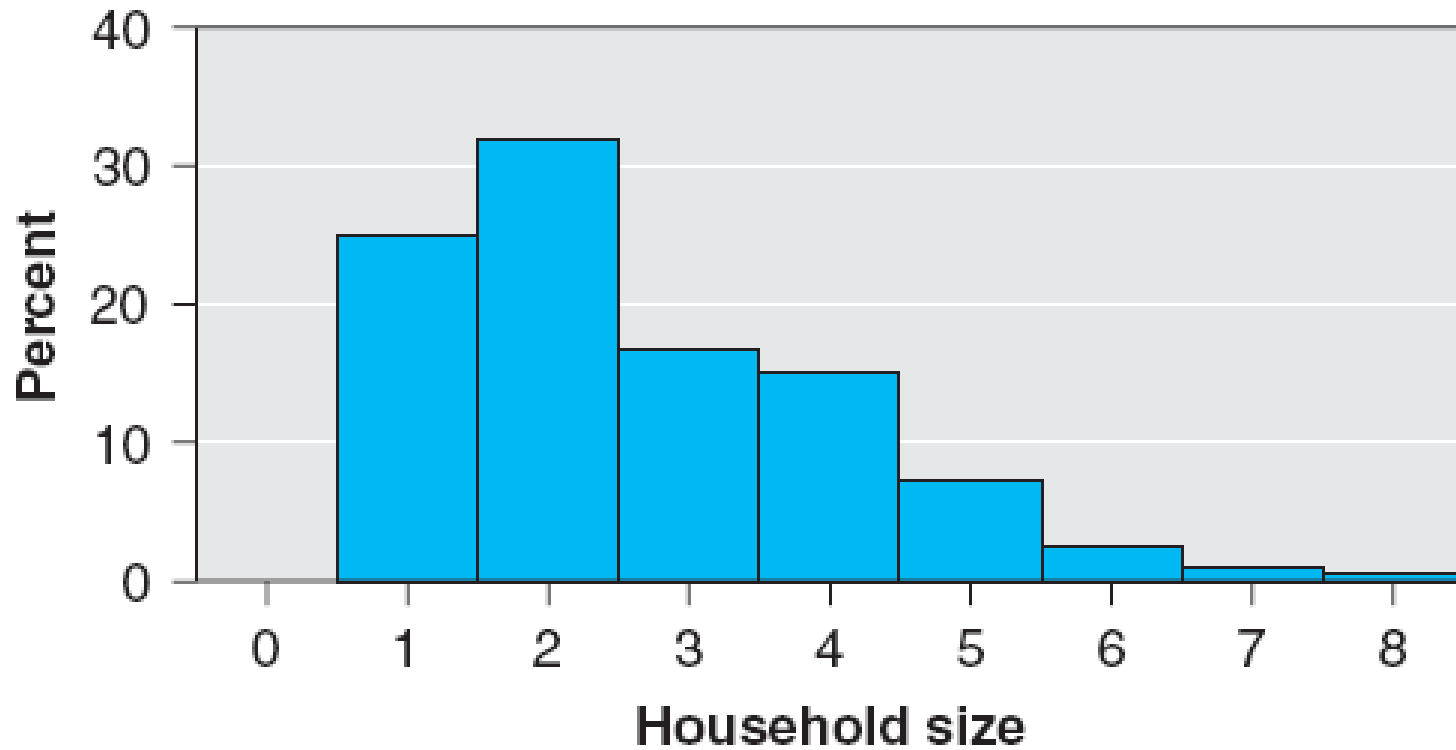
Stem-and-Leaf Example: Age at Death for Presidents

7

PRESIDENT	AGE	PRESIDENT	AGE	PRESIDENT	AGE
Washington	67	Fillmore	74	Roosevelt	60
Adams	90	Pierce	64	Taft	72
Jefferson	83	Buchanan	77	Wilson	67
Madison	85	Lincoln	56	Harding	57
Monroe	73	Johnson	66	Coolidge	60
Adams	80	Grant	63	Hoover	90
Jackson	78	Hayes	70	Roosevelt	63
Van Buren	79	Garfield	49	Truman	88
Harrison	68	Arthur	56	Eisenhower	78
Tyler	71	Cleveland	71	Kennedy	46
Polk	53	Harrison	67	Johnson	64
Taylor	65	McKinley	58	Nixon	81
				Reagan	93

Example (Percentage) Histogram

8

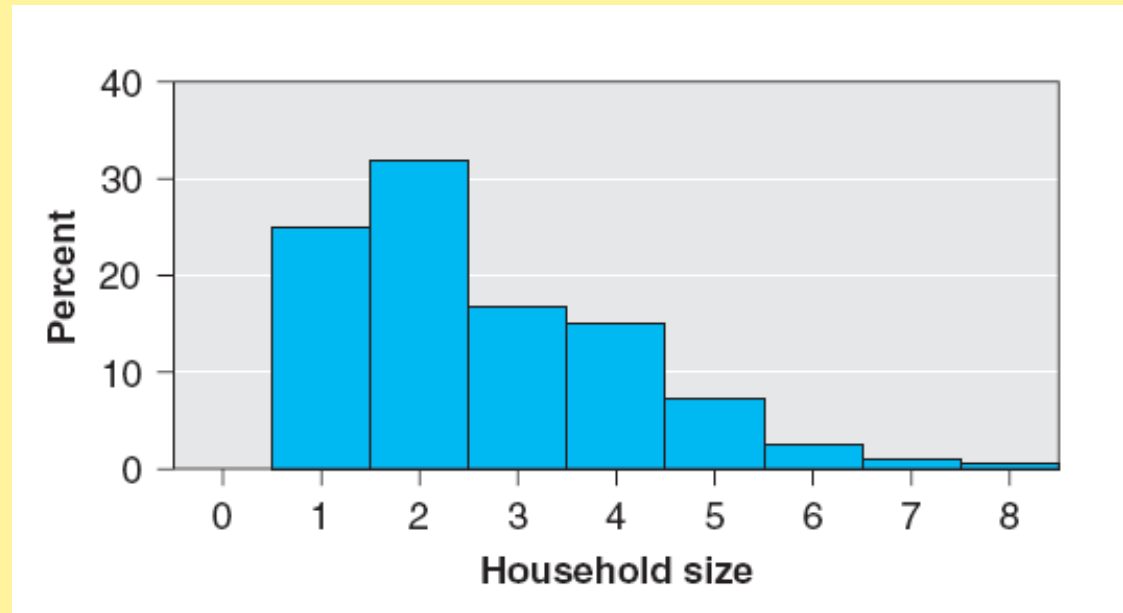


Side by side?

9

Annual Earnings (Thousands of Dollars)

Stem	Leaf
1	6,7
2	9,9,9,2,7,7,8,3,5
3	4,5,1,2,6,6,3,9
4	6,9,7,7,7,5,9,8,6,3,0,3,4,8,3
5	3,5,3,7,7,9
6	3,2,2,4,4,0,0
7	7,4



Similarities/differences?

Sample/Population Distribution

10

- Frequency distributions and histograms exist for the population as well as for the sample
- Population distribution vs. sample distribution
- As the sample size increases, the sample distribution looks more and more like the population distribution

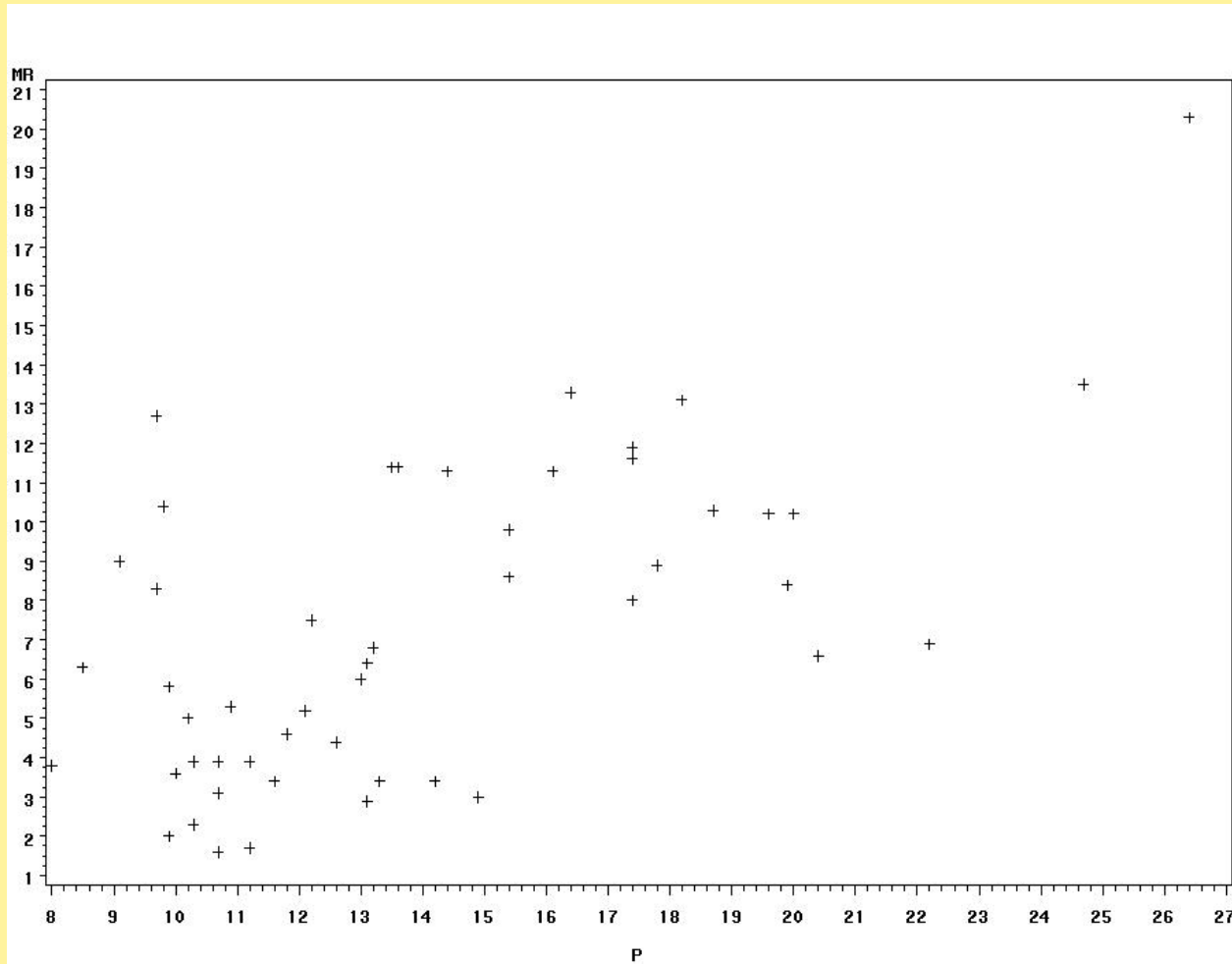
Describing Distributions

11

- Center, spread (numbers later)
- Symmetric distributions
 - Bell-shaped or U-shaped
- Not symmetric distributions:
 - Left-skewed or right-skewed

On to examining two variables for relationships . . .

12



Describing the Relationship Between Two Nominal (or Ordinal) Variables

13

Contingency Table

- Number of subjects observed at all the combinations of possible outcomes for the two variables
- Contingency tables are identified by their number of rows and columns
- A table with 2 rows and 3 columns is called a 2 x 3 table (“2 by 3”)

2 x 2 Contingency Table: Example

14

- 327 commercial motor vehicle drivers who had accidents in Kentucky from 1998 to 2002
- Two variables:
 - wearing a seat belt (y/n)
 - accident fatal (y/n)

		Accident Fatal		
		Yes	No	
Seat Belt	Yes	30	212	242
	No	33	52	85
		63	264	327

2 x 2 Contingency Table: Example, cont'd.

15

- How can we compare fatality rates for the two groups?
- Relative frequencies or percentages within each row
- Two sets of relative frequencies (for *seatbelt=yes* and for *seatbelt=no*), called **row relative frequencies**
- If seat belt use and fatality of accident are related, then there will be differences in the row relative frequencies

Row relative frequencies

16

- Two variables:
 - wearing a seat belt (y/n)
 - accident fatal (y/n)

		Accident Fatal		
		Yes	No	
Seat Belt	Yes			100
	No			100
				100

Describing the Relationship Between Two Interval Variables

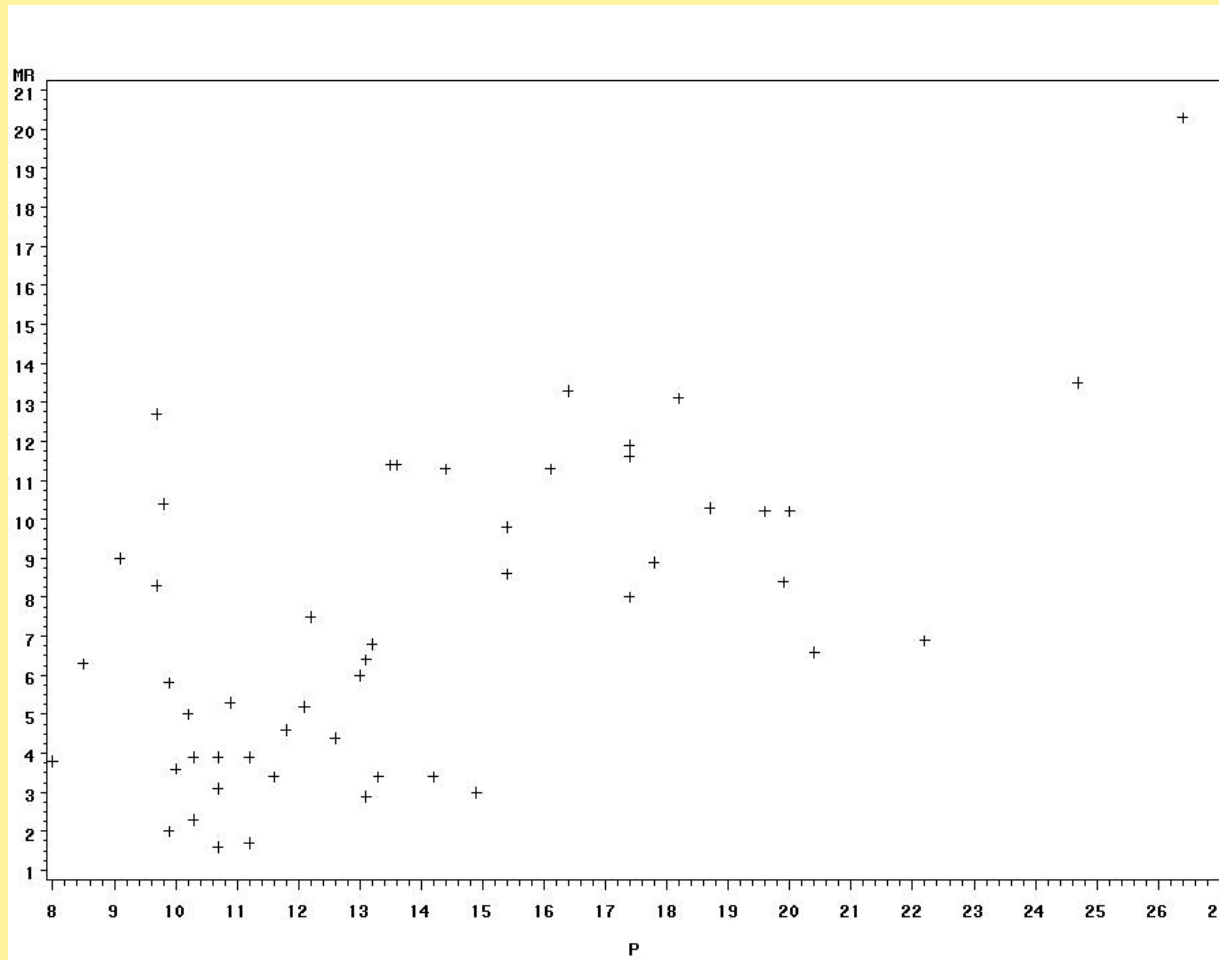
17

Scatter Diagram

- In applications where one variable depends to some degree on the other variables, we label the dependent variable Y and the independent variable X
- Example:
Years of education = X
Income = Y
- Each point in the scatter diagram corresponds to one observation

Scatter Diagram of Murder Rate (Y) and Poverty Rate (X) for the 50 States

18



3.1 Good Graphics ...

19

- ... present large data sets concisely and coherently
- ... can replace a thousand words and still be clearly understood and comprehended
- ... encourage the viewer to compare two or more variables
- ... do not replace substance by form
- ... do not distort what the data reveal
- ... have a high “data-to-ink” ratio

CARTE FIGURATIVE des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite.

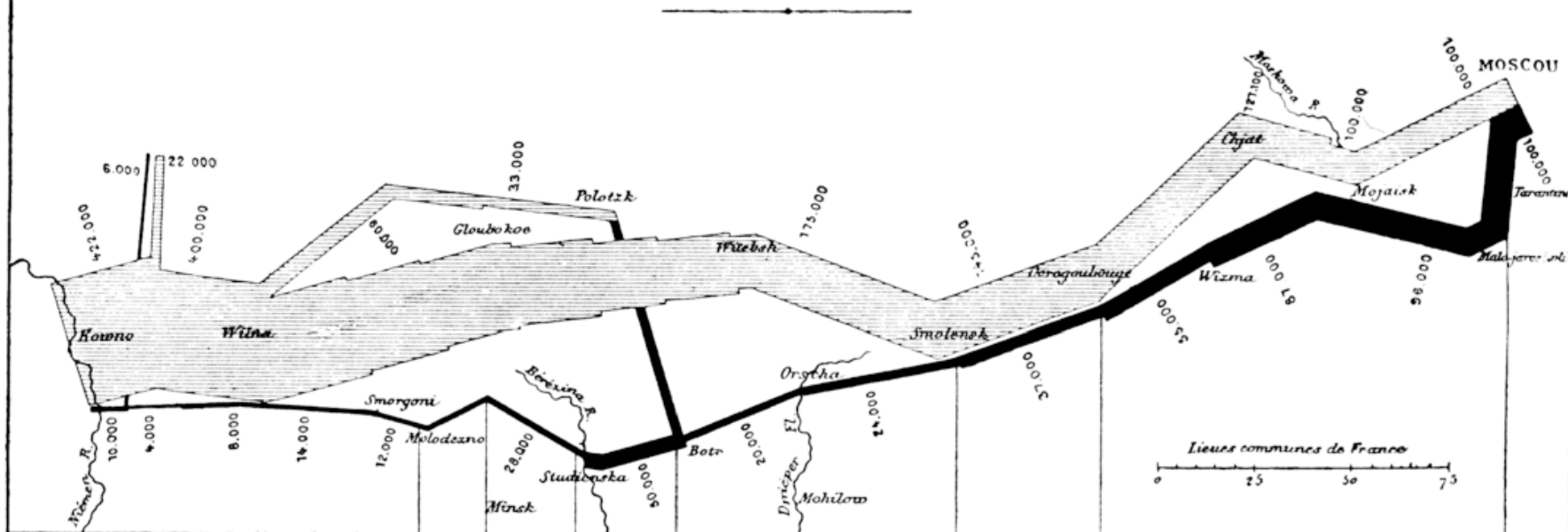
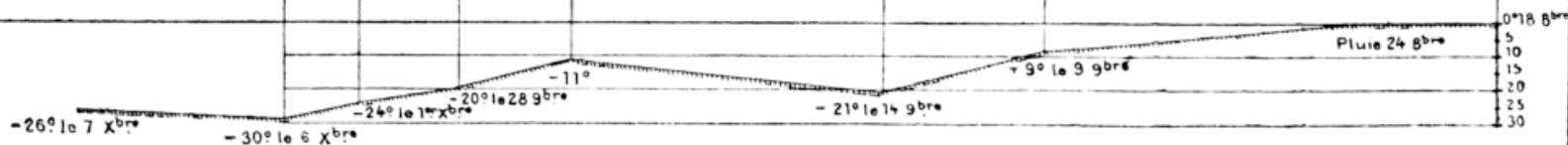


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro



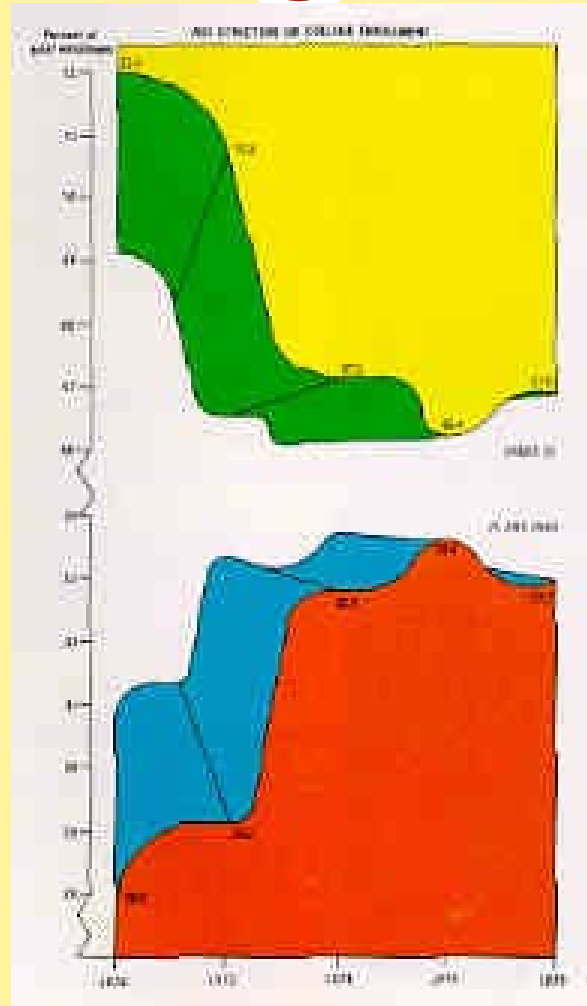
3.2 Bad Graphics...

21

- ...don't have a scale on the axis
- ...have a misleading caption
- ...distort by stretching/shrinking the vertical or horizontal axis
- ...use histograms or bar charts with bars of unequal width
- ...are more confusing than helpful

Bad Graphic, Example

22



Attendance Survey Question #5



- On an index card
 - Please write down your name and section number
 - Today's Question: