



# STATISTICS FOR DATA SCIENCE

## Data Visualization and Interpretation

---

**Prof. Uma D**

**Prof. Silviya Nancy J**

**Prof. Suganthi S**

Department of Computer Science and Engineering

# STATISTICS FOR DATA SCIENCE

---

## Data Visualization and Interpretation – Good Vs Bad Visualization

Prof. Uma D  
Prof. Silviya Nancy J  
Prof. Suganthi S

Data visualization is a great way to represent huge amounts of data in a simple and intuitive fashion.

All data visualizations have the same goal: help viewers easily grasp information to make quick inferences or decisions.

However, it is important that visualizations are not overdone and hit the sweet spot where they are catchy, informative, and easy to navigate.

**Color consistency:** One thing that's evident throughout this visualization is the consistency of the colors in the dashboard.

On the right is a single legend — Highlight Segment — that shows the legend for the color, which remains in both the bar chart at the bottom and the pie charts on the map.

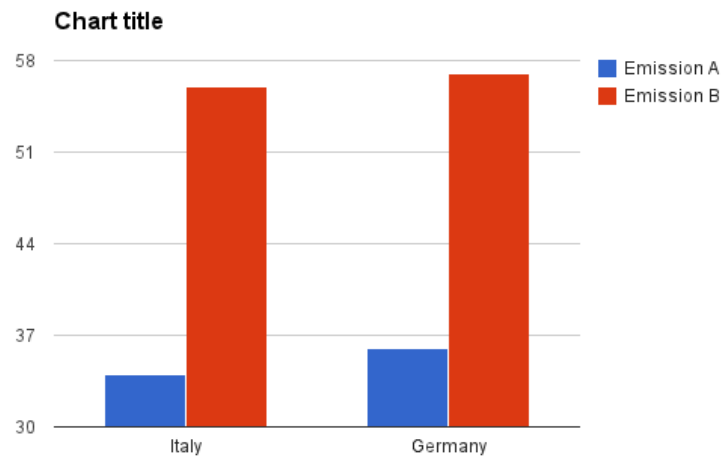
**Simplicity:** The two large charts make digesting the data easy.

**Interactivity:** The slider in the top-right corner controls the time period displayed on the charts. That interactive feature put users in control of what they view.

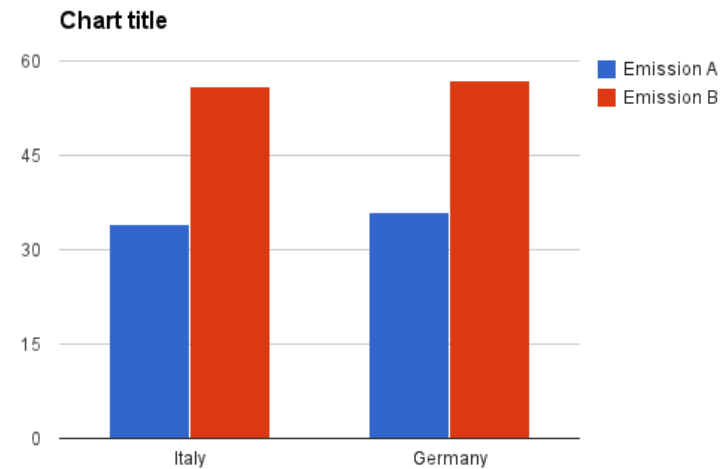
### Some Do's and Don'ts

- Use the full axis
- Avoid distortion
- Sort the data for ease of comparison
- Use consistent intervals on any axis or indicate a break
- Use the chart type wisely
- Don't use colors and effects without reason
- Don't use 3D

For bar charts, the numerical axis (often the y axis) must start at zero.



Wrong



Correct

### **Sort your data for easier comparisons**

The bar chart is a good example, where the chart x-axis is sorted on the y-values not on the alphabetic order of the country names.

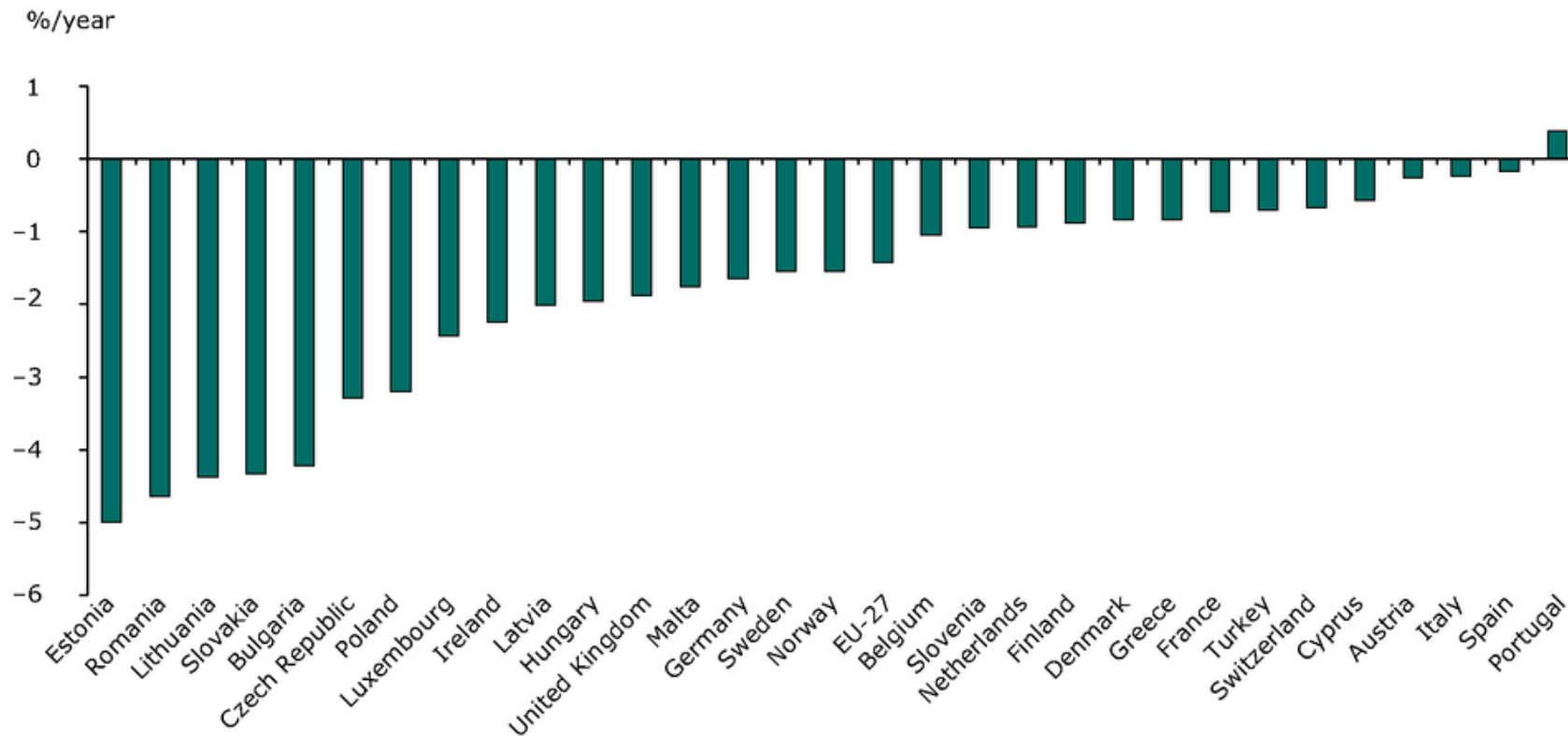
It will be otherwise very difficult if not impossible for users to do a proper comparison across the many bars.

It is in any case easy with a quick eye-scan to find your own country in the list.

# STATISTICS FOR DATA SCIENCE

## In Bar Charts???

Contd..

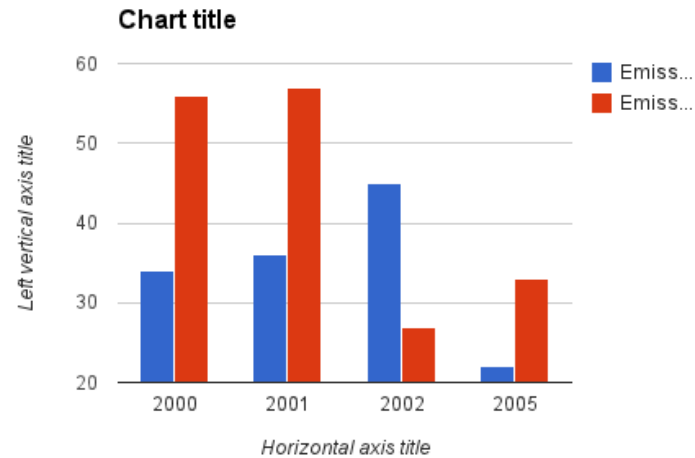




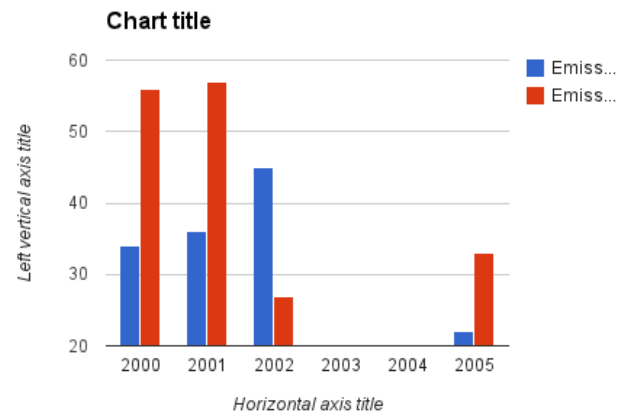
## In Bar Charts???

**Use consistent intervals on axis (be transparent on data gaps)**

- Be clear when some data is missing. Explain the reason why is missing. Use the full axis and do not skip values when you have numerical data.
- The x-axis in the "wrong example" below has a time-series with inconsistent intervals (missing years 2003 and 2004) giving a distorted view of data over time.



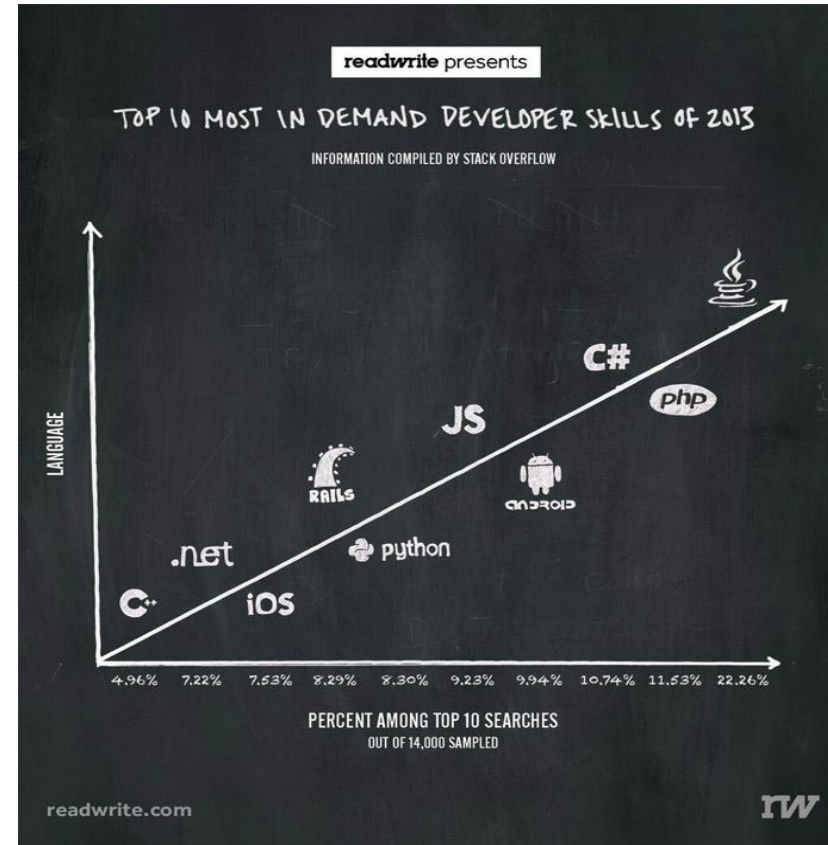
Wrong



Correct

## In Line Graph???

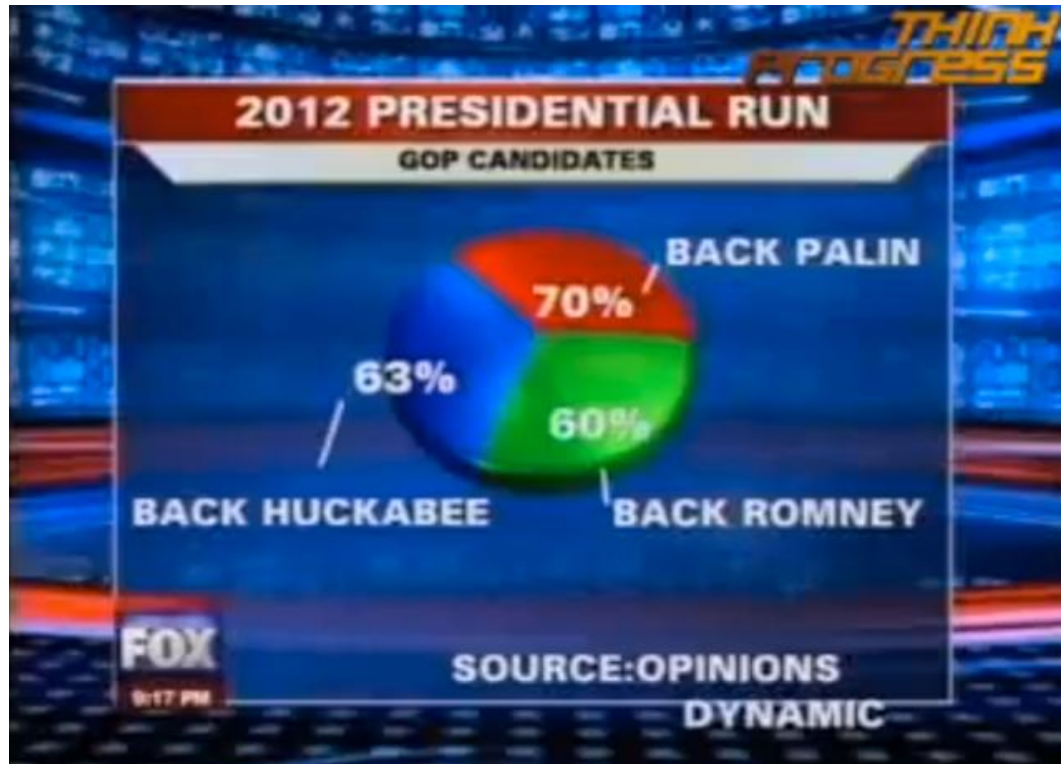
**An axis should have something on it**  
We would have thought this obvious, but a line graph should have something numerical on each axis. The graph below does not. Its vertical axis is labeled "Language" but even in the most generous interpretation this is a categorical variable and thus not appropriate for display using a line graph.



# STATISTICS FOR DATA SCIENCE

## In Pie Charts????

In this pie chart, the three sectors of the pie add up to 193%, which makes no sense. Such mistakes in data would render your final visualizations useless.



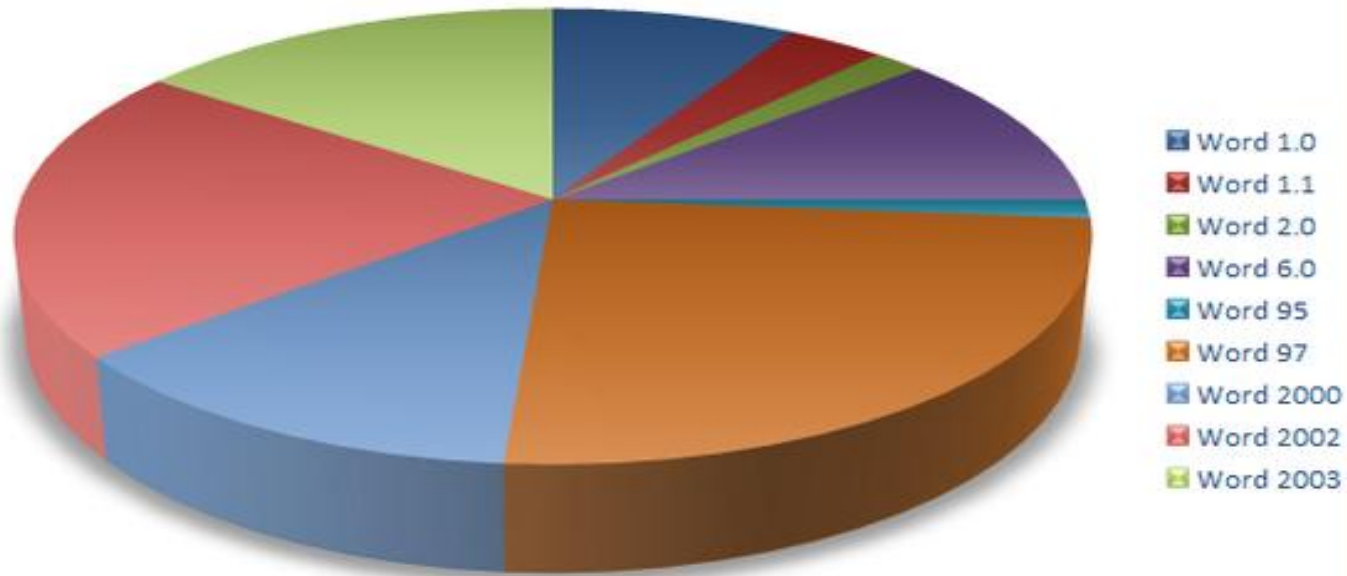
### Wrong Choice of Data Visualization

Once your data is ready, you should be careful about what type of visualization you use. For instance, in the visualization below, a pie chart was the wrong choice. The intention there was to show how many features a given Microsoft Word version has. The pie chart, on the other hand, shows the proportion of features in a particular version as a percentage of the total features in all versions. A bar chart would be a better choice for this data.

# STATISTICS FOR DATA SCIENCE

## Good vs Bad visualization

Microsoft Word Features By Version Added



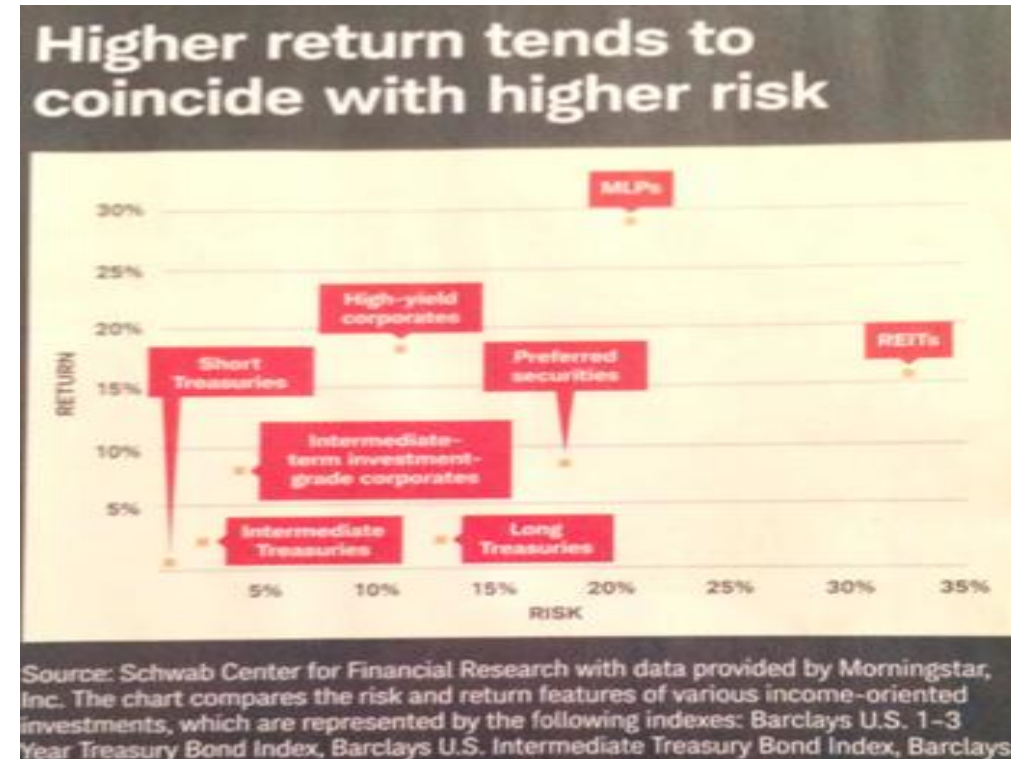
# STATISTICS FOR DATA SCIENCE

## In Scatter Plot????



When there are two variables, and their correlation is of interest, a scatterplot is usually recommended. But not here!

The text labels completely dominate this chart and the designer tried very hard to place them but a careful look reveals that some boxes are placed above the dots while others are placed to their right and the dot for "Short Treasuries" holds refuge quite a while away from the dot. This means the locations of the text boxes do not substitute for the dots.





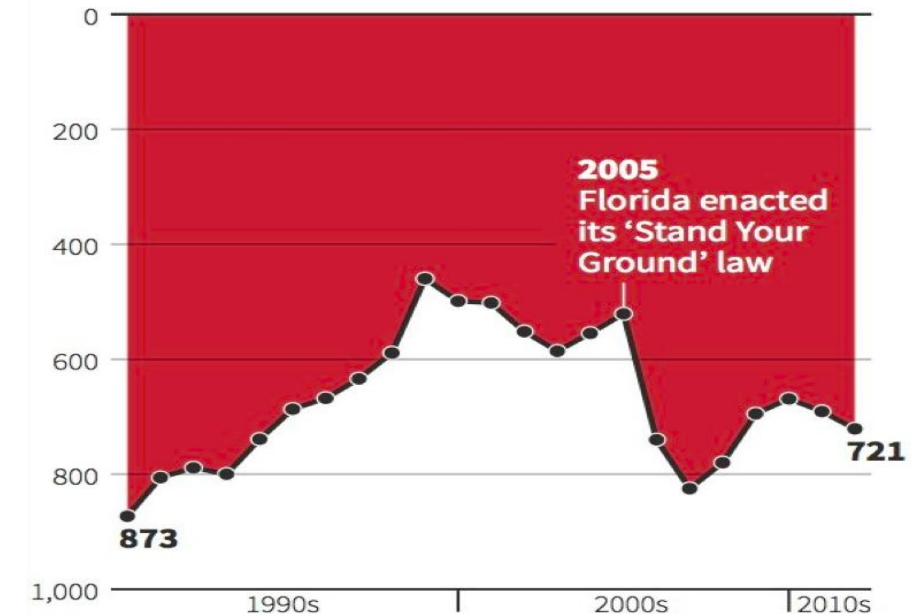
# STATISTICS FOR DATA SCIENCE

## In Scatter Plot????

At first glance, it looks like gun deaths are on the decline in Florida. But a closer look shows that the y-axis is upside-down, with zero at the top and the maximum value at the bottom. As gun deaths increase, the line slopes downward, violating a well established convention that y-values increase as we move up the page.

### Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

REUTERS

In summary, data visualizations tell stories.

Relatively subtle choices, such as the range of the axes in a bar chart or line graph, can have a big impact on the story that a figure tells.

When you look at data graphics, you want to ask yourself whether the graph has been designed to tell a story that accurately reflects the underlying data, or whether it has been designed to tell a story more closely aligned with what the designer would like you to believe.





**THANK YOU**

---

**Prof. Uma D**

**Prof. Silviya Nancy J**

**Prof. Suganthi S**

Department of Computer Science and Engineering