# Number Visualization

## Giuseppe Santucci

### University of Rome "La Sapienza"
### santucci@dis.uniroma1.it

Thanks to:
Ross Ihaka (very inspiring lectures)

# Number visualization ?

- Obviously information visualization is, in general, about numbers

- In some cases, however, the numerical part is relevant, and the use of tables, graphs and other visual means to communicate **quantitative** information is commonplace in business today (pie chart, diagrams, boxplots, scatterplots, etc.)

- Actual software applications allows for easy (?) development of different typologies of charts

- I will discuss the basic relationships and the logical steps that allows for moving from quantitative data to suitable visualizations.

# Types of Data

- ## Quantitative (allows arithmetic operations)
  - *123, 29.56, …*

- ## Categorical (group, identify & organize; no arithmetic)
  - ### Nominal (name only, no ordering)
    - *Direction: North, East, South, West*
  - ### Ordinal (ordered, not measurable)
    - *First, second, third …*
    - *Hot, warm, cold*
  - ### Interval (starts out as quantitative, but it is made categorical by subdividing into ordered ranges)
    - *0-999, 1000-4999, 5000-9999, 10000-19999, …*
  - ### Hierarchical (successive inclusion)
    - *Region: Continent > Country > State > City*
    - *Animal > Mammal > Horse*

- ## Relationships
  - Correlation
  - ....

# Table and graphs,

- **Table** and **graphs** are widely used to communicate quantitative information
- The goals of presenting quantitative data are
  - Analyzing
  - Monitoring
  - Planning
  - Communicating
- Remember that we are dealing with data that is
  - Quantitative
  - Categorical
- Not all numbers carry quantitative information
  - Categorical intervals
  - IDs (e.g., order number)
- The problem is to map such data to the right visualization, and clear indications about that exists

# uhmmm...



- Boring ?

- I do agree !

- I have changed my mind !

- It is plenty of books that teach about quantitative data and how to represent it (see references).

- Read all of them!  I'll go for another way...

# Outline
## (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues
  - Quantitative perception (basic rules)
  - The role of interaction
- Two examples for IR

# A lotto game

- Forms of lotto are played world-wide and many people have theories about how to make money at the game

- User task ? ---> Money !!!

- We will examine a particular lotto game, to see whether it might be possible to play it profitably

- The game we'll look at is the daily pick-it lottery run by the state of New Jersey in the USA

# Lotto rules

- Each player selects a number between 000 and 999

- A winning number is selected by independently picking three digits between 0 and 9 at random

- All players that hold the winning number split the prize money for the game

# Available data

- The results of the games (winning number and winning amount) are publicly available

- Does this data contain information which will enable us to choose a profitable strategy for this game?

- We will use the results of 254 consecutive games to look for a profitable strategy

# The data (254 values)

(winning number, winning amount)

- (810, $190.0), (156, $120.5), (140, $285.5), (542, $184.0), (507, $384.5),
- (972, $324.5), (431, $114.0), (981, $506.5), (865, $290.0), (499, $869.5),
- (020, $668.5), (123, $83.0), (356, $188.0), (015, $449.0), (011, $289.5),
- (160, $212.0), (507, $466.0), (779, $548.5), (286, $260.0), (268, $300.5),
- (698, $556.5), (640, $371.5), (136, $112.5), (854, $254.5), (069, $368.0),
- (199, $510.0), (413, $102.0), (192, $206.5), (602, $261.5), (987, $361.0),
- (112, $167.5), (245, $187.0), (174, $146.5), (913, $205.0), (828, $348.5),
- (539, $283.5), (434, $447.0), (357, $102.5), (178, $219.0), (198, $292.5),
- (406, $343.0), (079, $332.5), (034, $532.5), (089, $445.5), (257, $127.0),
- (662, $557.5), (524, $203.5), (809, $373.5), (527, $142.0), (257, $230.5),
- (008, $482.5), (446, $512.5), (440, $330.0), (781, $273.0), (615, $171.0),
- (231, $178.0), (580, $463.5), (987, $476.0), (391, $290.0), (267, $176.0),
- (808, $195.0), (258, $159.5), (479, $296.0), (516, $177.5), (964, $406.0),
- (742, $182.0), (537, $164.5), (275, $137.0), (112, $191.0), (230, $298.0),
- (310, $110.0), (335, $353.0), (238, $192.5), (294, $308.5), (854, $287.0),
- (309, $203.5), (026, $377.5), (960, $211.5), (200, $342.0), (604, $259.0),
- (841, $231.0), (659, $348.0), (735, $159.0), (105, $130.5), (254, $176.0),
- (117, $128.5), (751, $159.0), (781, $290.0), (937, $335.0), (020, $514.0),
- (348, $191.0), (653, $304.5), (410, $167.0), (468, $257.0), (077, $640.0),
- (921, $142.0), (314, $146.0), (683, $356.0), (000, $96.0), (963, $295.0),
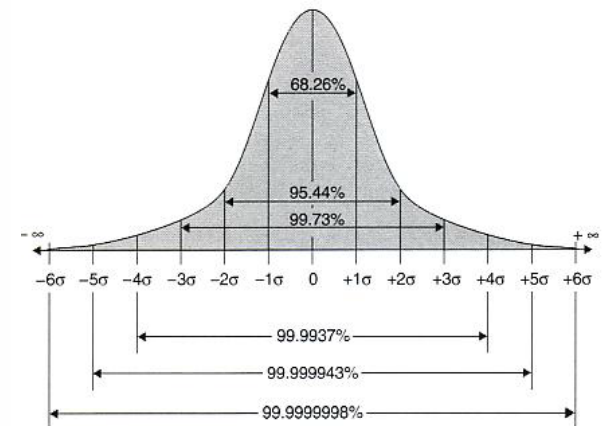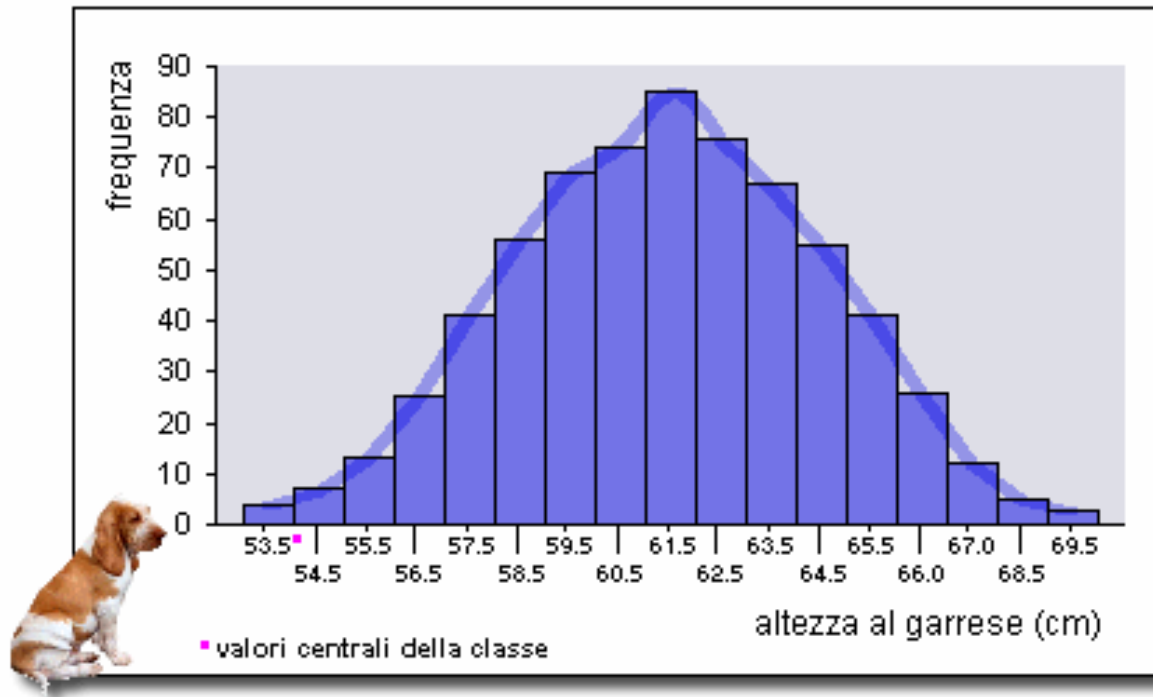
# Visualizing the data

- Humans can really only make sense of three or four numbers at a time

- By representing the values in a graphical form we make it easier to handle large numbers of values

- Using visualizations should make it possible to learn more about this data

- We have NOT to **lie** or make **noise** !!!

# User task and visualization

- One approach to making money at "Pick It" is to try to select numbers which are more likely to win

- We can look at the distribution of the winning numbers to see whether some ranges of values are more like to produce a winner than others

- One way to do this is to produce a histogram of the winning numbers

# Histogram example



Altezza al garrese di 659 cani di razza "Bracco italiano". Istogramma.

bin

# Data distribution



What can we infer from this histogram?

(Is the bin size ok?)

# Analysis

- It looks there tend to be more winners in the region from 100 to 300 than in other regions

- This suggests that we might be best to choose numbers in this range

Do you agree ?

# We are telling lies...

## (wrong number understanding)

- Even if the winning numbers are chosen randomly, we can expect some "random variability" in a sample

- To judge the significance of what we see in the histogram we have to recall some formal statistical theory

# The mean is not enough !

- There are 254 values. We would expect the number of values in each cell to be approximately: 25.4 = 254/10
- Such a number is a random variable as well, with normal distribution
- 95% of the observations fall within +/- $2\sigma$

# Better number visualization



- Variance <u>analysis</u> AND <u>visualization</u>

# Conclusions and new task

- Winning numbers are totally random

- It makes no sense to look for a " lucky " number

- However, we can change our task:

  - to increase the amount won !

- So we study the distribution of winning amount

# New visualization

# Looking for new insights

- The histogram shows that there is a wide (more than $2\sigma$) range amounts won in the game

- It *might* be possible to choose the numbers which win larger amounts

- We search for relationship between ticket number and winning amount

- A scatter plot is the natural way to look for such a relationship.

# New visualization

# Insights from the scatterplot

- The winning amounts in a band to the left of the plot appear to generally be higher than those in the rest of the plot

- We can investigate this further by separating the numbers into groups according to the first digit of the ticket number and drawing box plots for each group
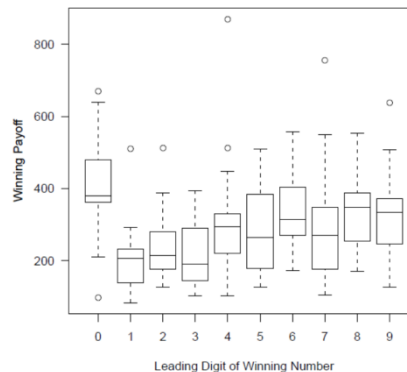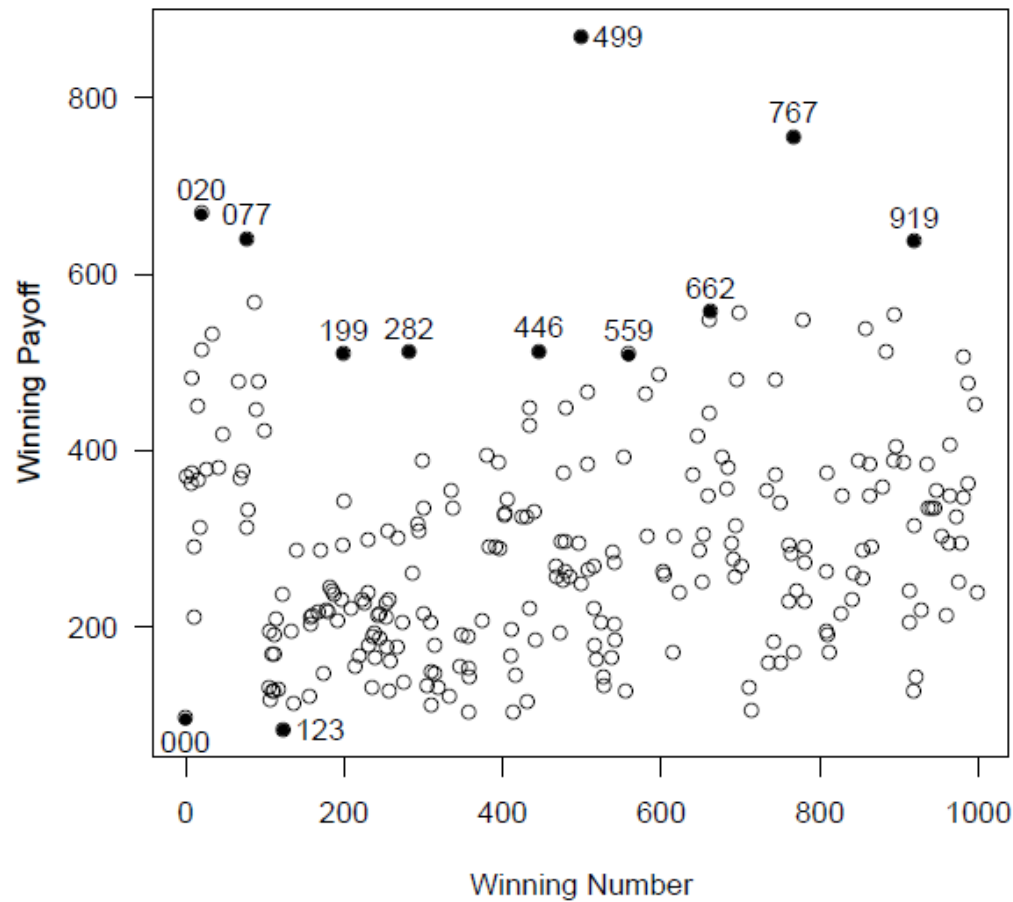
# Boxplot

# Lottery's boxplots



Leading Digit of Winning Number

# New insights

- Tickets with a leading zero digit clearly tend to produce larger winnings

- It is also apparent that there are some very large and some very small winning amounts

- It is probably of interest to identify the ticket numbers corresponding to these extremes

# High and low winning numbers

# Lotto strategy

- While winning numbers are non predictable, <u>players' choices are</u>!

- Choose numbers which are less likely to be chosen by other players
- Then, when you win (if), you will tend to win more
- Possible ways to choose:
  – Choose a number with a leading zero
  – Choose a number with repeated digits
  – Avoid "obvious" numbers like, e.g. 000, 123, 246, . . .

# Lessons learned

- Define clearly the task
- Use basic visualizations
  - bar charts
  - scatterplots
  - boxplots
- Be ready to switch among them
- Look for precise values when needed
- Do not lie !

# Outline
## (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues
  - Quantitative perception (basic rules)
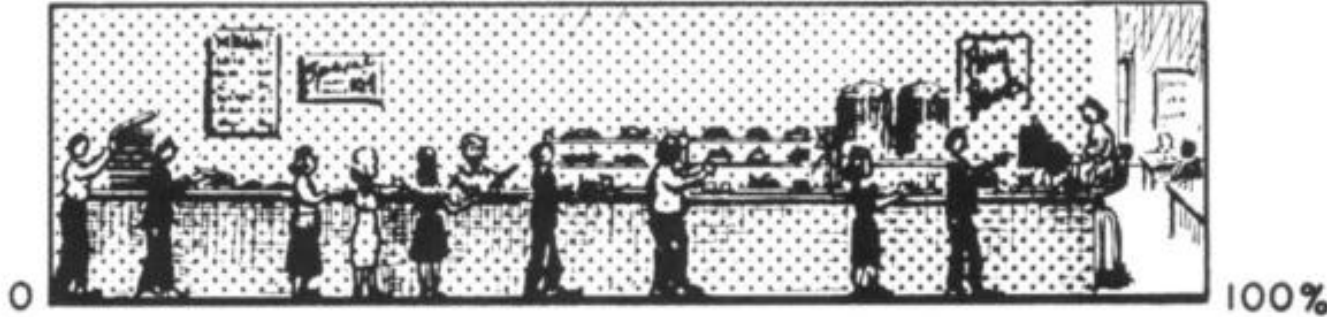  - The role of interaction
- Two examples for IR

# Rule 0:
# Do not use diagrams when handling few numbers

- It does not make sense to use graphs to display very small amounts of data

- The human brain is quite capable of grasping one two, or even three values

# Rule 0 violation (and also rule 2)



The Company Cafeteria was used by 9 Out of 10 Employees during the Fiscal Year 1949

0                                                    100%

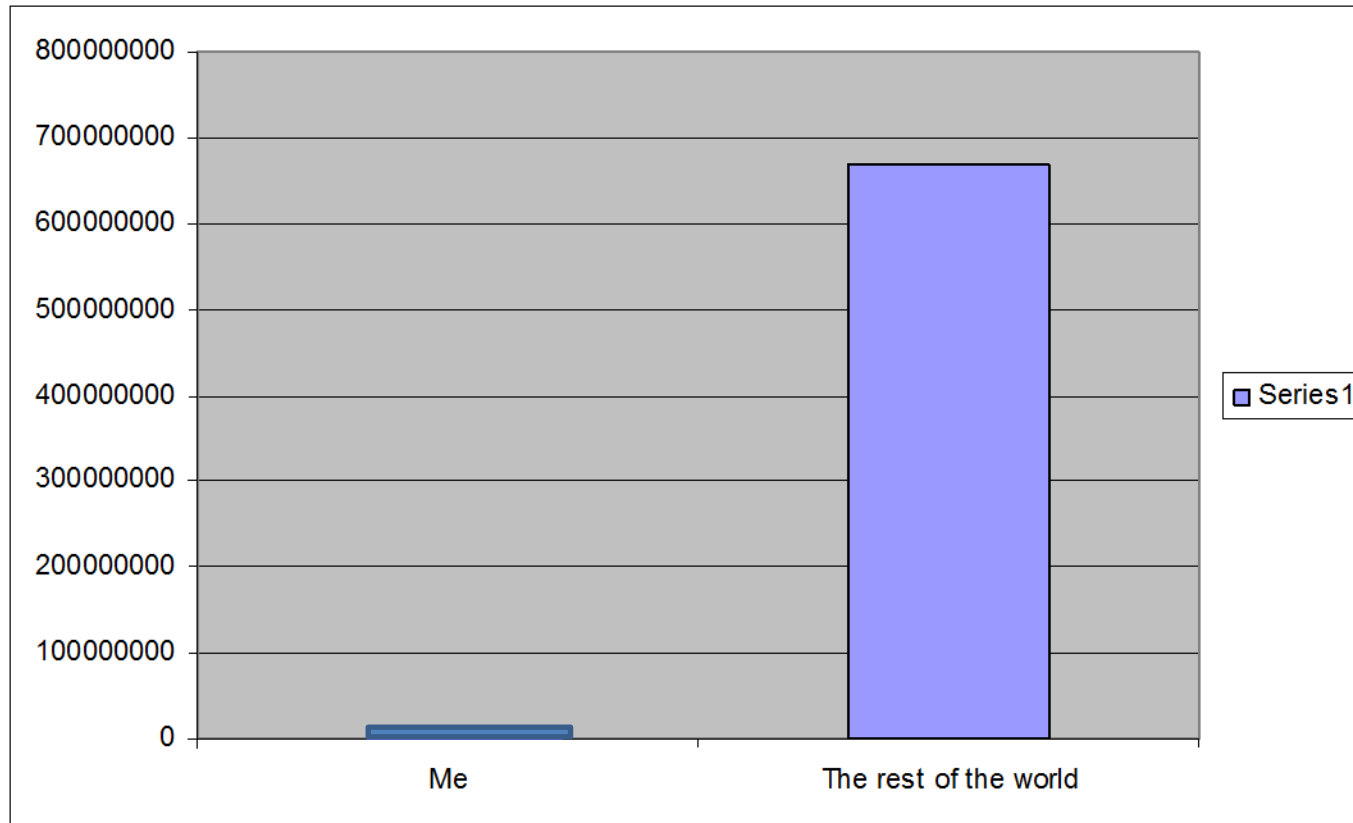Source: COMPANY REPORTS

**90%**

# Rule 0 violation

**Class Gender Breakdown**



**Male      60%**
**Female 40%**

# Rule 1:
# Insure data quality / significance

- Graphs are only as good as the data they display

- No amount of creativity can produce a good graph from dubious or non relevant data

# Rule 1 violation

# Rule 1 violation (and also rule 0)



Not very significant data but good example of distortion

# Rule 2:
# Insure chart simplicity

- Graphs should be no more complex than the data which they portray

- Unnecessary complexity can be introduced by
  - irrelevant decorations
  - colors
  - 3d effects
  - ...
- These are collectively known as "chartjunk"

- For a very comprehensive set of chartjunk effects look at Microsoft Excel
  - the later the version the larger the set !

**Age structure of College enrollment
(percentage of enrolled people above 25 years)**

# Rule 2 violation (and also rule 3)

- A very good bad example!
- Only 5 numbers on it but
  - 4 meaningless colors
  - useless 3D
  - useless axes split
  - confusing and wrong visual attributes (size)
  - nonsense interpolation
- Designers of this graph are now working in the Microsoft Excel's team, inspiring the new Excel's versions ...

*American Education Magazine*

38

# The same data...



**Age Structure of College Enrolment**

Percent of Total Enrolment, Aged 25 and Over

# The same data...

| Year | Percentage above 25 |
|------|---------------------|
| 1972 | 28.0 |
| 1973 | 29.2 |
| 1974 | 32.8 |
| 1975 | 33.6 |
| 1976 | 33.0 |

# The same data...



**Age Structure of College Enrolment**

Percent of Total Enrolment, Aged 25 and Over

# Rule 2 violation



Earnings Per Share And Dividends
(Dollars)

- Why 3D?
- The extra dimension used in this graph has confused even the person who created it..

*The Washington Post*, 1979

# The same data...



**Earnings Per Share and Dividends**

Legend:
- Earnings
- Dividends

| Year | Earnings | Dividends |
|------|----------|-----------|
| 1972 | 1.53 | 1.02 |
| 1973 | 1.71 | 1.08 |
| 1974 | 1.63 | 1.16 |
| 1975 | 1.72 | 1.16 |
| 1976 | 1.82 | 1.28 |
| 1977 | 1.7 | 1.34 |

Dollars

# Rule 3:
# Do not distort data

- Graphs should not provide a distorted picture of the values they portray

- Distortion can be:
  - deliberate
  - accidental

- Of course, it could be useful to know how to produce a graph which bends the truth...

# Rule 3 violation



FACULTIES

At a very quick glance:
- balanced faculty population
- most male students

What's wrong with this graph?

# The truth : population size



**Faculty Size**

# The truth : male /female ratio



**Percentage of Female Students**

# In other cases distortion is ok...

# The lie factor

- Edward Tufte of Yale University has defined the "lie factor" as a measure of the amount of distortion

$$\text{Lie Factor} =$$

size of effect in graphic / size of effect in data

- If the lie factor is greater than 1, the graph is exaggerating the size of the effect

# Measuring distortion through the lie factor (miles per gallon across years)



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.

1978 '79 '80 '81 '82 '83 '84 '85

18 19 20 22 24 26 27 27½

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

$$\text{Data Effect} = \frac{27.5 - 18}{18} = 0.53, \qquad \text{Graph Effect} = \frac{5.3 - .6}{.6} = 7.83,$$

$$\text{Lie Factor} = 14.8$$

# The same data with lie factor=1 (and following the previous roles)



**Required Fuel Economy Standards: New Cars Built from 1978 to 1985**

# Common sources of distortion

- The use of 3 dimensional "effects" is a common source of distortions in graphs (and of occlusion)

- Another common source is the inappropriate (or deliberate?) use of linear scaling when using area or volume to represent values
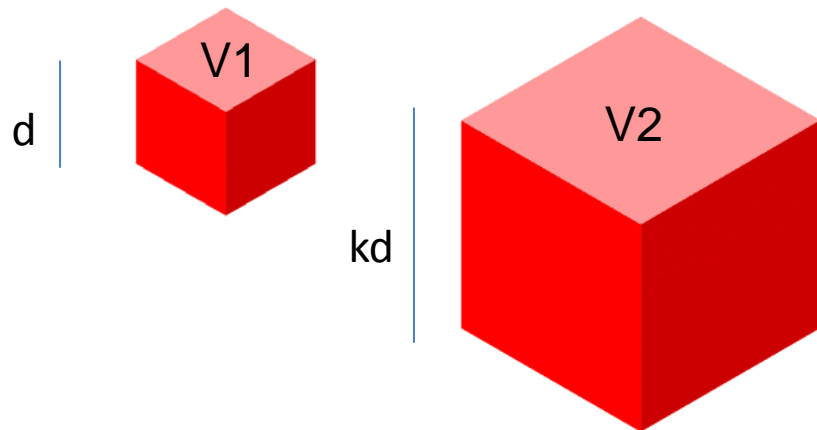
# Distortion through volumes



IN THE BARREL...
Price per bbl. of light crude, leaving Saudi Arabia on Jan. 1

April 1 $14.55

$13.34
$12.70
$12.09
$11.51
$10.95
$10.46
$2.41

73 | 74 | 75 | '76 | 1977 | 1978 | 1979

Lie factor= ~9

$V1 = d^3$
$V2 = k^3 d^3$

$V1/V2 = k^3$
$kd/d = k$



d

V1

kd

V2

Lie factor $\sim= k^3/k = k^2 =$
**size_of_effect_in_data²**

# The same data

# Distortion through areas



kd

Lie factor $\sim= k^2/k = k =$
**size_of_effect_in_data**

d

Is the bottom dollar roughly
half the size of the top one?

# The same data with lie factor = 1

Note that in a histogram you are comparing **lengths**, not **areas**



This is why it is better to use thin bars...
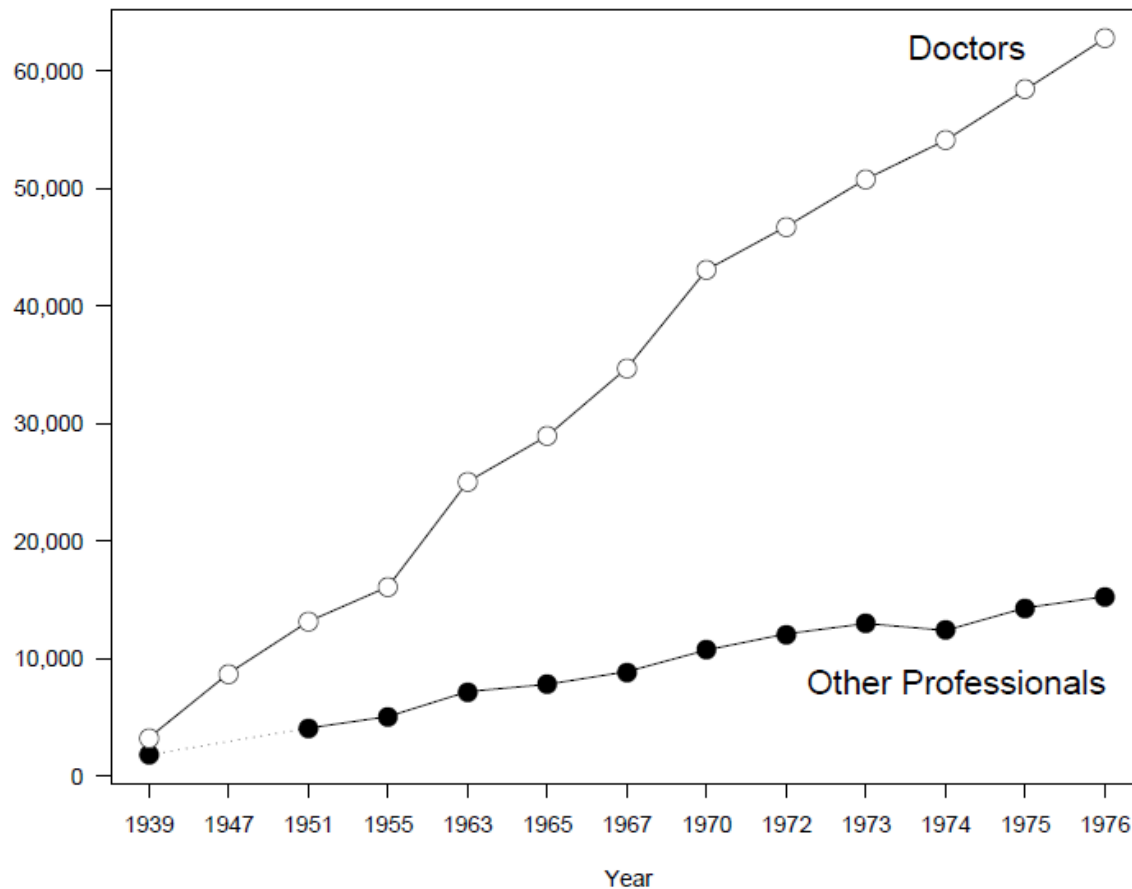
# Distortion (deliberate?)



What's wrong with this graph?
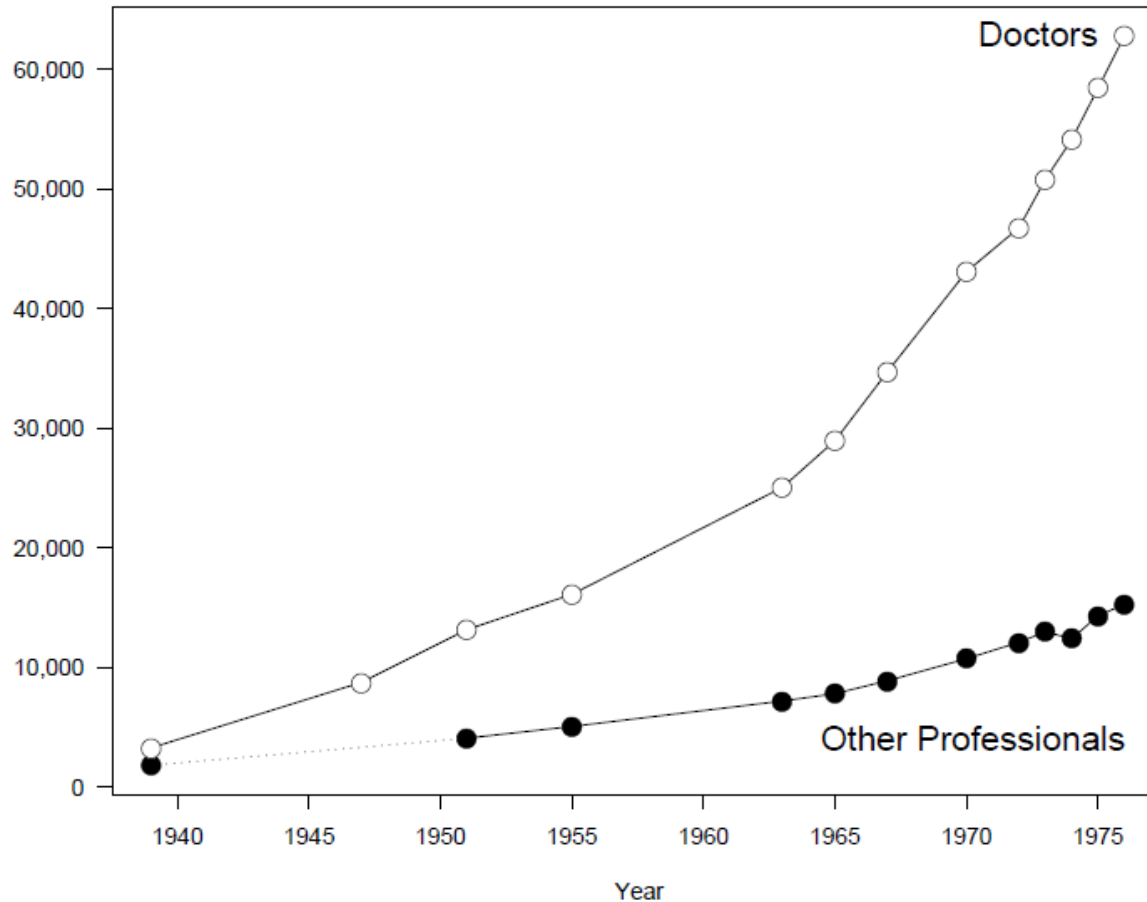
Neglecting chartjunk...

# Removing chartjunk
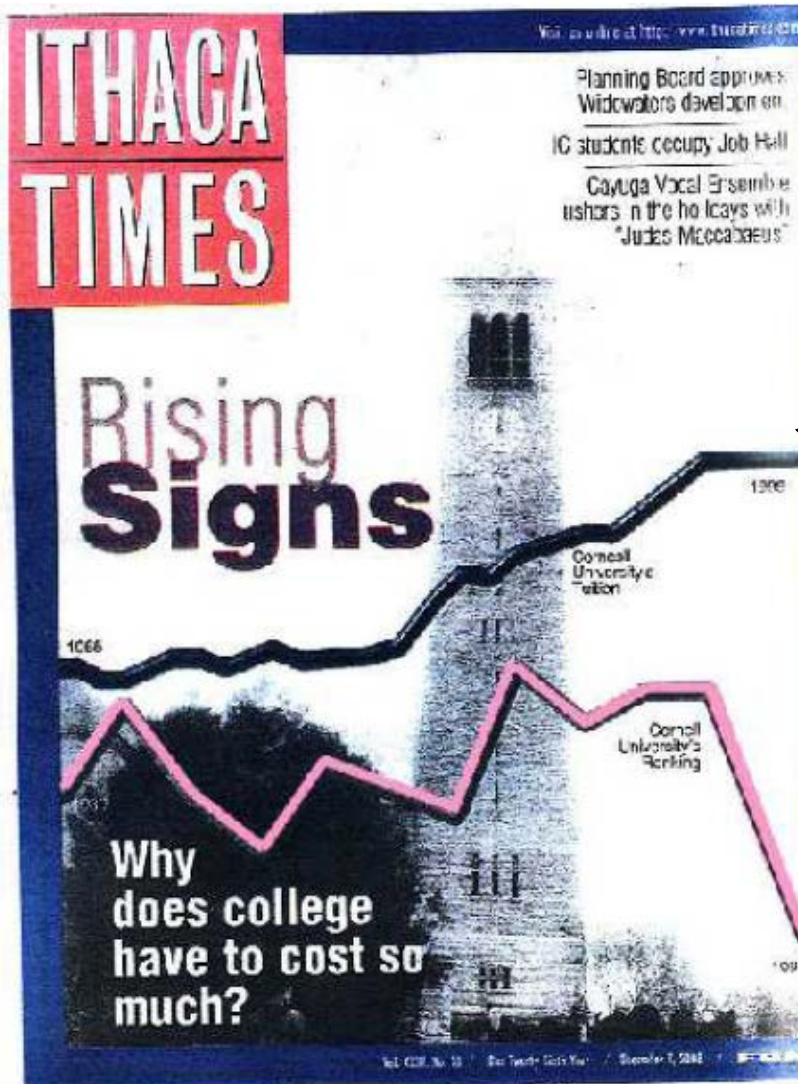


It suggests
a linear trend

# Real data...

## Median Net Incomes



The time scale was not consant!

Exponential trend !
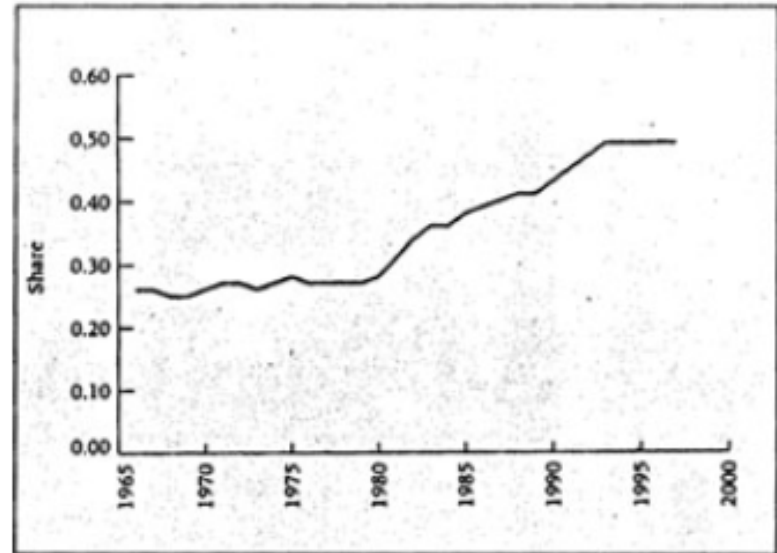
# One of the best graph lie...



- The cover story, "Why does college have to cost so much?" shows a large graph superimposed on a scene from the Cornell campus. There are two jagged lines running across the graph

    - "Cornell's Tuition" = MONEY

    - "Cornell's Ranking" = QUALITY

- The clear impression is that students are paying more for far less
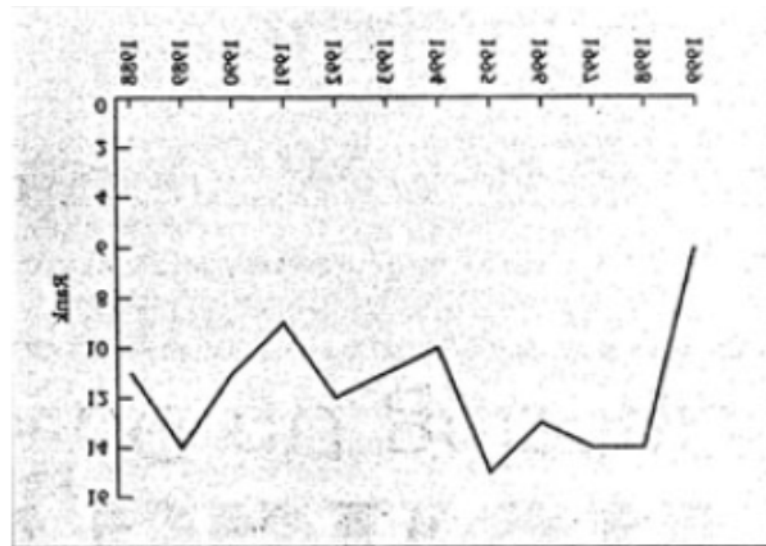
- What is wrong with it?

# The lie

- The ranking graph covers an 11 year period, the tuition graph 35 years, yet they are shown simultaneously (the same apparent width) on the same horizontal "scale".

- The vertical scale for tuition and ranking could not possibly have common units, but the ranking graph is placed under the tuition graph creating the impression that cost exceeds quality.

- And here is the masterstroke: the sharp "drop" in the ranking graph over the past few years actually represents the fact that Cornell's rank has IMPROVED from 15th TO 6th ...
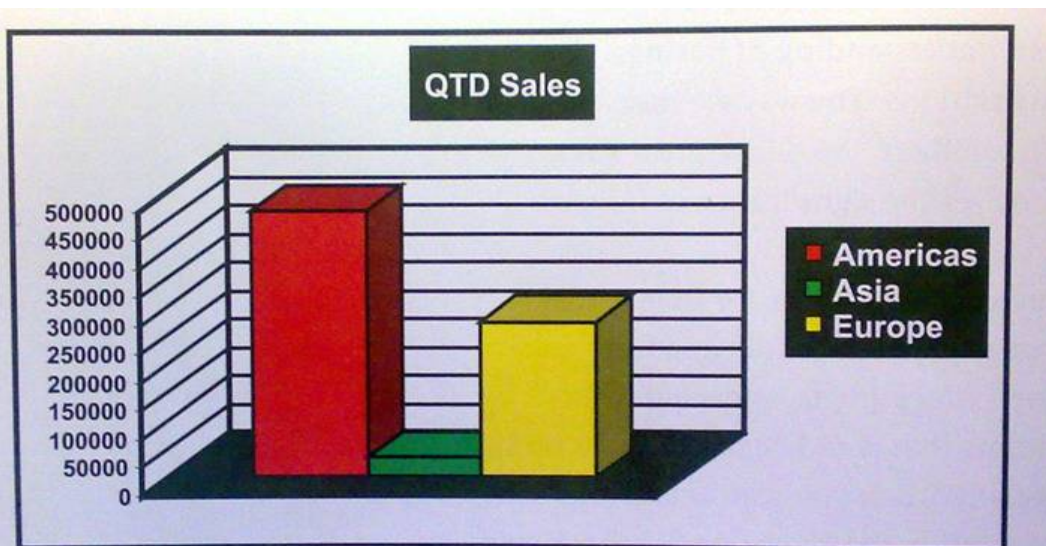
# The real data



Money



Rank

# Outline
## (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues

# Another bad example

- You are a manager of a big company
- You need to control and to report, every Monday, the current state of quarterly sales in the Americas, Asia, and Europe, with the goal of verifying your forecast
- Someone presents you with this graph
- Are you happy with it? (disregarding chartjunk)

•YOU MISS :

•Units !
•The actual date !
•Some additional summarizing information (e.g., percentages)
•Planned sales v.s. actual sales



QTD Sales

500000
450000
400000
350000
300000
250000
200000
150000
100000
50000
0

■ Americas
■ Asia
■ Europe

# All the needed information
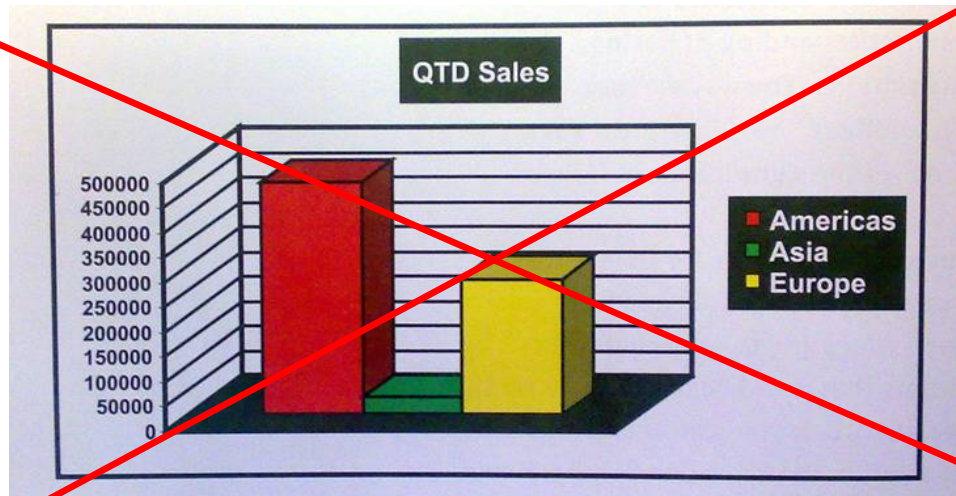
**2003 Q1-to-Date Regional Sales**

March 15, 2003

| | Sales (U.S. $) | Percent of Total Sales | Current Percent of Qtr Plan | Projected Sales (U.S. $) | Qtr End Projected Percent of Qtr Plan |
|---|---|---|---|---|---|
| Americas | 469,384 | 60% | 85% | 586,730 | 107% |
| Europe | 273,854 | 35% | 91% | 353,272 | 118% |
| Asia | 34,847 | 5% | 50% | 43,210 | 62% |
| | $778,085 | 100% | 85% | $983,212 | 107% |

Note: To date, 83% of the quarter has elapsed.

# Always remember  to
# (in addition to rules 0..3):

- Label your axes
- Make your units clear
- Use appropriate and readable label values
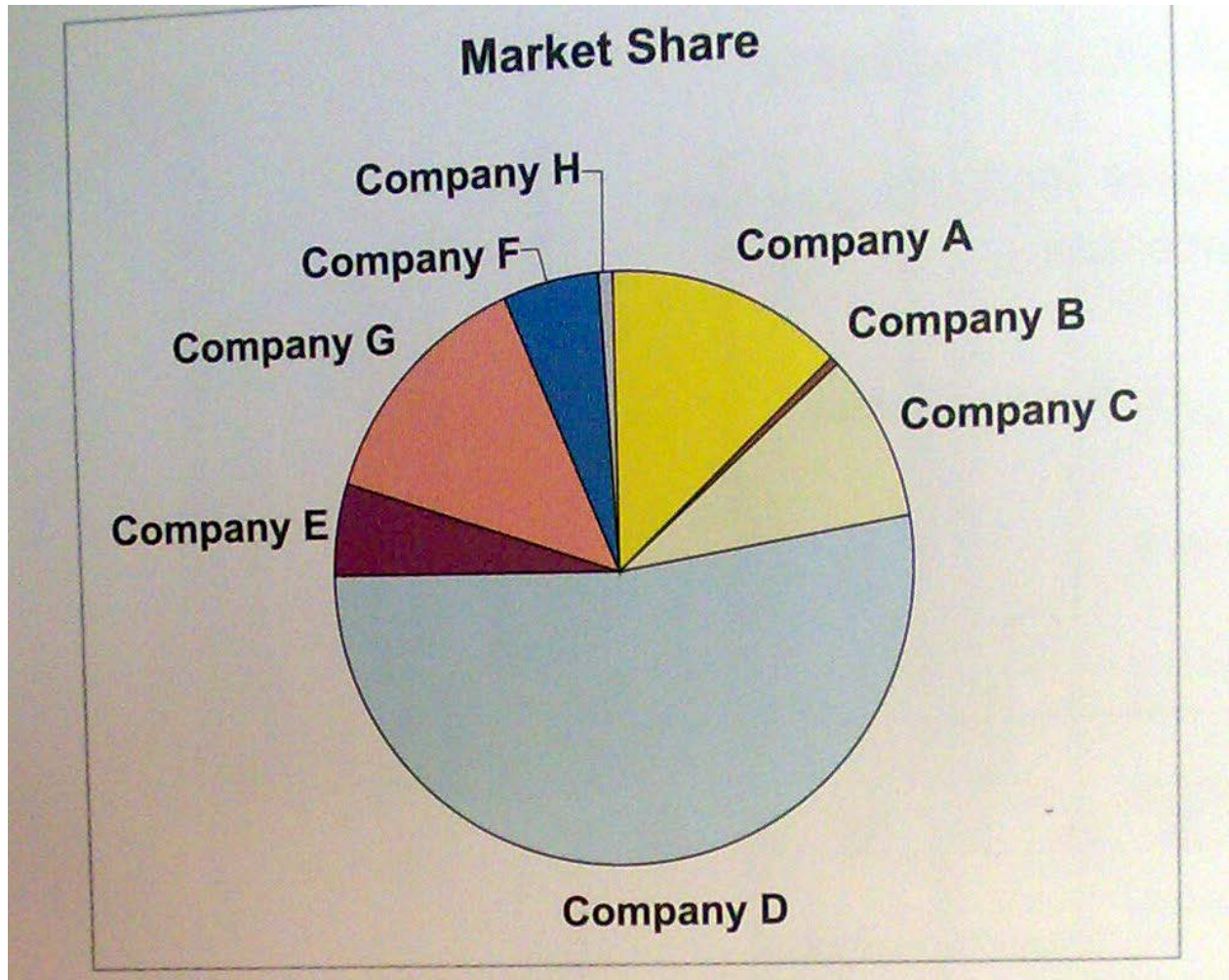- Add useful ancillary pieces of information

# The last example: our company against the world!



**Market Share**

Company H
Company F
Company G
Company E
Company A
Company B
Company C
Company D

- What is the purpose of this chart?
- Comparison !
- What is wrong whit it?

# The last example



Market Share

Company H
Company F
Company G
Company E
Company A
Company B
Company C
Company D

- Is the order clear?
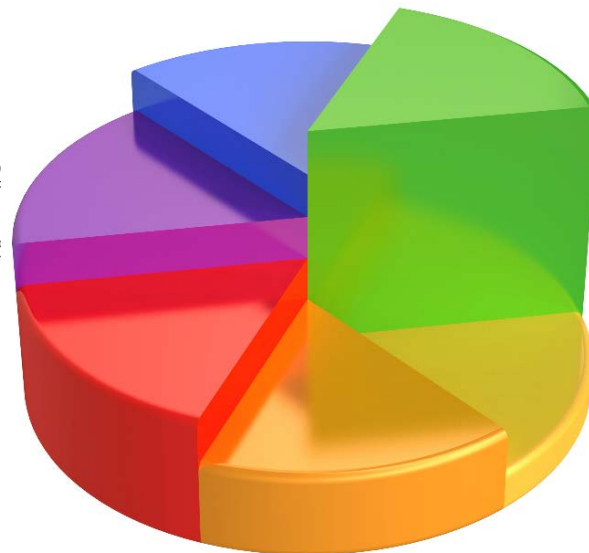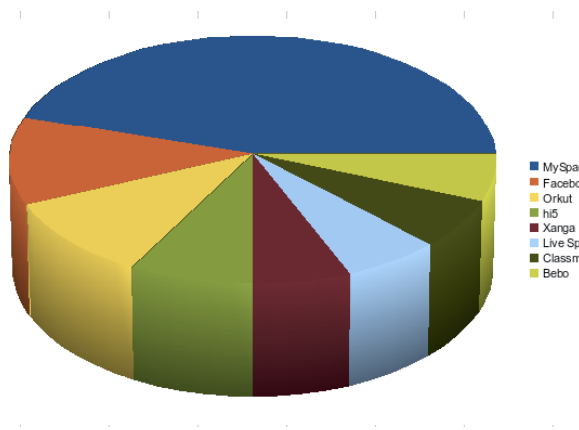- Which is my company?
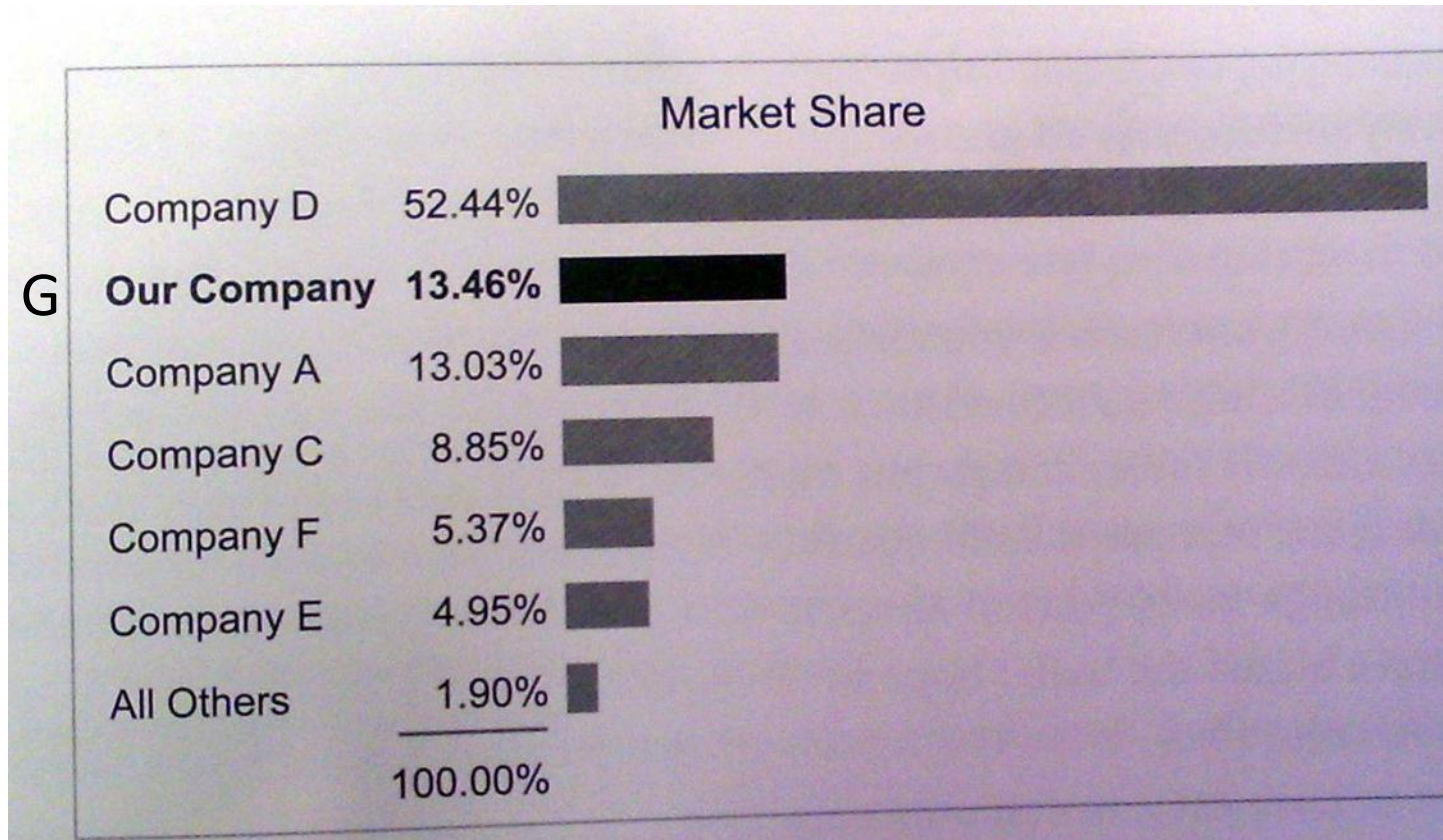- Who is bigger G or A?

# Even worst : 3D!!!

I HATE pie charts!

# A better solution



If you have ordering (ranking) alternatives think about that!

# Chartjunk is not the unique enemy...

- Before PCs, building graphs was a matter of paper and pencil
  - requiring time and effort
  - pushing you to better understand :
    - the meaning of numbers
    - the graph purpose
    - the graph organization
    - …

- now, with Excel you can produce graphs so fast that you might loose control...
  - you select predefined solutions
  - you might not understand how the graph is built (row, columns, headings, ...)
  - you can make mistakes (e.g., missing a row...)

# So...

1. Look at the numbers (plus statistics) and at the task

2. Plan a graph (even on the paper!)
   – kind of graph(s) / or even plain numbers
   – label your axes
   – units
   – scale

3. Look for an Excel implementation of your design

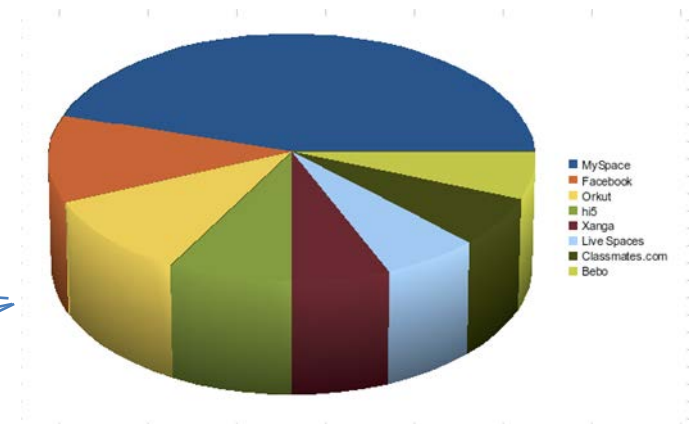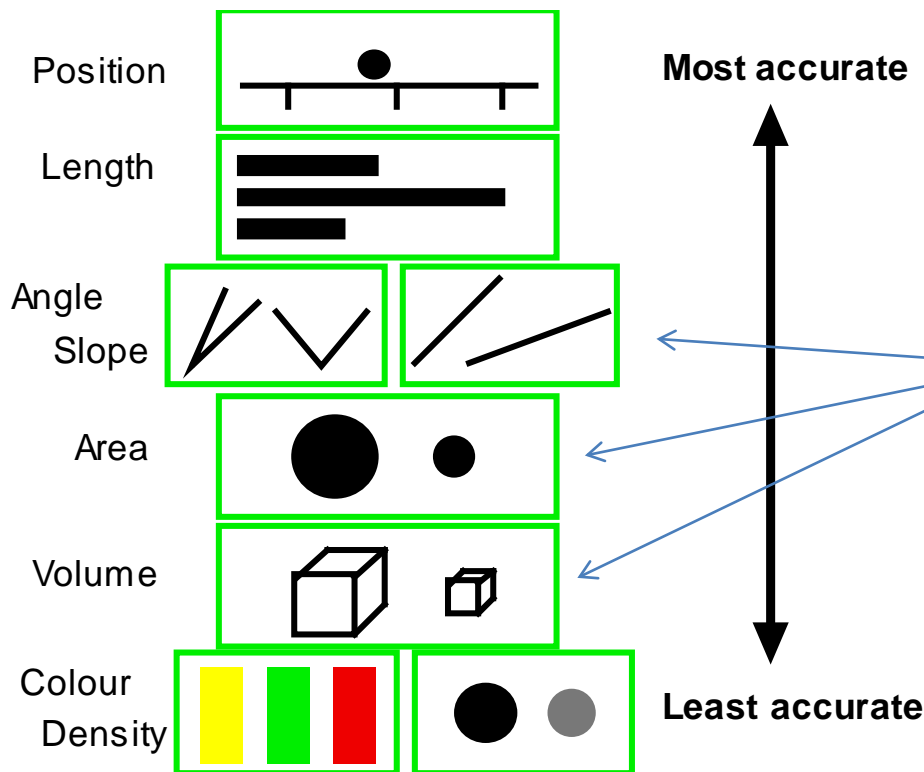4. If step 3 fails, proceed without Excel ! You can also consider more serious visualization tools, e.g., R (http://cran.r-project.org/bin/windows/base/).

# Outline
# (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues
  - Quantitative perception (basic rules)
  - The role of interaction
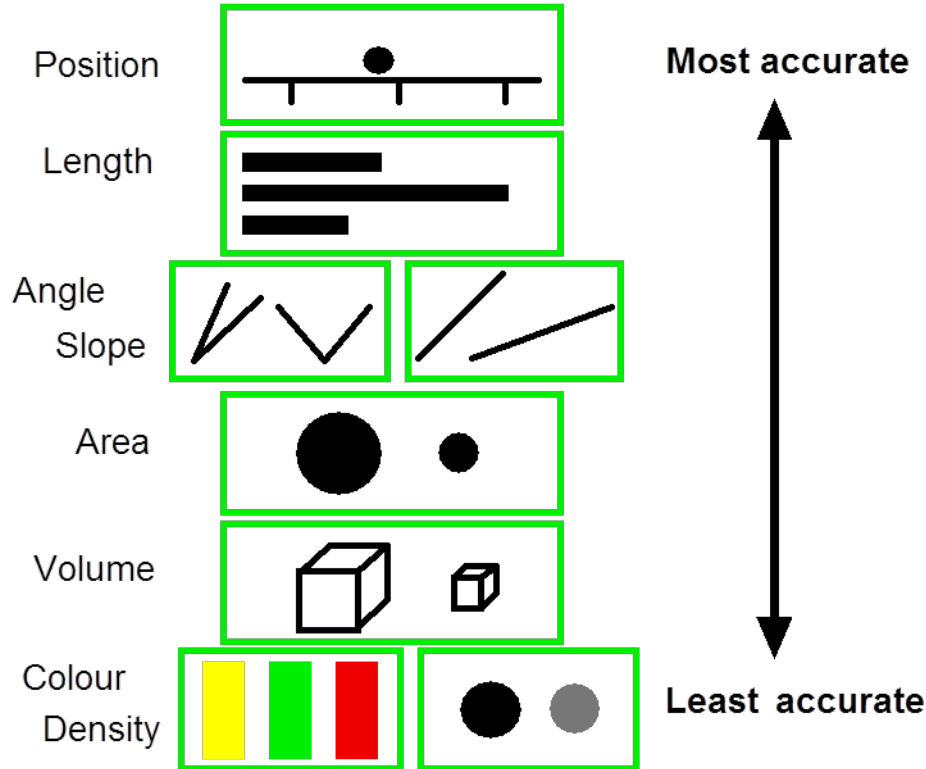- Two examples for IR

# Why do I  pie-charts?

The relative difficulty of assessing **quantitative** value as a function of visual encoding mechanism, as established by Cleveland and McGill

Position

Length

Angle
Slope

Area

Volume

Colour
Density

**Most accurate**

**Least accurate**



- MySpace
- Facebook
- Orkut
- hi5
- Xanga
- Live Spaces
- Classmates.com
- Bebo

Pie-charts discards the two first choices

I do NOT see ANY reason to use them

# What about quantitative comparison?
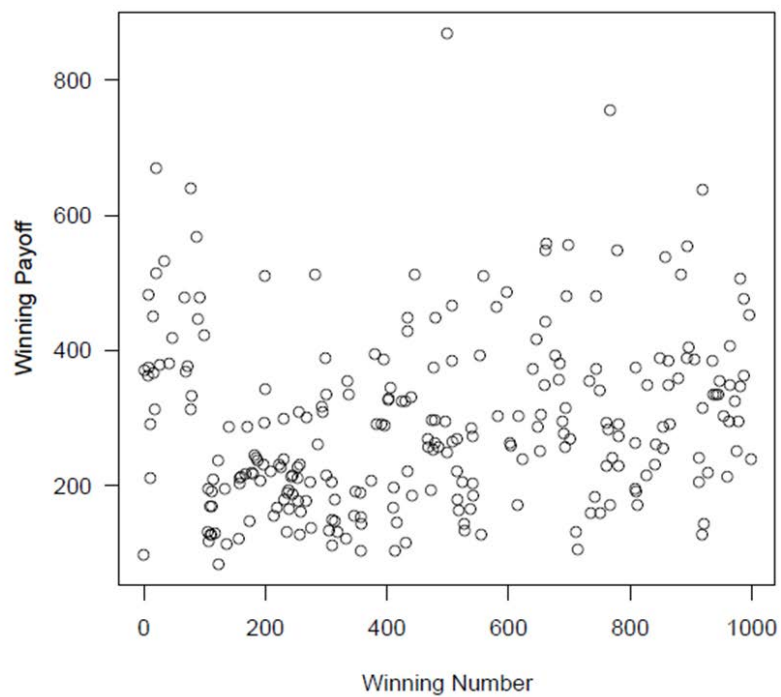


Use position and length
Avoid angles
Avoid areas
Avoid volumes
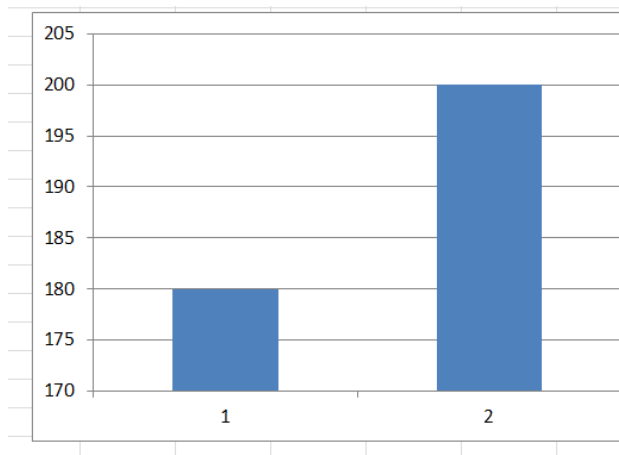
Use colors carefully

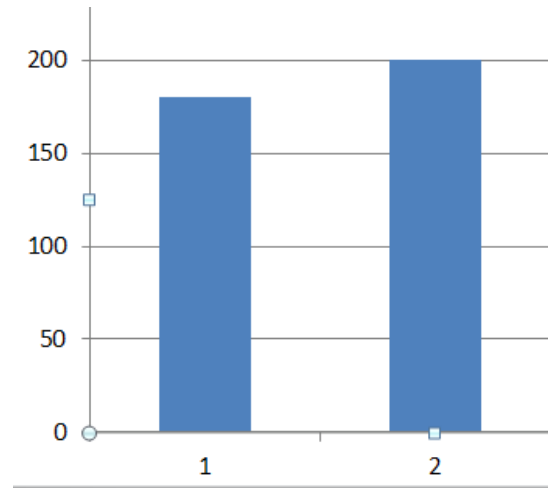# Position

- It works fine

# Length?

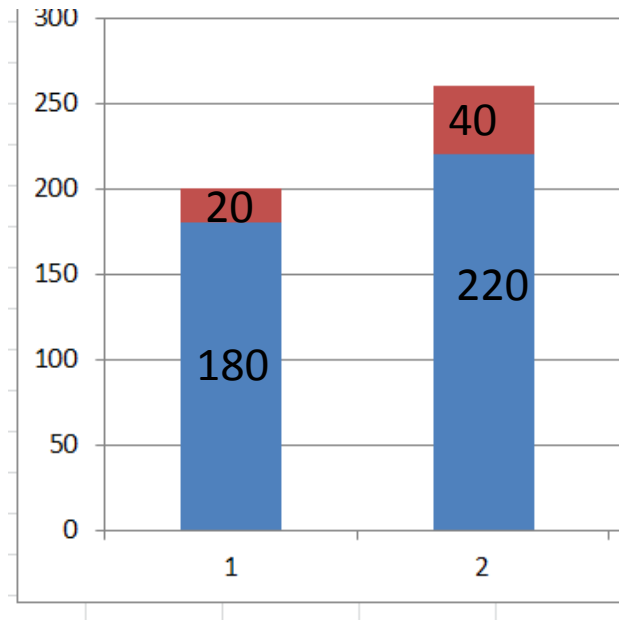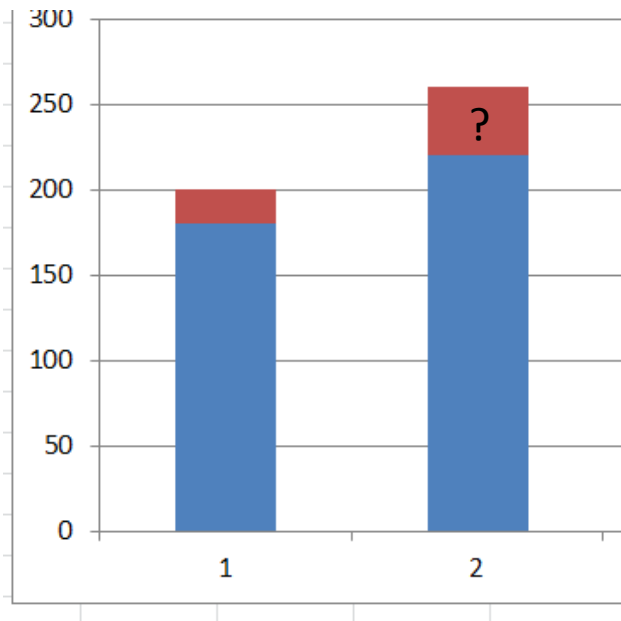- Length is fine as well , but use the right scale!



Automatically produced
by Excel



The reality

# Length?

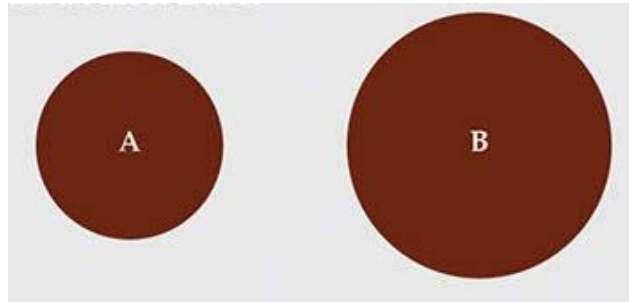- The lookup of precise number might be difficult if the position is not evident (e.g., stacked bar chart)



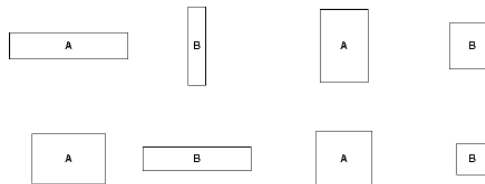It makes sense to explicitly add figures

# Areas: some new surprising issues

- Human being are very bad in estimating area ratios



- What is the ratio between this two circles?
  35% 40% 45% 50% 55% 60%   ?
- What is the shape that produces the biggest error?



- The square!
- Perceptual Guidelines for Creating Rectangular Treemaps (Nicholas Kong et al., Infovis 2010)

# Colors / Numerical data

- Someone already thought how to associate quantitative values to colors and different choices are available
- Do not reinvent the wheel
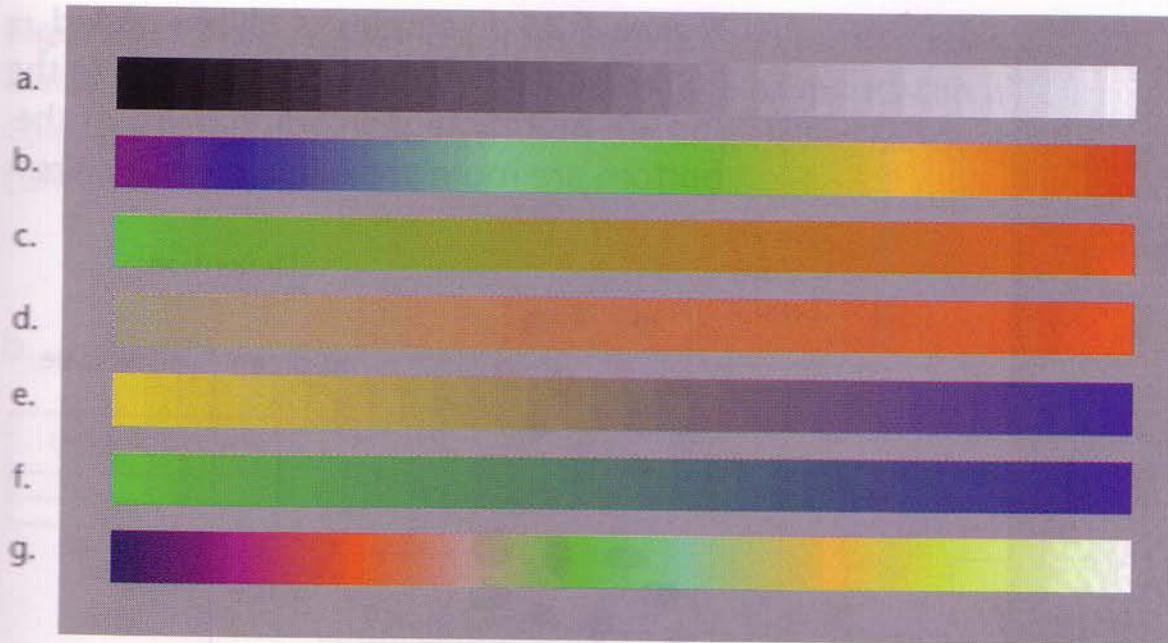- (The rainbow scale does not work)



rainbow scale



HSI color model
(Keim and Kriegel) - Issues in visualizing
large databases. Proc. of the IFIP working conference
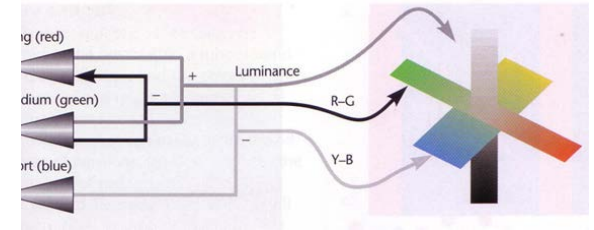on Visual database Systems, 1995
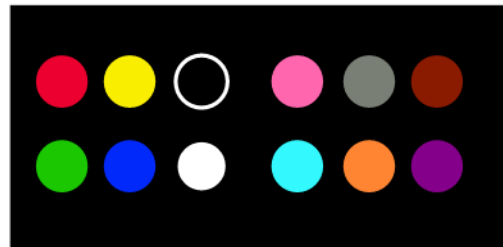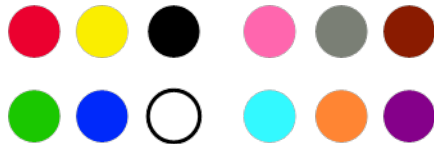
# Other choices (Colin Ware)



**Figure 4.24**
Seven different color sequences: (a) Gray scale. (b) Spectrum approximation. (c) Red-green. (d) Saturation. (e) and (f) Two sequences that will be perceived by people suffering from the most common forms of color blindness. (g) A sequence of colors in which each color is lighter than the previous one.

# Colors /Categorical data

- Colors are fine with categorical data
- Do not reinvent the wheel (again)
- The Ewald Hering idea is that there are only 6 elementary colors arranged in three pairs
- That gives us up to 12 (6+6) colors easily distinguishable (11!)

12 Colors for labeling

# Outline
## (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues
  - Quantitative perception (basic rules)
  - The role of interaction
- Two examples for IR

# Interaction?



Ad–Hoc TEL Monolingual English Task Top 5 Participants – Comparison to Median Average Precision by Topic (Topics 701–AH to 725–AH)

- Average precision (y axis) compared to the topic median (5 experiments)

# Interaction ?



Zoom in/out
Reordering,
Brushing,
...

AH-719
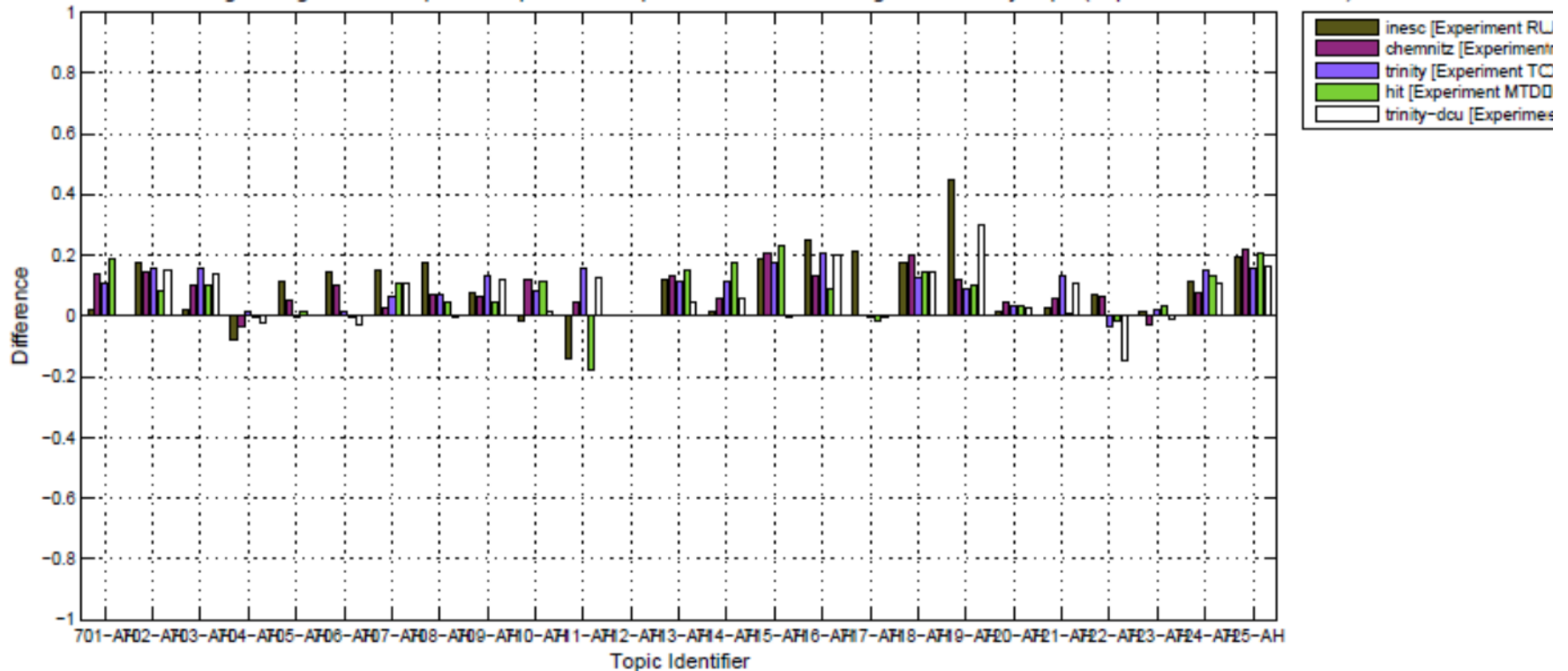AP=0,67
Median=0,25

AH-720
AP=0,92
Median=0,91

# Outline
## (basically what you have NOT to do)

- An introductive example
- Good and bad graphs
  - Basic rules
  - Some additional considerations
- Visual issues
  - Quantitative perception (basic rules)
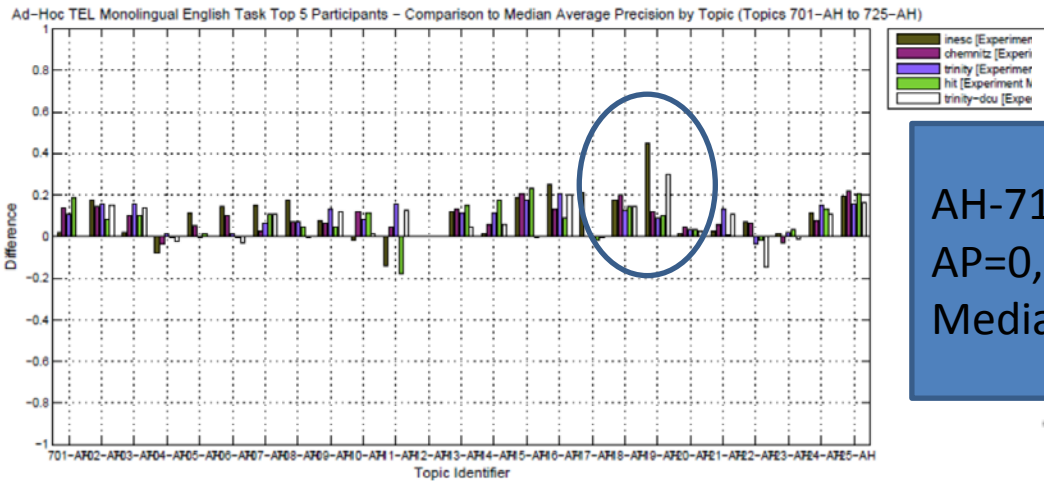  - The role of interaction
- Two examples for IR

# Parallel coordinated views



**Table**
- Topic
- Recall
- AvgPrecision

**Scatterplot**
- Y  Recall
- X  AvgPrecision

**Histogram**
- X  Number of relevant docs
  - bin = 15
- Y Number of topics in the interval

# Rank analysis

# Rank analysis (relevance 0-3)

## The actual result
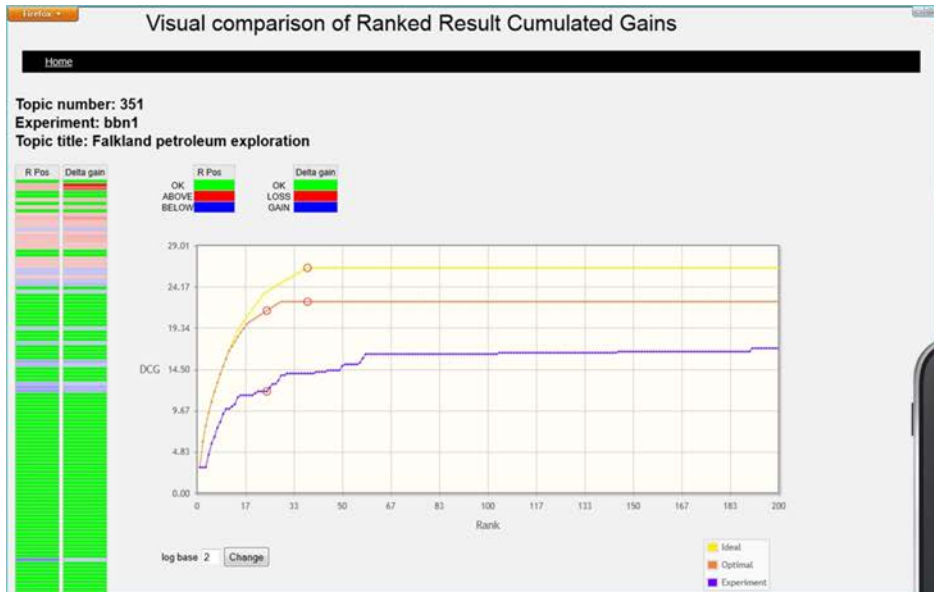
| GT(V) | DF | DCG |
|---|---|---|
| 3 | 3,00 | 3,00 |
| 1 | 1,00 | 4,00 |
| 2 | 1,26 | 5,26 |
| 3 | 1,50 | 6,76 |
| 2 | 0,86 | 7,62 |
| 2 | 0,77 | 8,40 |
| 3 | 1,07 | 9,47 |
| 2 | 0,67 | 10,13 |
| 0 | 0,00 | 10,13 |
| 1 | 0,30 | 10,43 |
| 0 | 0,00 | 10,43 |
| 3 | 0,84 | 11,27 |

| | |
|---|
| OK |
| ABOVE |
| BELOW |

## The optimal result

| GT(O) | DF | DCG |
|---|---|---|
| 3 | 3,00 | 3,00 |
| 3 | 3,00 | 6,00 |
| 3 | 1,89 | 7,89 |
| 3 | 1,50 | 9,39 |
| 2 | 0,86 | 10,25 |
| 2 | 0,77 | 11,03 |
| 2 | 0,71 | 11,74 |
| 2 | 0,67 | 12,41 |
| 1 | 0,32 | 12,72 |
| 1 | 0,30 | 13,02 |
| 0 | 0,00 | 13,02 |
| 0 | 0,00 | 13,02 |

Winter School 2012
Zinal, Valais - Switzerland
23 - 27 January 2012

# Books worth to read

- Stephen Few - Show me the number - Analytic press

- Stephen Few - Now You See It: Simple Visualization Techniques for Quantitative Analysis - Analytic press

- Robert Spence - Information Visualization: Design for Interaction (2nd Edition) - Addison-Wesley (ACM Press)

- Edward Tufte - The Visualization of quantitative information - Graphics Pr

- Colin Ware - Information Visualization, Third Edition: Perception for Design (Interactive Technologies) - Morgan Kaufmann