## Reference

- Ward, M., Grinstein, G. G., Keim, D. **Interactive data visualization foundations, techniques, and applications**. Natick, Mass., A K Peters, 2010.

- Binary distance :
http://people.revoledu.com/kardi/tutorial/Similarity/BinaryVariables.html

140

140

# Visual Representation of Data

Álvaro Figueira, PhD.
MSc in Data Science - Data Visualization

[dcc]   **DEPARTAMENTO DE CIÊNCIA DE COMPUTADORES**
**FACULDADE DE CIÊNCIAS** DA UNIVERSIDADE DO PORTO

FC FACULDADE DE CIÊNCIAS
UNIVERSIDADE DO PORTO   U. PORTO

141

141

## Data Visualization goals

Data visualization refers to the techniques used to communicate data or information by encoding it as visual objects contained in graphics.

The goal is to **communicate information clearly and efficiently to users**. It is also one of the steps in data analysis or data science:

> *The main goal of data visualization is to communicate information clearly and effectively through graphical means.*
>
> *It doesn't mean that data visualization needs to look boring to be functional or extremely sophisticated to look beautiful. To convey ideas effectively, both aesthetic form and functionality need to go hand in hand, providing insights into a rather sparse and complex data set by communicating its key-aspects in a more intuitive way.*
>
> *Yet, designers often fail to achieve a balance between form and function, creating gorgeous data visualizations which fail to serve their main purpose — to communicate information.*
>
> **Friedman*, Data Visualization and Infographics,* 2008**

**Obs:** however, an ideal visualization should not only communicate clearly, but stimulate viewer engagement and attention.

143

## The current main issue

Graphical visualizations began as a mean to communicate numbers and quantities.

But, information visualizations are also executing particular analytical tasks such as **making comparisons** or **determining causality**.

The design principle of the information graphic should *support that analytical task*, showing the comparison or causality.
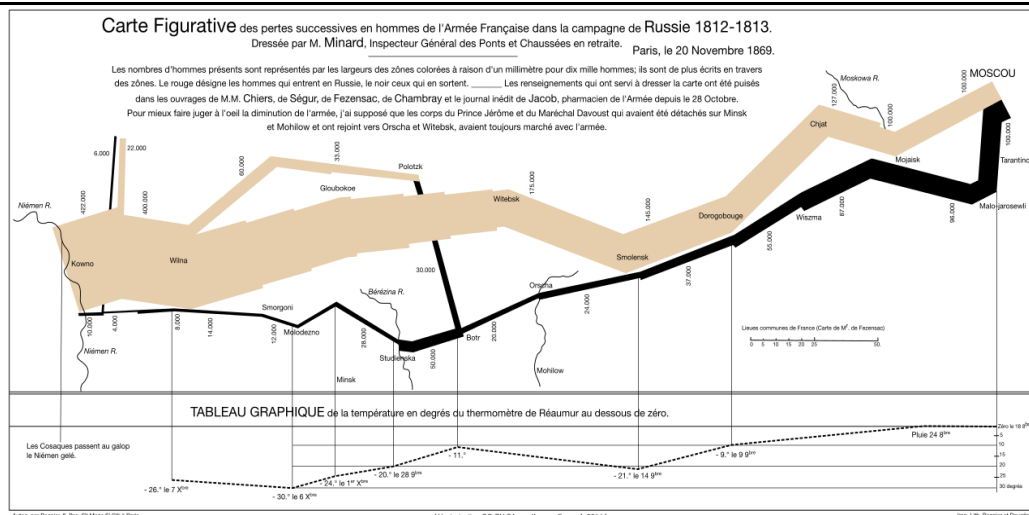
However, more than often this doesn't happen with many graphics...

144

2

# Graphical displays should

- **Show the data**
- Induce the viewer to think about the **substance** rather than about the methodology, the graphic design, or even the technology of graphic production
- Present many numbers in a **small space**
- **Avoid distorting** what the data has to "say"
- Make **large data sets coherent**
- Encourage the eye to **compare different pieces** of data
- Reveal the data at **several levels of detail**, from a broad overview to the fine structure
- Serve a **clear purpose**: description, exploration, tabulation or decoration
- Be closely **integrated with the statistical and verbal** descriptions of a data set.

145



The **Minard diagram** shows the losses suffered by Napoleon's army in the 1812–1813 period. Six variables are plotted: the size of the army, its location on a two-dimensional surface, time, direction of movement, and temperature. The line width illustrates a comparison (size of the army at points in time) while the temperature axis suggests a cause of the change in army size. This multivariate display on a two-dimensional surface tells a story that can be grasped immediately while identifying the source data to build credibility.

[Tufte wrote in 1983 that: "It may well be the best statistical graphic ever drawn."]

146

## Quantitative Information Types

- Simple Quantities
- Time-series
- Rankings
- Part-to-whole
- Deviation
- Frequency distribution
- Correlation
- Nominal comparison
- Geographic or geospatial

147

## Taxonomies for visual representation of data quantities

Christian Behrens, 2008



**Correlations**
- Scatterplot
- Bubble chart

**Continuous Quantities**
- Simple line chart
- Multiset line chart
- Stacked area chart
- Sparklines

**Discrete Quantities**
- Simple bar chart
- Multiset bar chart
- Dot matrix
- Stacked bar chart
- Isometric bar chart
- Span chart

**Proportions**
- Simple pie chart
- Ring chart

**Flows**
- Sankey diagram
- Thread arcs

**Hierarchies**
- Tree diagram
- Treemap

**Networks**
- Tree diagram
- Relation circle
- Pearl necklet

**Space**
- Topographic map
- Thematic map

148

# Some current techniques

**Abstract structures**
- Proportions: pie chart and ring chart
- Correlations: scatterplot, bubble chart
- Discrete quantities:
- bar chart, dot matrix, stacked bar charts
- Continuous quantities:
- line chart, stacked chart, sparklines
- Multidimensional: parallel coordinates

**Hierarchical structures: trees**
- Node-link layout (cartesian and polar)
- Treemaps: rectangular, circular and Voronoi
- Sunburst.

**Relational structures: networks**
- Node-link diagrams
- Layouts: matrix, linear, force directed, Sankey, circular, polar, geographical

**Temporal structures**
- Timelines: linear and polar
- Flows (Sankey diagrams)

**Spatial structures: maps**
- Dot distribution maps
- Isometric maps: isolines and heatmaps
- Choropleth maps
- Magnification and fish-eye views
- Cartograms: Dorling's, area-value, isochronic

**Temporal structures**
- Animated maps
- Representation of trajectories
- Temporal flows

**Textual structures**
- Word clouds
- Textual trees

149

---

# Proportions

Pie Chart

Ring Chart



When to avoid?

150

5

## Correlations: Scatter plot

| | S1 | S2 |
|---|---|---|
| Rui | 43 | 20 |
| Pedro | 25 | 25 |
| Ana | 34 | 15 |
| Jorge | 10 | 15 |
| Rita | 30 | 26 |
| João | 35 | 25 |



Scatter Plot

What does the x-axis represent?

151

## Correlations: Bubble chart



How many properties must have each record in this chart?

152

## Box Plot



Can you explain a box plot?

153

# Waterfall graphs



A waterfall chart helps in understanding the cumulative effect on the initial value which is increased or decreased by a series of intermediate values leading to a final value

154

154

# Continuous quantities and Sparklines

Tufte's (Tufte, 2006) own Sparklines are **data-intense**, **design-simple**, **word-sized graphics**. Sparklines, can **display the temporal evolution of variables**, its most recent value, its name and the out-of-the-norm values, **everything in a highly condensed graphic**.

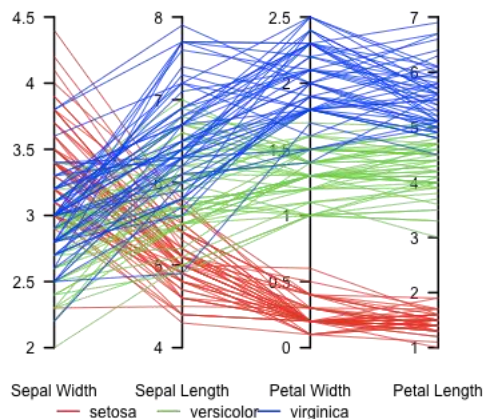**Obs:** Sparklines can run in a text layout.



| | Region ⇕ | magnitude_trend ⇕ | count ▾ | avg(Magnitude) ⇕ |
|---|---|---|---|---|
| 1 | Fox Islands, Aleutian Islands, Alaska | | 14 | 3.271429 |
| 2 | Island of Hawaii, Hawaii | | 14 | 3.035714 |
| 3 | Puerto Rico region | | 14 | 3.035714 |
| 4 | Southern Alaska | | 10 | 2.880000 |
| 5 | Andreanof Islands, Aleutian Islands, Alaska | | 8 | 2.712500 |
| 6 | Central California | | 8 | 2.925000 |
| 7 | Baja California, Mexico | | 7 | 2.957143 |
| 8 | Virgin Islands region | | 7 | 3.185714 |
| 9 | Kodiak Island region, Alaska | | 6 | 2.733333 |
| 10 | Central Alaska | | 5 | 2.920000 |

Note: the aspect ratio of line chart or a sparkline are crucial for good reading. The visual average of the hill-slopes within the line should be ideally 45 degrees.

155

# Multidimensions and Parallel Coordinates
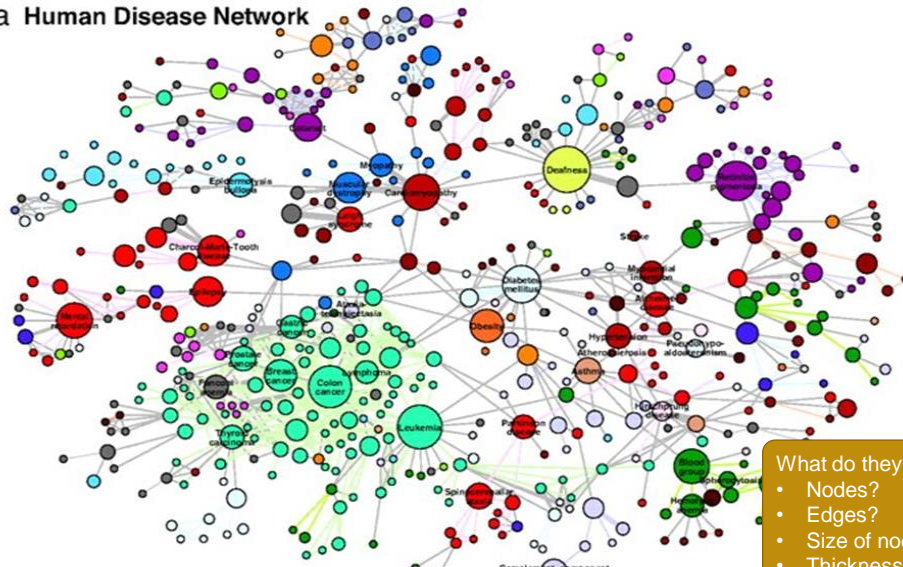
**Parallel coordinate plot, Fisher's Iris data**



Parallel coordinates can be arranged in order to show certain correlations among attributes for the same type of object

Sepal Width    Sepal Length    Petal Width    Petal Length
— setosa    — versicolor — virginica

156

8

# Hierarchical structures: Trees



157

# Relational structures with linear layout



158

## Relational structures: networks

a **Human Disease Network**

What do they represent:
- Nodes?
- Edges?
- Size of nodes?
- Thickness of edges?

source: Goh et al. The human disease network

159

## Relational structures: Sankey diagrams

Estimated U.S. Energy Use in 2011: ~97.3 Quads

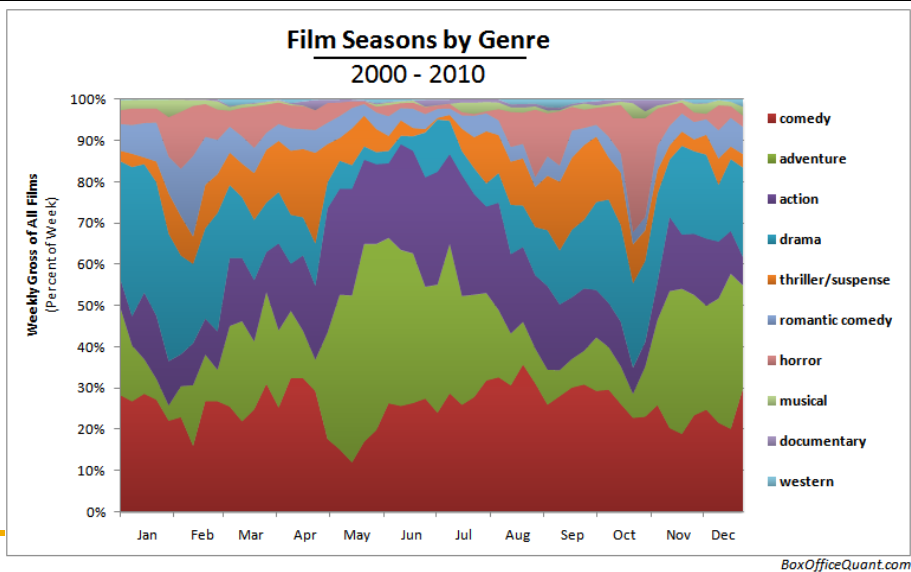**Lawrence Livermore National Laboratory**

160

## Relational structures: Sankey diagrams
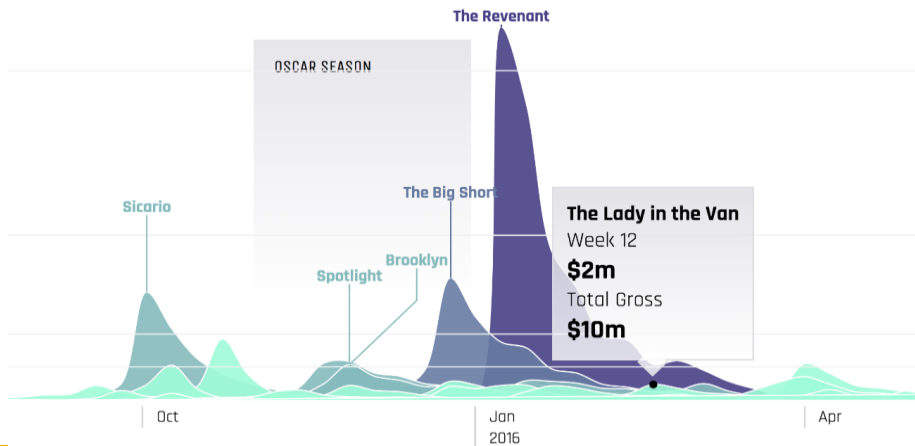


161

## Time series: stacked graphs



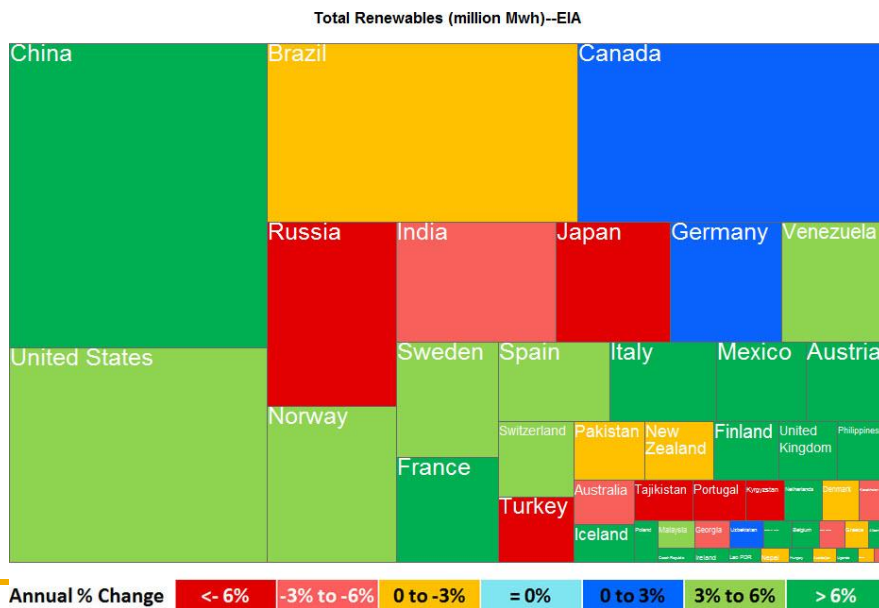162

## Filled regions

This is just a teaser for you all...
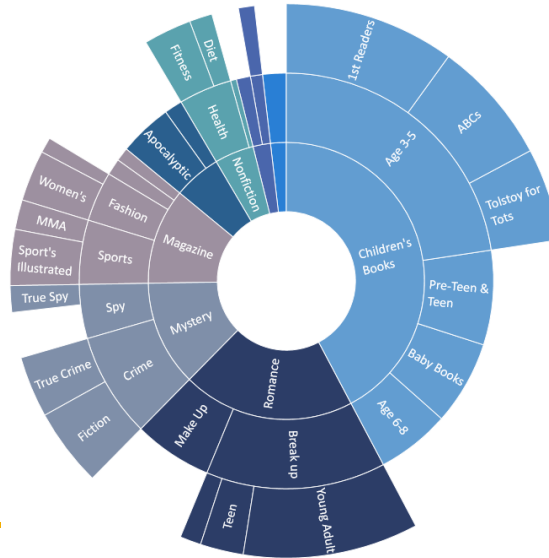
163

## Hierarchies: Treemaps



164

## Hierarchies: circular treemaps

165



## Heatmaps

### Value (%)

| | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov | Dec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | -9.8% | 84.5% | -79.5% | 1.8% | 23.4% | 54.7% | 48.2% | 75.4% | -61.4% | -76.7% | -35.8% | |
| 2010 | -66.5% | 41.2% | 38.9% | -41.1% | 20.5% | 55.7% | -97.9% | 77.8% | 80.3% | -88.7% | -66.0% | 10.6% |
| | 87.6% | -12.5% | 86.7% | -21.6% | 56.8% | 2.0% | 57.5% | 10.1% | -86.2% | 8.6% | 66.4% | 82.2% |
| | 95.8% | 42.3% | 64.8% | 0.1% | 50.8% | -70.0% | -6.7% | 32.9% | 47.0% | -16.9% | 40.7% | -13.3% |
| | -35.0% | -8.2% | 36.9% | 38.4% | -3.2% | -76.1% | 61.2% | 9.8% | 10.4% | 5.2% | 93.2% | 92.8% |
| | -49.2% | 24.5% | -46.3% | 54.1% | -18.9% | -38.9% | 3.1% | -59.4% | 15.2% | 67.5% | -73.7% | 20.8% |
| 2005 | 52.1% | -87.2% | 76.2% | -26.4% | -42.3% | 65.3% | -9.0% | 88.4% | 60.7% | 34.8% | -59.2% | 37.0% |
| | 89.1% | -17.7% | 5.6% | 81.4% | 11.1% | 81.7% | -90.3% | -58.4% | 42.7% | -69.4% | -0.3% | 18.8% |
| | -61.5% | 42.3% | -70.0% | 2.9% | 17.2% | 85.1% | -20.1% | -43.7% | 52.5% | 70.1% | -21.0% | 55.4% |
| | 55.0% | 4.7% | -68.1% | 76.2% | -29.1% | 33.8% | -95.7% | 41.9% | 69.8% | 5.4% | -16.2% | 43.7% |
| | 8.0% | 27.5% | 17.3% | 39.0% | 26.1% | -14.5% | -56.4% | 66.5% | 66.5% | 67.2% | -17.6% | 63.1% |
| 2000 | 43.4% | 53.5% | -34.6% | 49.9% | -56.9% | -49.9% | 41.2% | 77.7% | 35.5% | -50.9% | -10.5% | 72.6% |
| | -34.9% | -19.1% | -73.3% | 44.7% | 0.2% | 70.5% | 85.0% | -54.2% | -40.0% | 13.4% | 27.3% | -20.6% |
| | 83.8% | -12.9% | -48.2% | 65.4% | 60.8% | 50.0% | -31.0% | -39.7% | -5.2% | 59.3% | 46.7% | 87.0% |
| | -33.0% | -37.6% | -86.2% | 10.7% | -1.4% | 91.2% | -49.8% | -31.2% | -91.5% | -38.6% | -49.8% | 18.5% |
| | 37.2% | 43.2% | -45.4% | -32.2% | 45.3% | 25.2% | 90.9% | -69.2% | 42.3% | -40.8% | 96.7% | 10.9% |
| 1995 | -33.2% | -47.2% | 98.1% | 1.8% | -26.5% | 72.6% | -21.8% | 90.5% | 20.6% | 76.0% | 50.6% | -42.0% |

What's the meaning of the colors?
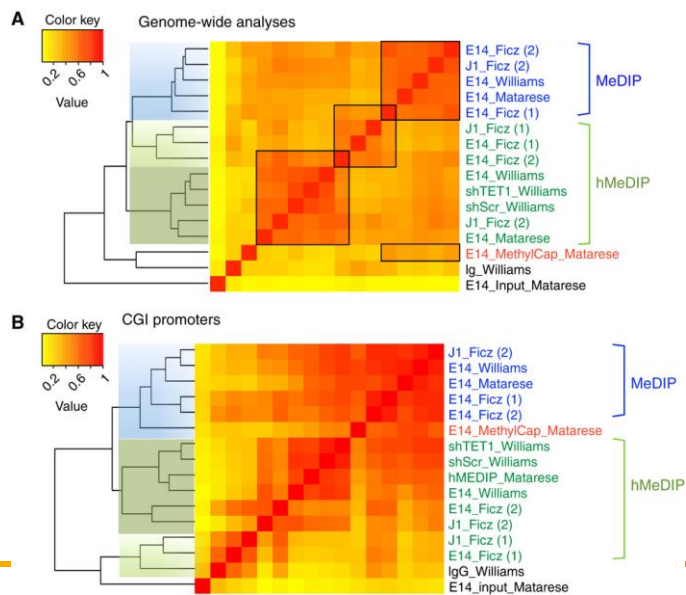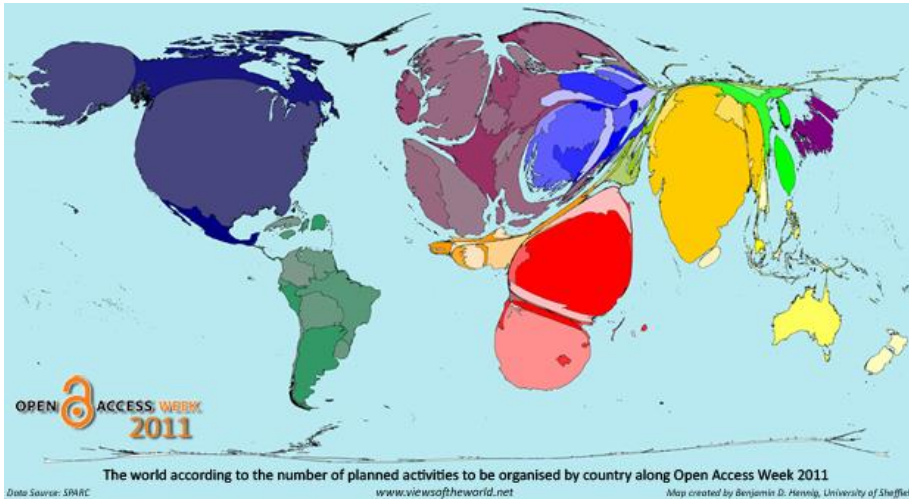
166

13

## Heatmaps



167

# Heatmaps with trees



168

**Maps and area cartograms**

169



**Maps and area cartograms**

170

15

## Textual structures: Word Clouds



What can you infer from this?

171

## Textual Structures



172

## Algorithmic Flow Charts



173

Recently, we have seen incredible, big data, real-time, dashboards
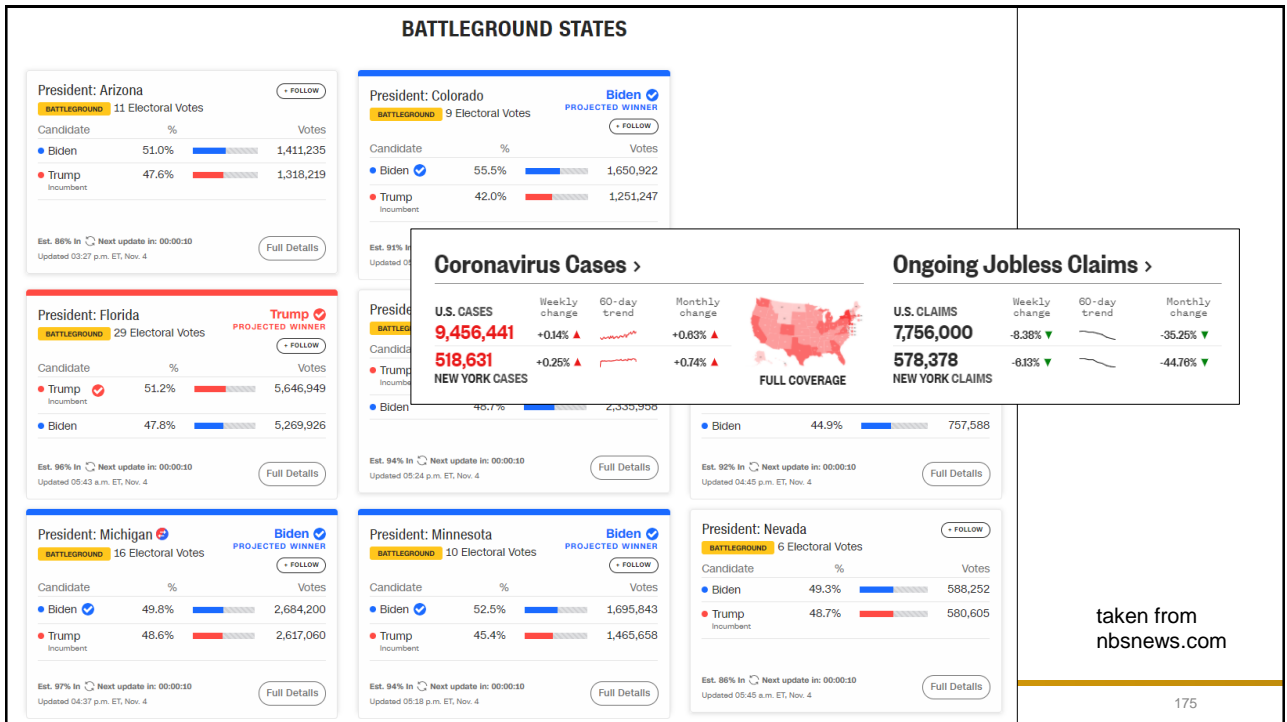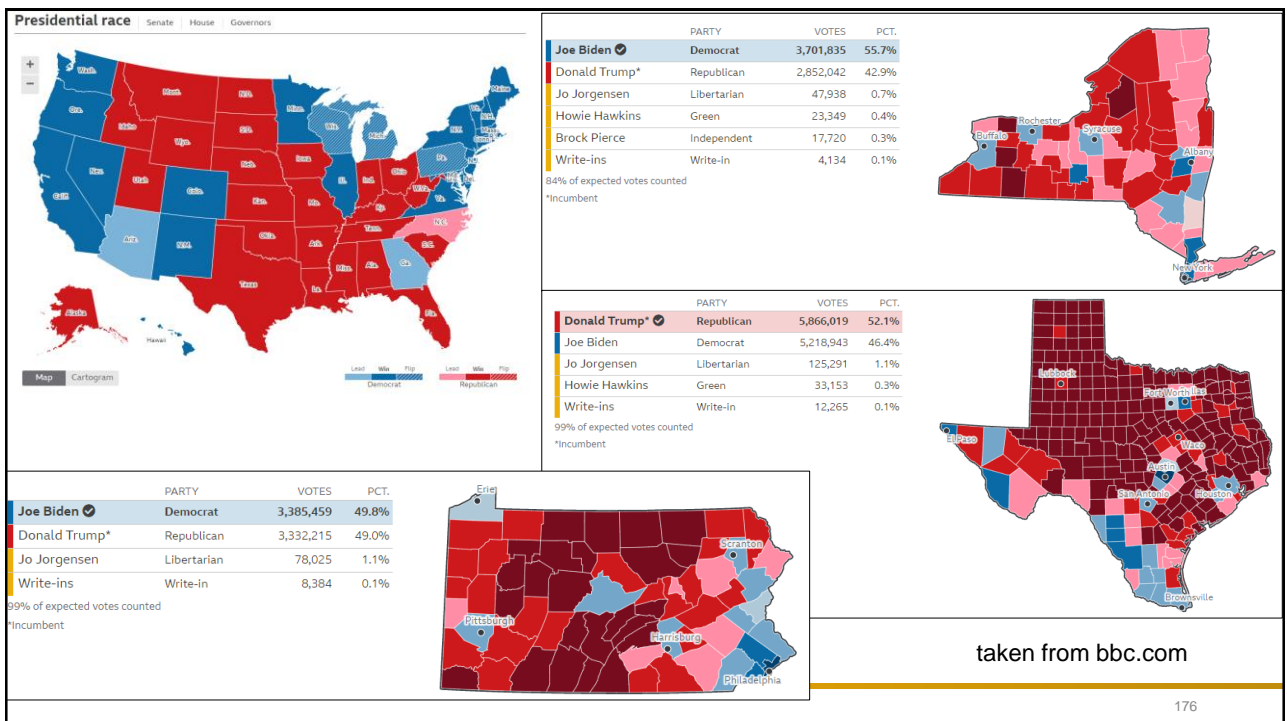


taken from bbc.com

174

174

175



176

https://www.nytimes.com/interactive/2020/11/03/us/elections/results-president.html

177

177