

The Graphs They Are a-Changin'

Principles, Examples, Software for Data Visualization

Constantin Manuel Bosancianu **and** Joost van Beek

*Doctoral School and Center for Media and Communication Studies, Central European
University*

April 26, 2012

Plan

Things to speak about:

- ① Basics of *good* data visualization;
- ② “The *good*, the *bad*, and the *ugly*” when it comes to data visualization - examples;
- ③ Software (open-source, web-based...);
- ④ Discussion time.

Importance

There is more data than ever waiting to be analyzed, mined for patterns, summarized, or linked to other data.

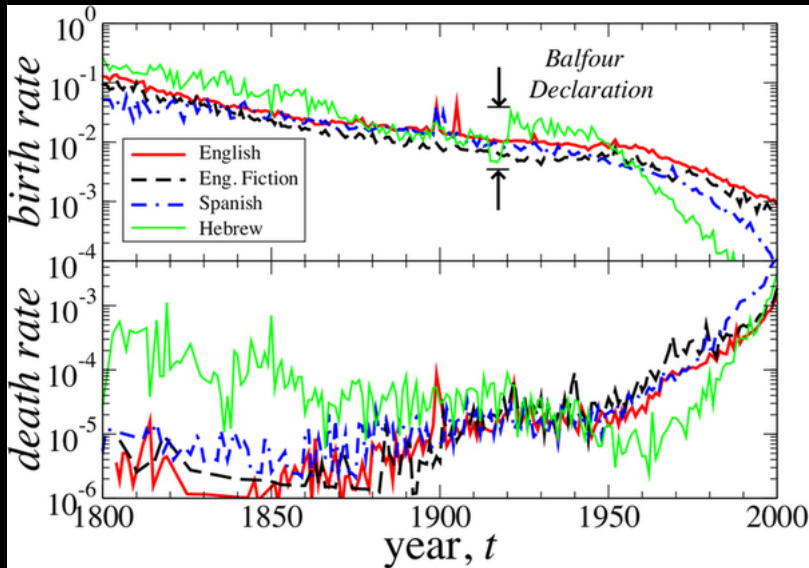


Figure: Word birth and death.

(<http://www.nature.com/srep/2012/120315/srep00313/full/srep00313.html>)

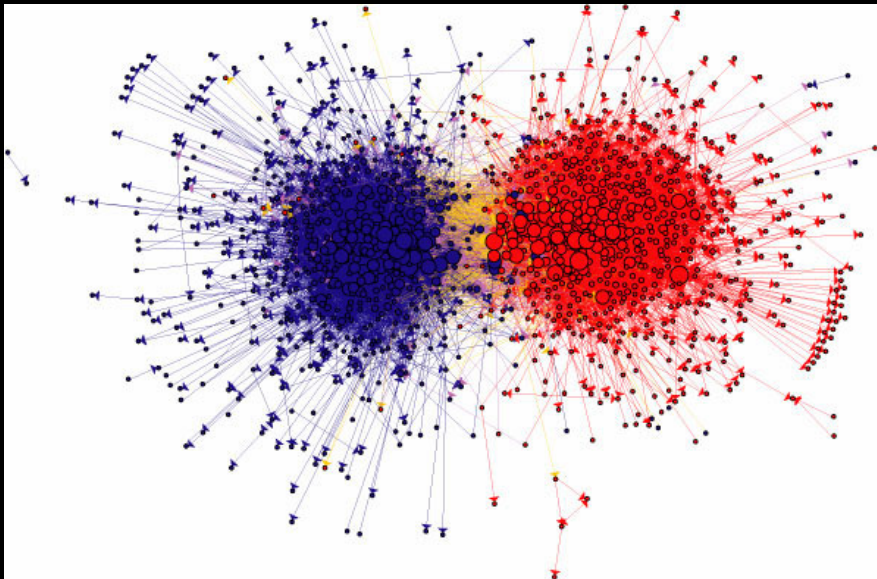


Figure: Linking patterns between US political blogs

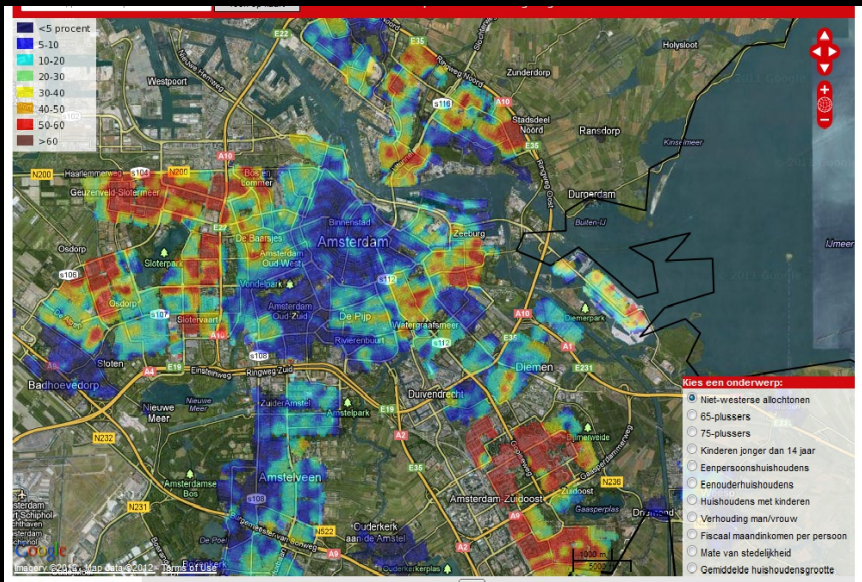


Figure: Immigrant clusters in Amsterdam

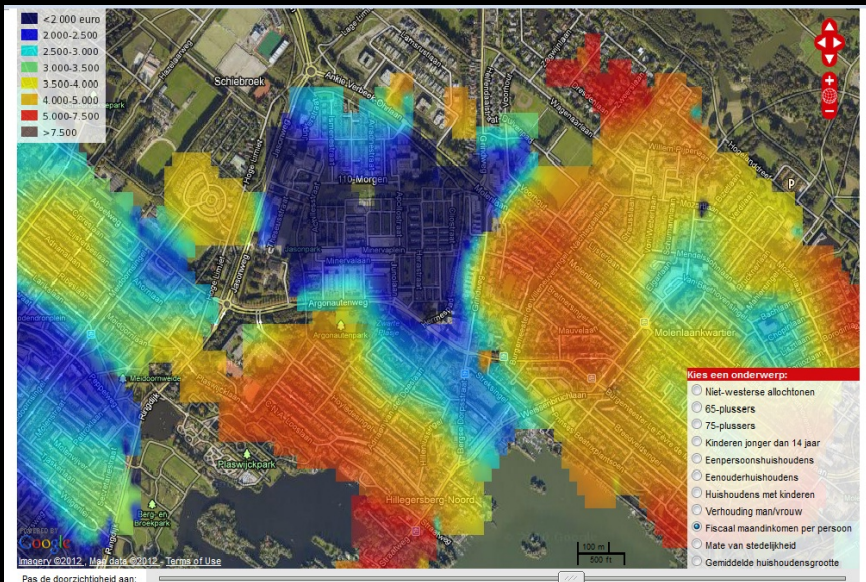


Figure: Income clusters in Rotterdam

Importance

We also observe a phenomenal level of growth in individual-level data: Internet, smartphones, automated sensors etc.

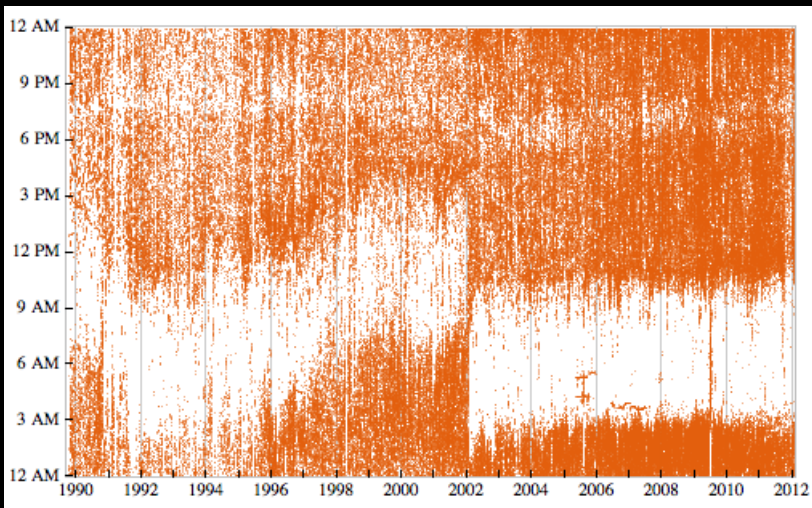


Figure: Stephen Wolfram's outgoing e-mail (approximately 300.000)

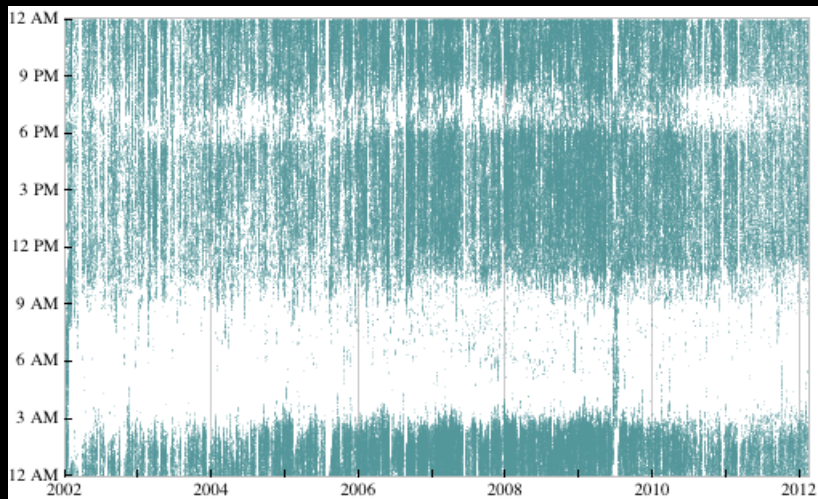


Figure: Stephen Wolfram's keystrokes (approximately 100 million)

Importance

Presenting this information in an accurate and intuitive way for the purpose of highlighting causal connections will be crucial for our ability to make adequate choices in a democracy.

1

Data visualization (DV)

- At the confluence between statistics and design, dealing with the search for the most effective and graphically intuitive way of making an argument on the basis of data.
- In 2000, an estimated 900 billion ($9 * 10^{11}$) to 2 trillion ($2 * 10^{12}$) graphs were generated every year (Tufte 2001).

Goals of DV

Multiple:

- Making an argument;
- Minimizing any distractions from the central argument;
- Ensuring the integrity of the argument;¹
- Summarizing a lot of information in a reduced space;
- Encouraging comparison.

¹“Making a presentation is a moral act as well as an intellectual activity.” (Tufte 2006, 141)

Principles of DV

- The overarching purpose is to show the data;
- Minimize the data-ink ratio, as much as possible;
- Erase non-data-ink, as much as possible;
- Minimize redundant data-ink, as much as possible;
- Revise and edit;
- Mobilize every graphical element needed.²

²Adapted from Tufte (2001)

ACCENT principles I

- **A**pprehension: Ability to correctly perceive relations among variables
- **C**larity: Ability to visually distinguish all the elements of a graph
- **C**onsistency: Ability to interpret a graph based on similarity to previous graphs

ACCENT principles II

- **Efficiency:** Ability to portray a possibly complex relation in as simple a way as possible
- **Necessity:** The need for the graph, and the graphical elements
- **Truthfulness:** Ability to determine the true value represented by any graphical element by its magnitude relative to the implicit or explicit scale³

³Source: D. A. Burn (1993), "Designing Effective Statistical Graphs".
In C. R. Rao, ed., *Handbook of Statistics*, vol. 9, Chapter 22.

Variable	Model 1	Model 2
Age	.027*** (.005)	.031*** (.006)
Gender	.094 (.174)	.074 (.215)
Education	.191*** (.044)	.055 (.056)
Marital status	.135 (.181)	.095 (.222)
Mobilized	-	.049 (.117)
Political interest	-	.733*** (.150)

Table: Estimates from a logistic regression model predicting likelihood of turnout (Sweden, EES 2009)

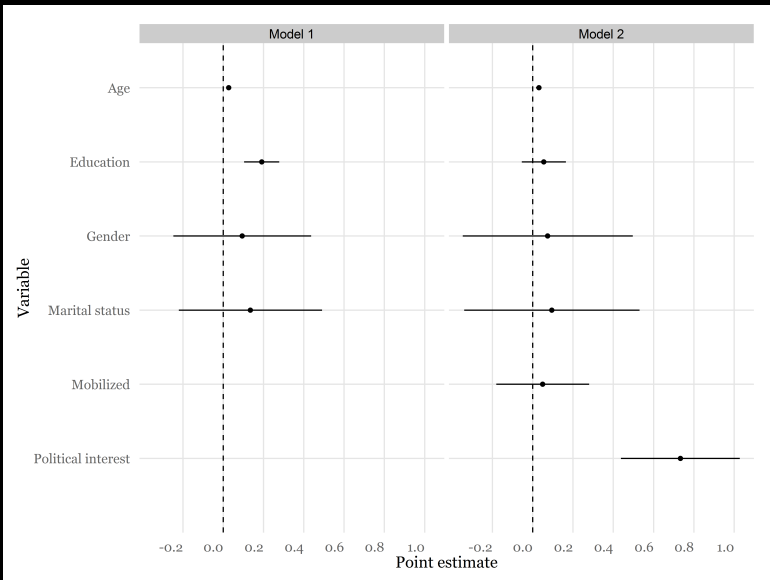


Figure: Estimates from the regression model in graphical form

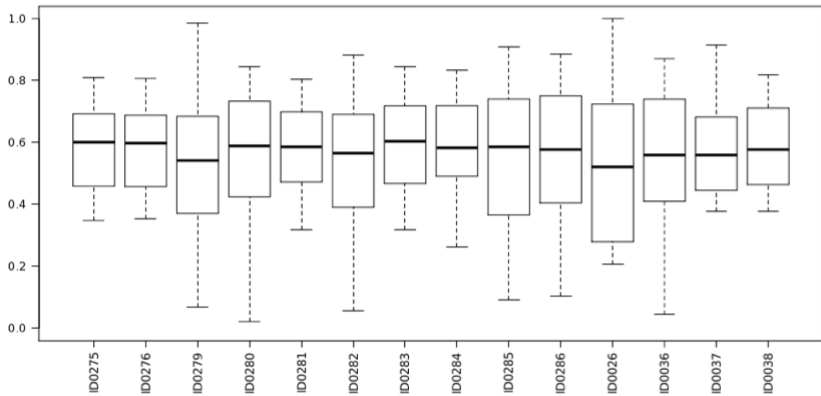


Figure: Traditional boxplot

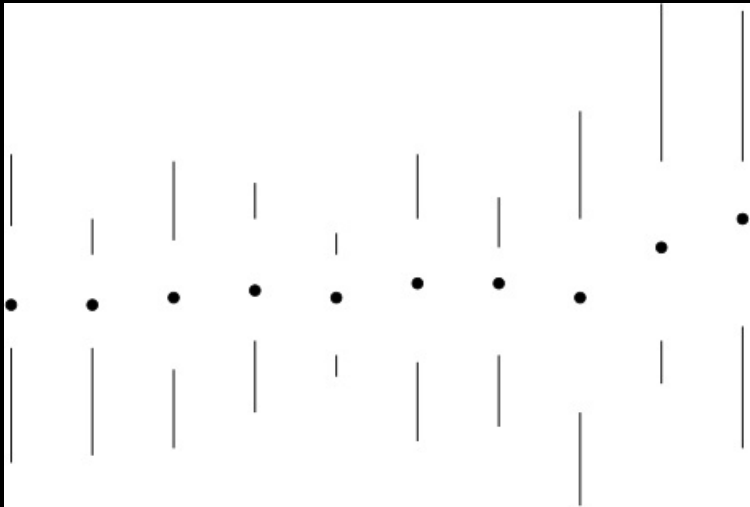
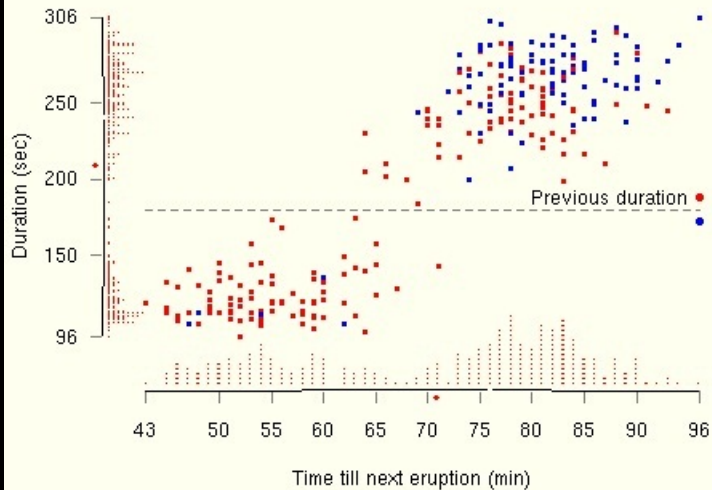


Figure: Quartile plot

Old Faithful Eruptions (271 samples)



2

2.1

Napoleon's 1812-1813 Russian campaign - Charles Joseph Minard.

Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.

Deviné par M. Misset, Inspecteur Général des Ponts et Chaussées en retraite, Paris, le 20 Novembre 1869.

Les arabes d'hommes peints sont représentés que les langues des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en lettres des zones. Le tracé dirige les hommes qui entrent en Russie; le noir ceux qui en sortent. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Cligny, de Tessier, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Leur mieux faite jugée à l'œil la diminution de l'armée, j'ai supposé que les corps de Blücher, Suvorov et de M. de Wittgenstein qui avaient été attachés sur Minsk et Mohilew n'avaient pas quitté ces lieux, ainsi que Wittgenstein, ainsi que Wittgenstein, ainsi que Wittgenstein.

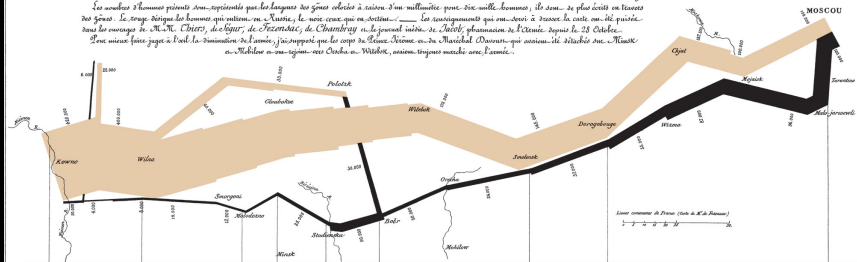
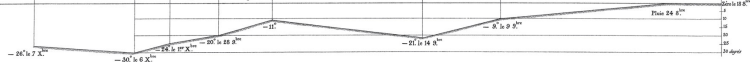


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Courbes peuvent se yuler la même nuit.



Imprimé par Bachelier, à Paris, 179, Rue de la Harpe, N° 179.

Dep. Lit. Répertoire de la Bibliothèque.

Figure: Campaign map

Attrition of Napoleon's Army in Russia, 1812

By John Boykin, based on CJ Minard

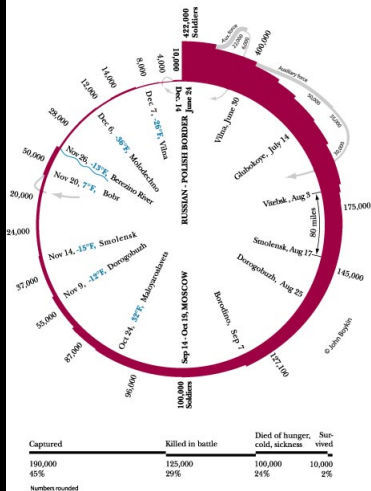


Figure: Alternative to the map

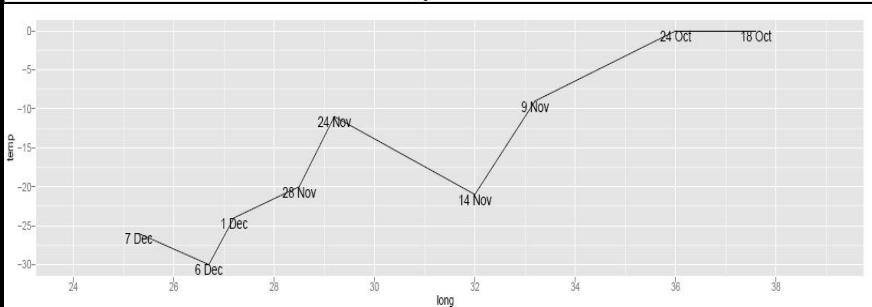
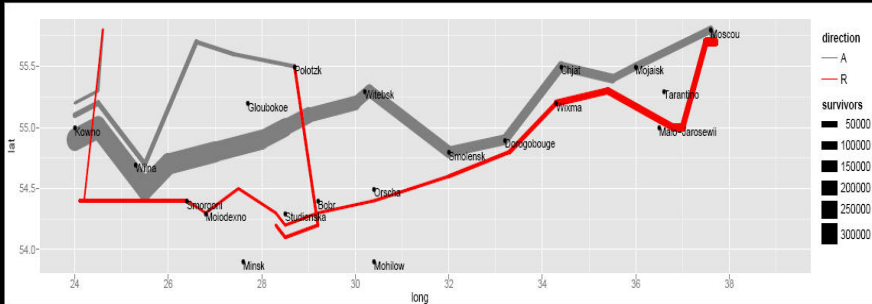
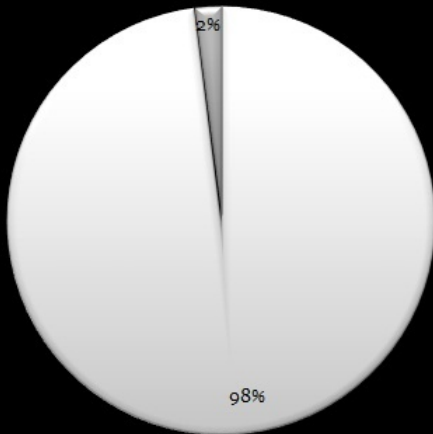


Figure: Alternative to the map

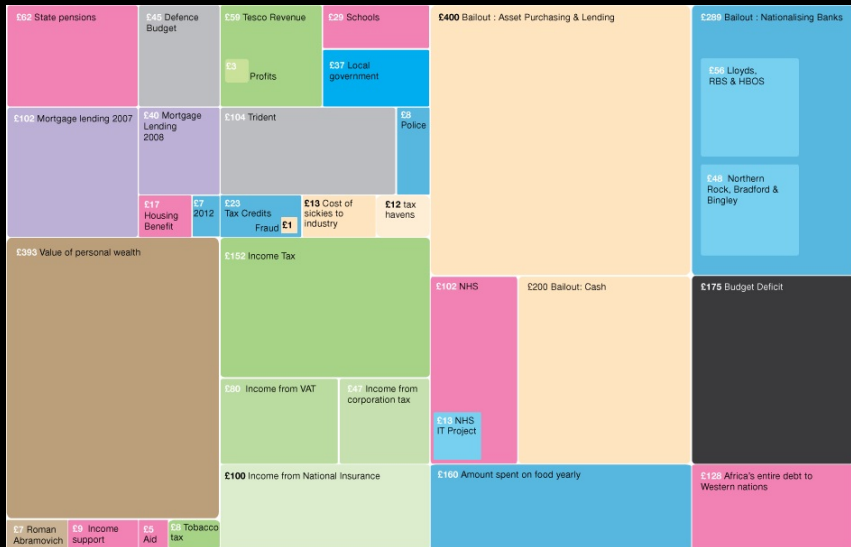
Napoleon's 1812-1813 Russian Campaign

■ Died ■ Survived



2.2

The UK Budget - David McCandless.



The Billion Pound-O-Gram

David McCandless / InformationIsBeautiful.net

● Giving
 ● Spending
 ● Fighting
 ● Hoarding
 ● Lending
 ● Bailing
 ● Earning

Source: UK Treasury, Guardian

2.3

Commuters in the US - SENSEable City Laboratory, MIT.

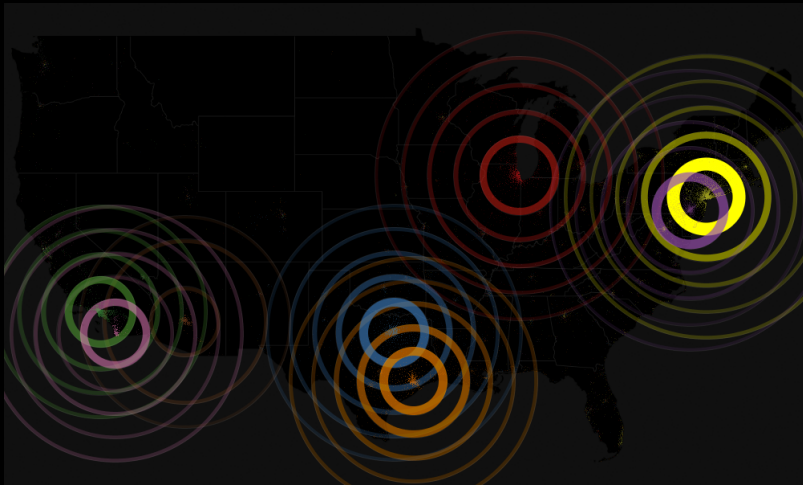


Figure: Commuters - July 2010, AT&T cell phone data

2.4

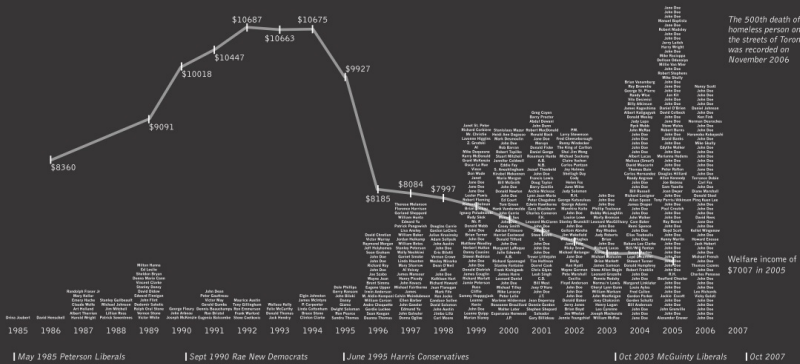
Welfare benefits in Ontario

Common Sense Revolution

Ontario Welfare Income for a Single Person in 2005 Constant Dollars & Homeless Persons Who Have Died on the Streets of Toronto 1985–2006

(National Council of Welfare & the Toronto Disaster Relief Committee)

© Scott Seirli



The 500th death of a homeless person on the streets of Toronto was recorded on November 2006

Welfare income of \$7007 in 2005

May 1985 Peterson Liberals
Sept 1990 Rae New Democrats
June 1995 Harris Conservatives
Oct 2003 McGuinty Liberals
Oct 2007



2.5

Web-based and interactive

The new frontier

- New York Times' *Mapping America*
- Washington Post's *Top Secret America*
- Wall Street Journal's *What They Know*
- Harvard's Berkman Center for Internet & Society
Mapping the Persian Blogosphere

3

3.1

'Chartjunk'

MONSTROUS COSTS

Total House and Senate
campaign expenditures,
in millions



Figure: Prominent example

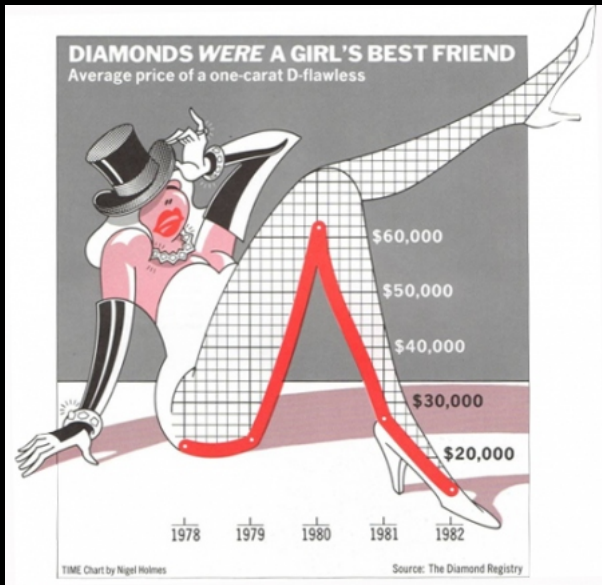
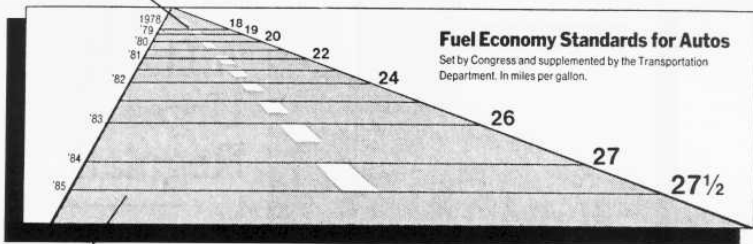


Figure: Prominent example

3.2

Misleading graphs

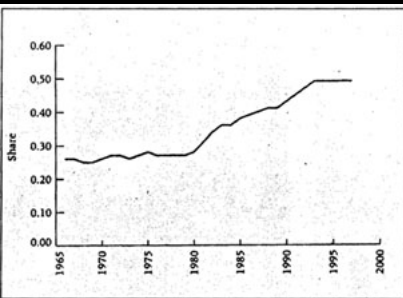
This line, representing 18 miles per gallon in 1978, is 0.6 inches long.



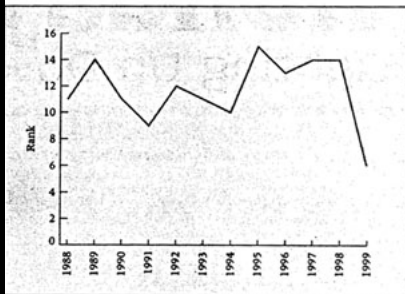
This line, representing 27.5 miles per gallon in 1985, is 5.3 inches long.

New York Times, August 9, 1978, p. D-2.

Figure: First example



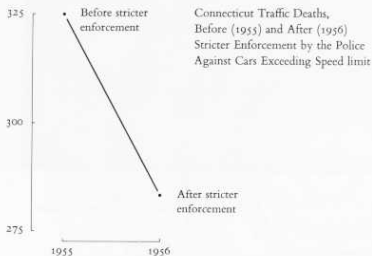
BY THE NUMBERS: OVER 35 YEARS, CORNELL'S TUITION HAS TAKEN AN INCREASINGLY LARGER SHARE OF ITS MEDIAN STUDENT FAMILY INCOME.



PECKING ORDER: OVER 12 YEARS, CORNELL'S RANKING IN *US NEWS & WORLD REPORT* HAS RISEN AND FALLEN ERRATICALLY.

Graphics must not quote data out of context.

Nearly all the important questions are left unanswered by this display:



A few more data points add immensely to the account:

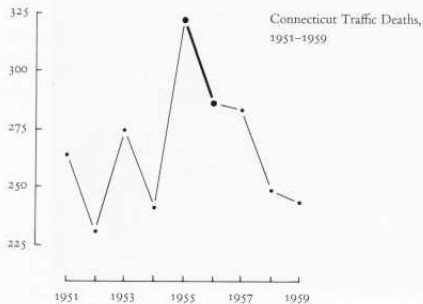


Figure: Third example

3.3

Poor understanding of statistics

2012 PRESIDENTIAL RUN

GOP CANDIDATES



SOURCE: OPINIONS

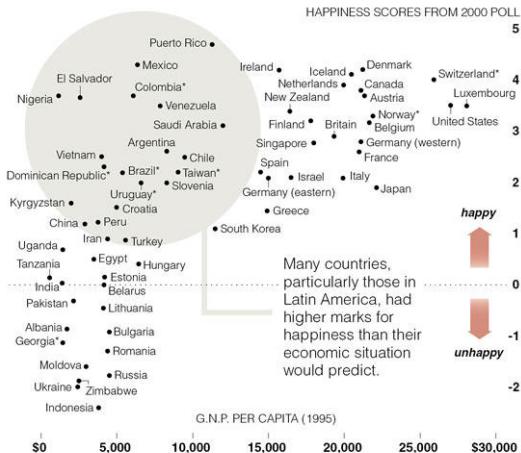
DYNAMIC

Figure: First example

A Plateau of Happiness

A country's wealth may not always dictate the happiness of its people.

As part of the World Values Survey project, inhabitants of different countries and territories were asked how happy or satisfied they were. Below is a sampling of happiness rankings, along with economic status.



*Poll results for these countries were from 1995.

Source: Ronald Inglehart, "Human Beliefs and Values : A Cross-Cultural Sourcebook Based on the 1999-2002 Values Surveys"

Figure: Second example

3.4

Poor choice of graphical display

Average Voltage in Seawater

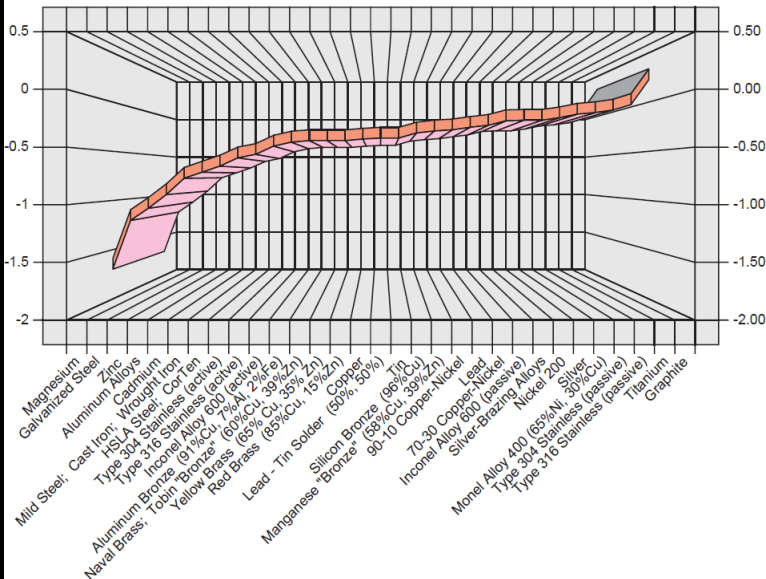


Figure: First example

Sotheby's / Christie's

Worldwide Sales Market Share Analysis

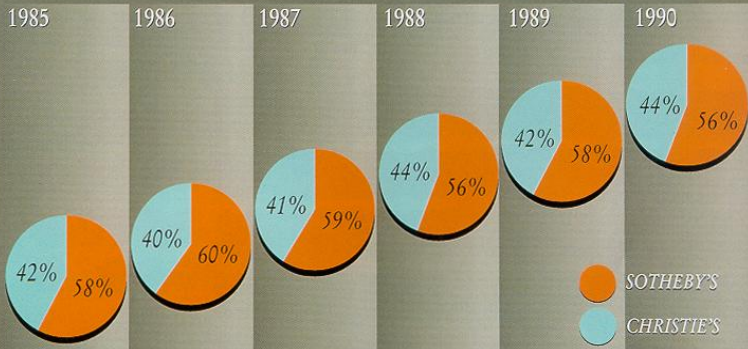


Figure: Second example

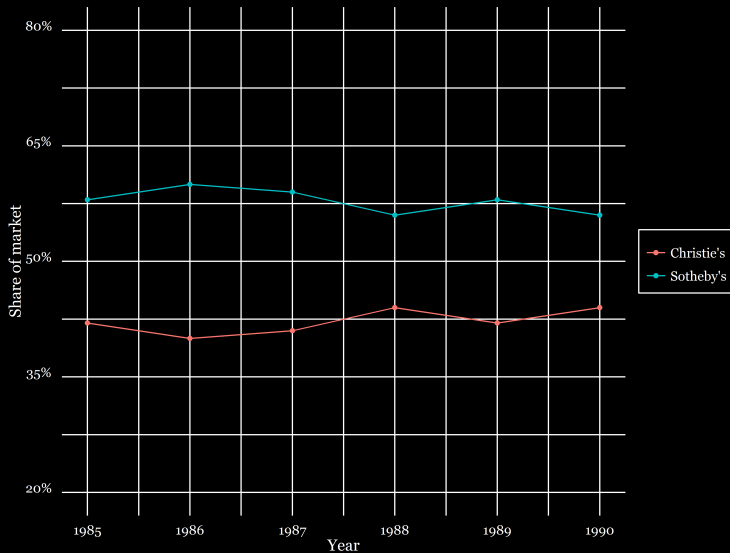


Figure: Alternative to second example

Chart 2 - Total Expenditures on Health as a Percentage Share of GDP, by OECD Country, 2004



Source: OECD Health Data 2007.

Note: For the United States the 2004 data reported here do not match the 2004 data point for the United States in Chart 1 since the OECD uses a slightly different definition of "total expenditures on health" than that used in the National Health Expenditure Accounts.

Figure: Third example

Expenditures on Health as Percentage of GDP for OECD Countries, 2004

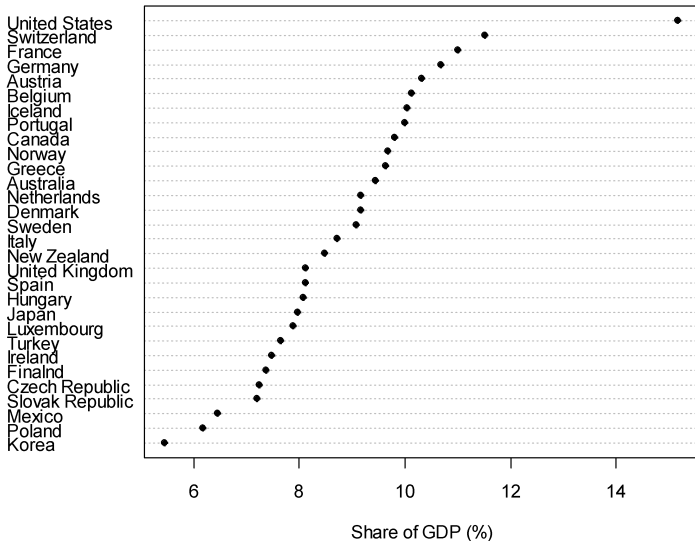


Figure: Reworked graph

4

Tools

To cover in the remaining minutes:

- Gapminder;
- IBM's Many Eyes;
- Web interface for ggplot2;

4.1

IBM's Many Eyes

<http://www-958.ibm.com/software/data/cognos/manyeyes/>

A “shared visualization and discovery” service, still in experimental phase

4.2

Hans Rosling's Gapminder project



Figure: Hans Rosling, Professor of International Health, Karolinska Institute, Stockholm, Sweden

Gapminder

- The problem he identifies: there is an abundance of yearly indicators for phenomena, scattered in the public domain
- Creates Gapminder Foundation and develops the Trendalyzer software (later sold to Google)
- Recently: Gapminder Desktop

Gapminder

Google develops, on the basis of Trendalyzer, Google
Public Data Explorer
(<http://www.google.com/publicdata/directory>)

4.3

Jeroen Ooms' ggplot2 interface

ggplot2

- R package developed by Hadley Wickham, on the basis of Leland Wilkinson's ideas regarding visualization (*The Grammar of Graphics*)
- Heavily code-based
- Jeroen Ooms adds a simple web-based interface to the package (other packages: IRT, lme4)

Honorable mentions

Still worthy to explore for a bit:

- Drillet (basic, but free)
- StatSilk (maps with indicators)
- GNU Octave (high-level interpreted language for numerical computations)
- IBM's Many Bills (specialized)
(<http://manybills.researchlabs.ibm.com/>)
- Wordle (word clouds)

5

Conclusion

Good data visualization involves thinking about the argument to be made, making choices among alternatives, and taking into consideration issues such as audience, parsimony, integrity. It will rarely result from canned routines and default options found in statistical packages.

Thank *you!*

References I

Books used for ideas or graphs:

- Tufte, Edward R. 1997. *Visual Explanations - Images and Quantities, Evidence and Narrative*. Cheshire, CT: Graphics Press.
- Tufte, Edward R. 2001. *The Visual Display of Quantitative Information*. Cheshire, CT: Graphics Press.
- Tufte, Edward R. 2006. *Beautiful Evidence*. Cheshire, CT: Graphics Press.
- Wickham, Hadley. 2009. *ggplot2 - Elegant Graphics for Data Analysis*. New York: Springer.
- Wilkinson, Leland. 2005. *The Grammar of Graphics*. New York: Springer.

References II

Internet sources where some of the graphs can be found:

- <http://www.informationisbeautiful.net/> (David McCandless, UK)
- <http://www.datavis.ca/gallery/index.php> (Michael Friendly, York University)
- <http://flowingdata.com/>
- <http://www.infosthetics.com/>
- <http://senseable.mit.edu/> (SENSEable City Laboratory, MIT)
- <http://chartporn.org/2012/03/02/improving-on-minard/>
- <http://igraphicsexplained.blogspot.com/>

References III

Web-based software:

- Gapminder Desktop
(<http://www.gapminder.org/downloads/>)
- IBM's Many Eyes (<http://www-958.ibm.com/software/data/cognos/manyeyes/>)
- Jeroen Ooms' ggplot2 interface
(<http://rweb.stat.ucla.edu/ggplot2/>)
- StatSilk (<http://www.statsilk.com/>)
- Wordle (<http://www.wordle.net/>)