

Multivariate Data Visualization



Multivariate Data

- Multiple variables (columns) of data
 - 2 variables (x, y) is bivariate
 - 3 variables (x, y, z) is trivariate...
- Mostly concerned with 4 or more variables
- Variables can be of any type (or combination of types)

Data Types

- **Numerical**

- Continuous (-11.2, -1.3, 0.4, 4.8, 14.9, ...)
- Discrete (-2, -1, 0, 1, 2, ...)

- **Categorical**

- Ordered (January, February, March, April, ...)
- Unordered (Red, Blue, Green, Purple, ...)

Data Types

- **Special/Structured**

- Time (1998-01-27 12:00, 1999-04-13 22:10, ...)
- GPS (37.77679 latitude, -122.45117 longitude, ...)
- Social Security Numbers (555-55-5555, ...)

- **Unstructured/Semi-Structured**

- Free-Form Text (Twitter Feed, Screenplay, ...)

Outline

- Bubble Chart
- Heatmap
- Scatterplot Matrix
- Small Multiples
- Parallel Coordinates
- Other Techniques
- *Encoding Data*
- *Examples*
- *Discussion*

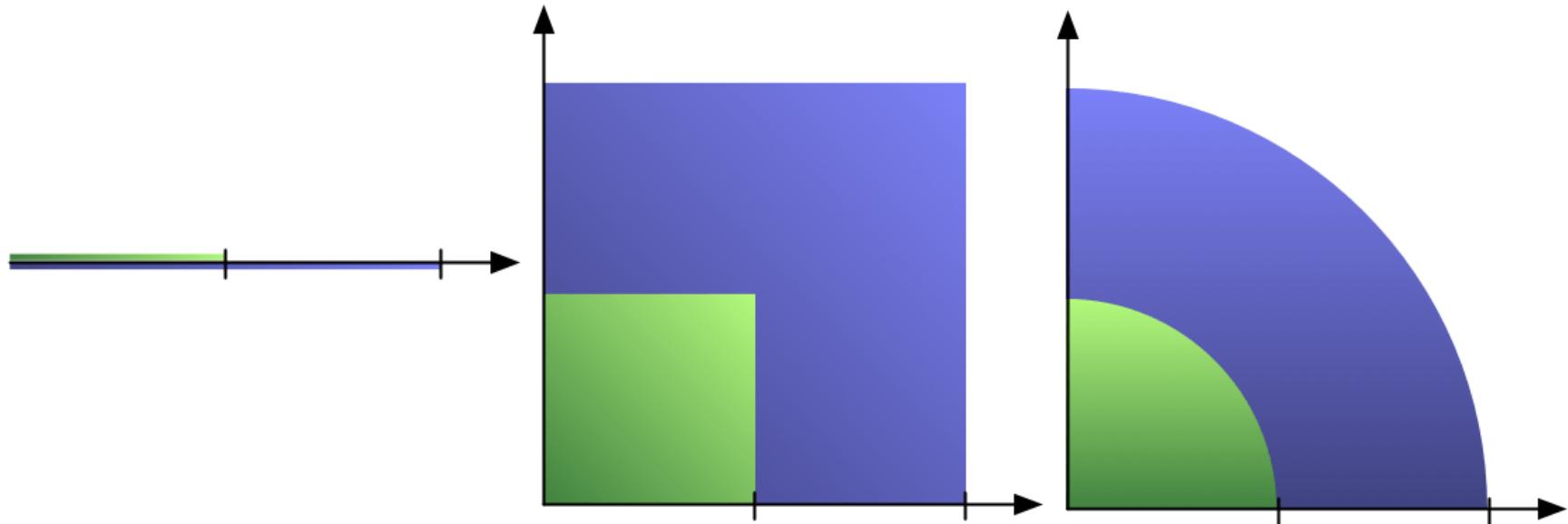
BUBBLE CHART

Encoding Data

- **Horizontal position**
 - Continuous data
- **Vertical position**
 - Continuous data
- **Circle area (*not radius*)**
 - Numerical data
- **Circle color**
 - Numerical or categorical

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

Linear vs Quadratic Change



<https://eagereyes.org/blog/2008/linear-vs-quadratic-change>

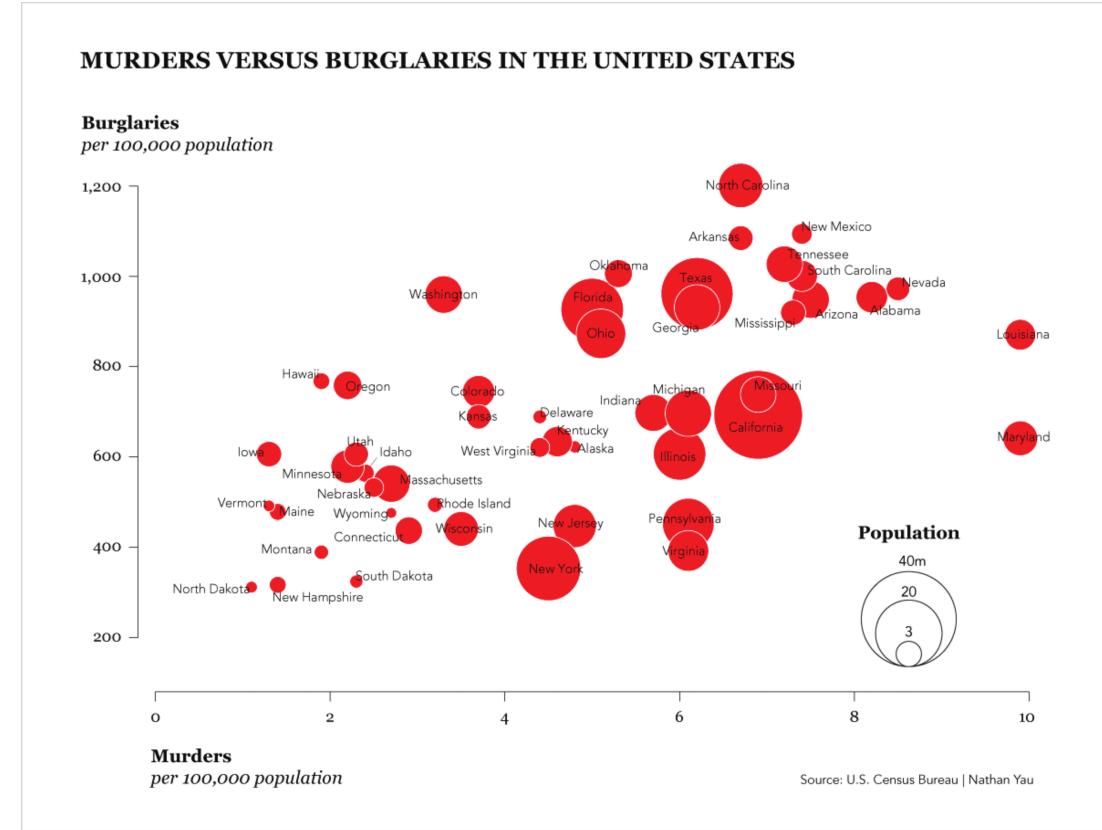
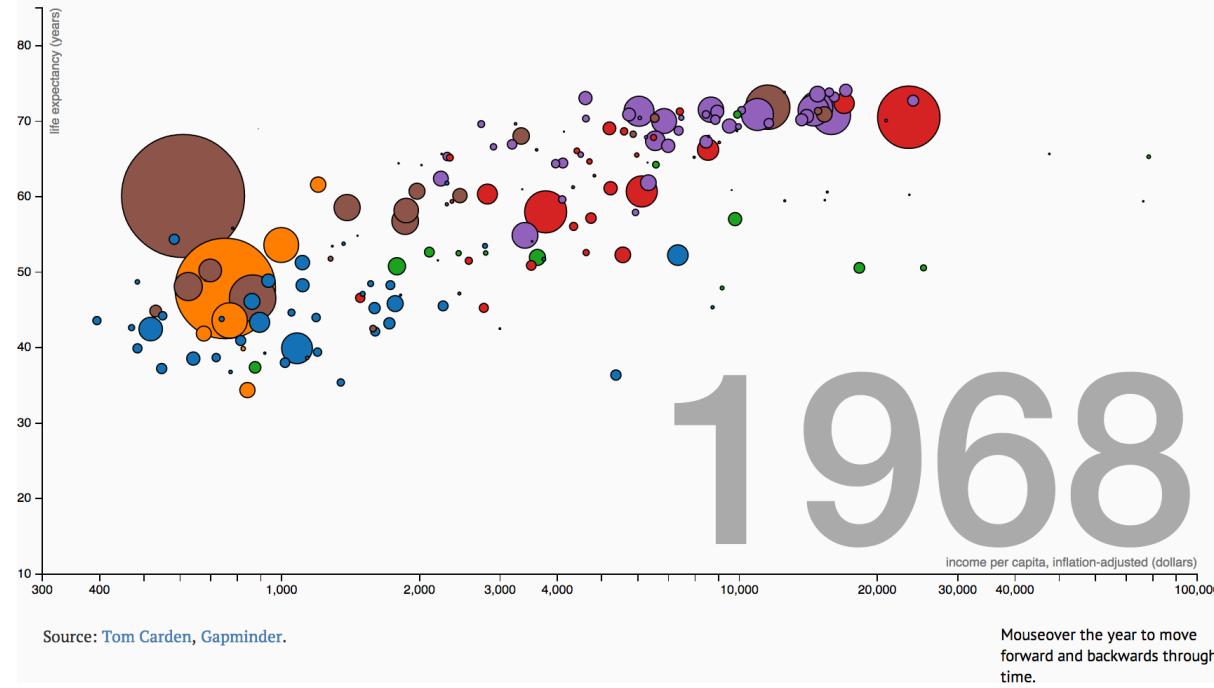


FIGURE 6-15 Bubble plot showing crime in the United States

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

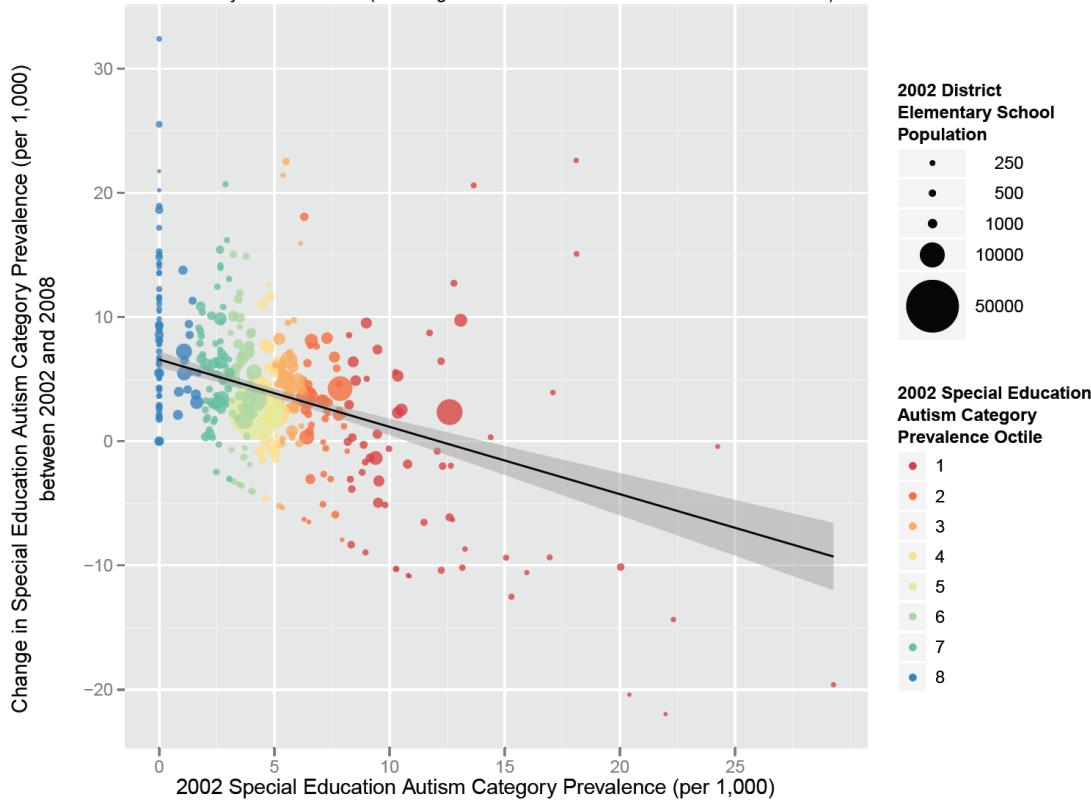
The Wealth & Health of Nations

March 13, 2012

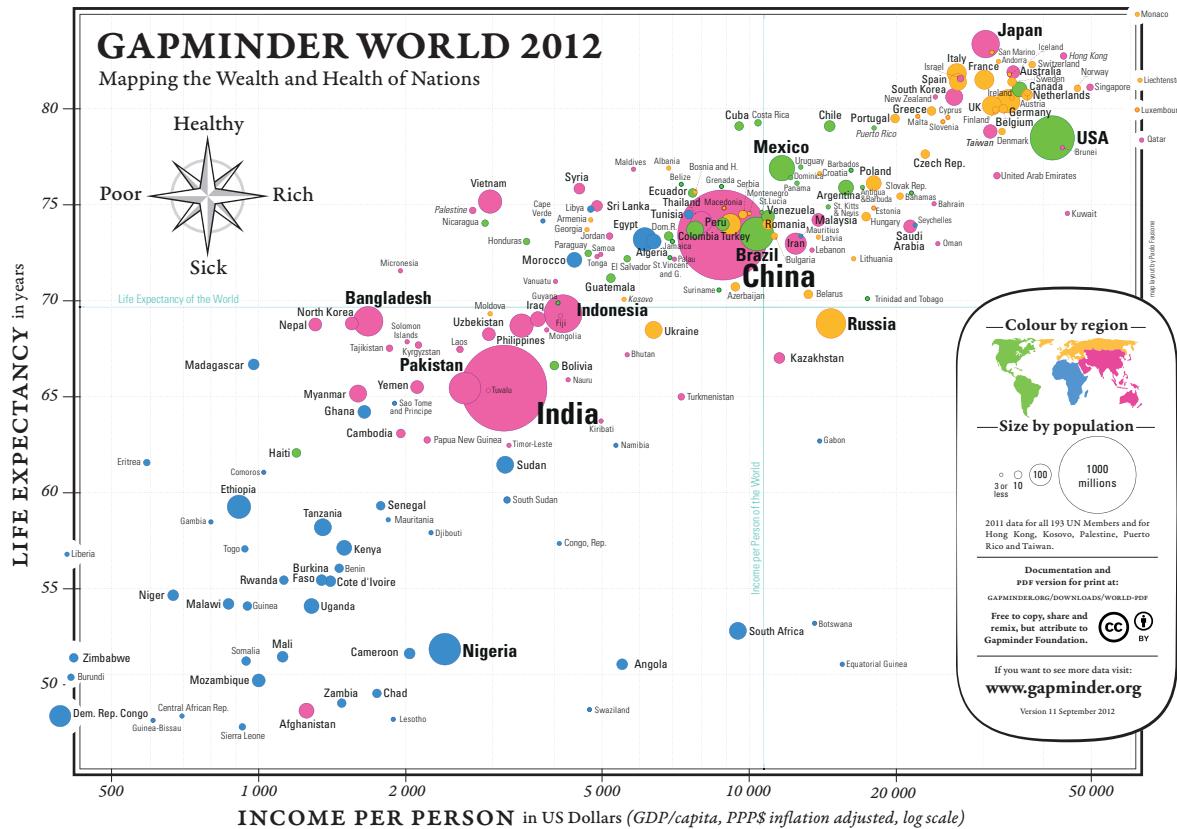


<http://bostocks.org/mike/nations/>

Figure 3. Change in Special Education Autism Category Prevalence between 2002 and 2008 vs Baseline (2002) Prevalence, Wisconsin Elementary School Districts (with weighted linear best-fit line and 95% confidence band)



<http://www.matthewmaenner.com/blog/?p=150>



<http://www.gapminder.org/downloads/world-pdf/>

Discussion

- Able to encode four dimensions of data
 - Ideal if one dimension is categorical (color)
- Rough comparison possible
 - Beware comparing circle areas
- Obscuring data may be an issue
 - Large circles should be behind smaller ones
 - Issues increases with density

HEATMAP

Encode Data

- **Horizontal and vertical position**
 - Row/column position in data table
 - Discrete or categorical variable
- **Box color**
 - Value from row, column in dataset
 - May be numerical or categorical

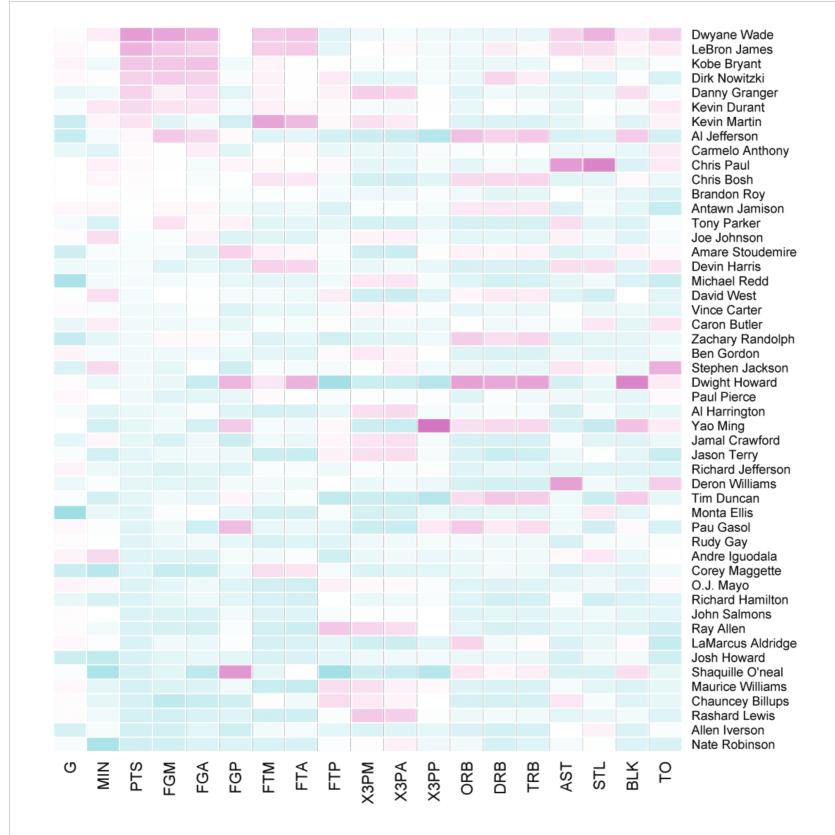
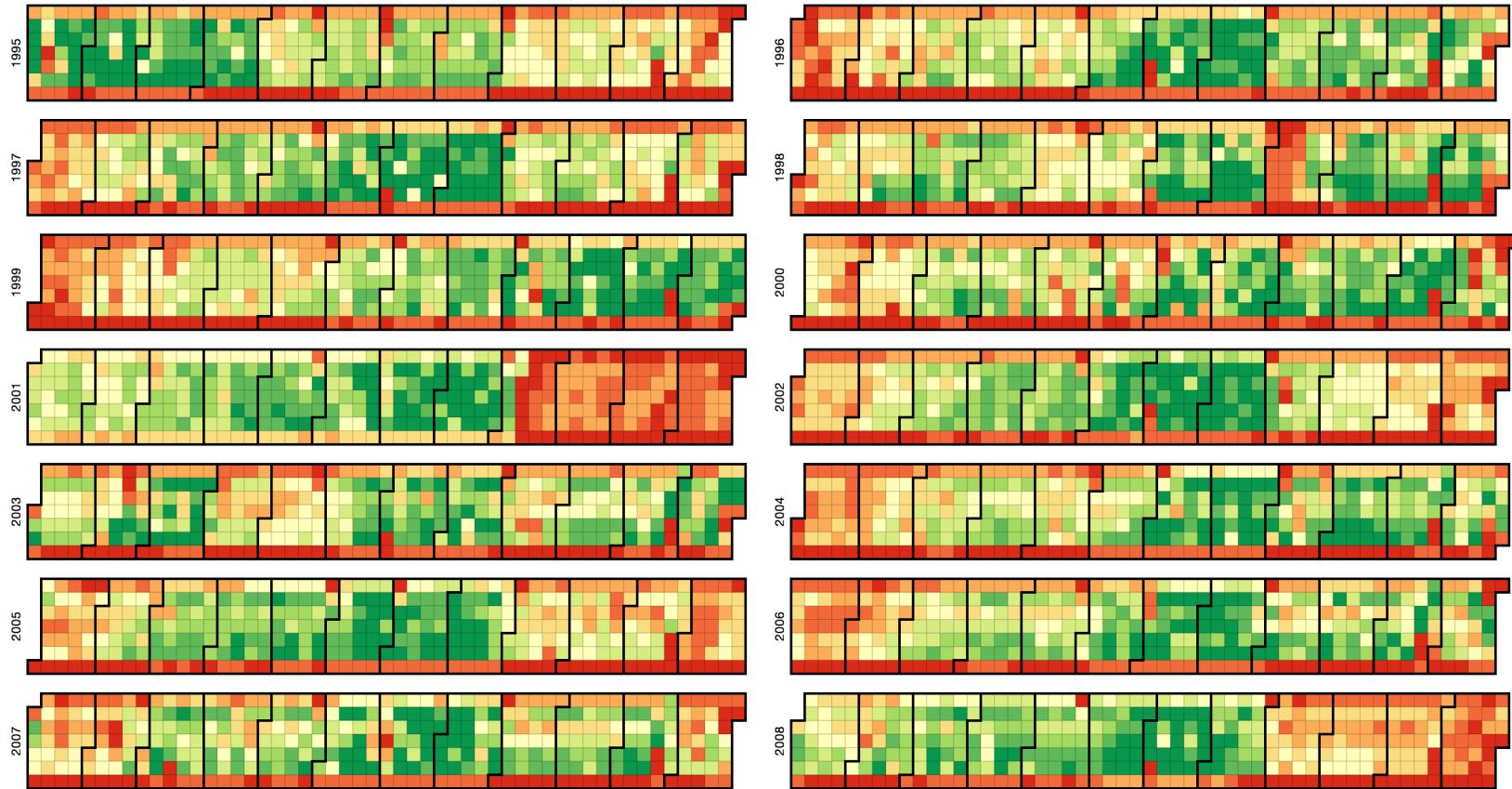
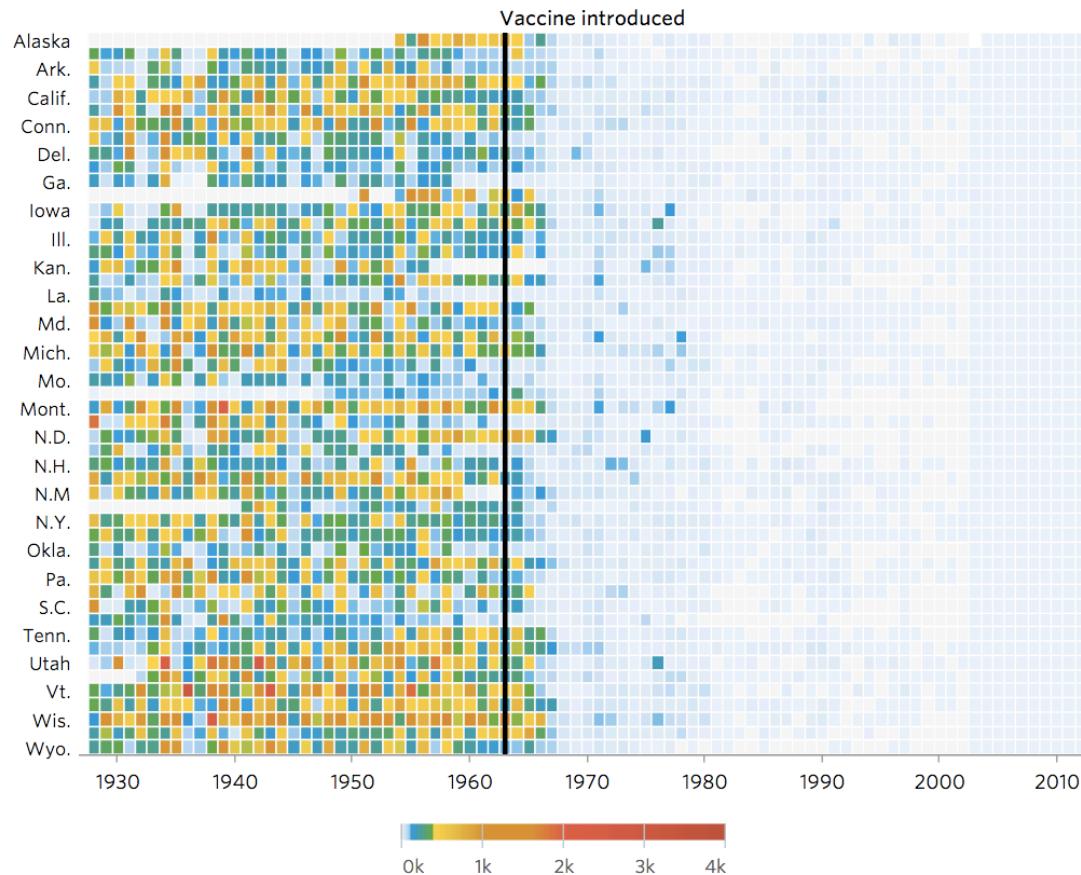


FIGURE 7-3 Default heatmap ordered by points per game

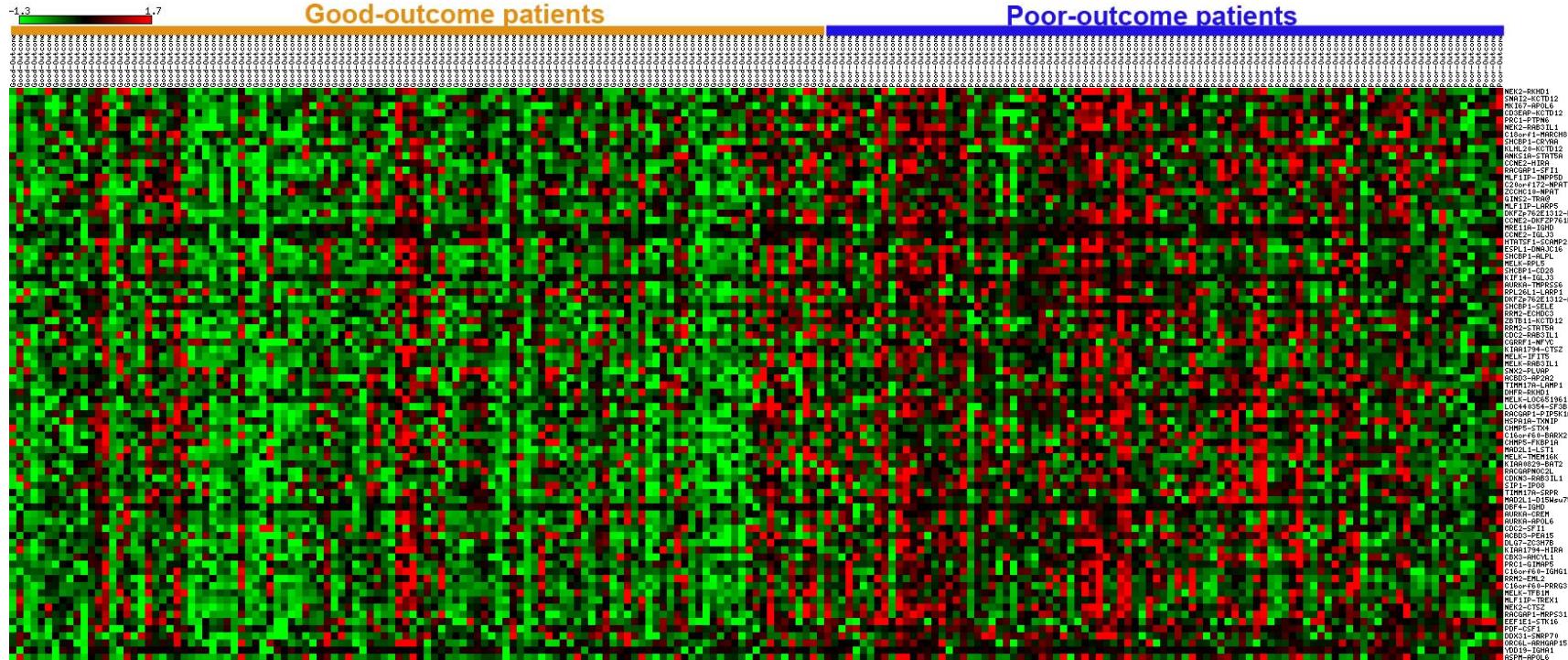
"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.



<http://mbostock.github.io/d3/talk/20111018/calendar.html>



<http://graphics.wsj.com/infectious-diseases-and-vaccines/>



<http://www.biomedcentral.com/1471-2105/9/125>

Discussion

- Essentially showing the data set, but replaces numbers by color values
- Good for certain types of data
 - Continuous data well-suited for cell values
 - Unordered categorical limited to 7 categories
- Encourages comparison and pattern finding
 - Sorting changes patterns!

SCATTERPLOT MATRIX

Encode Data

- **Scatterplot**
 - Horizontal position maps to one variable
 - Vertical position maps to another variable
- **Matrix of scatterplots**
 - Each scatterplot focuses on one pair
 - Which pair is determined by row and column

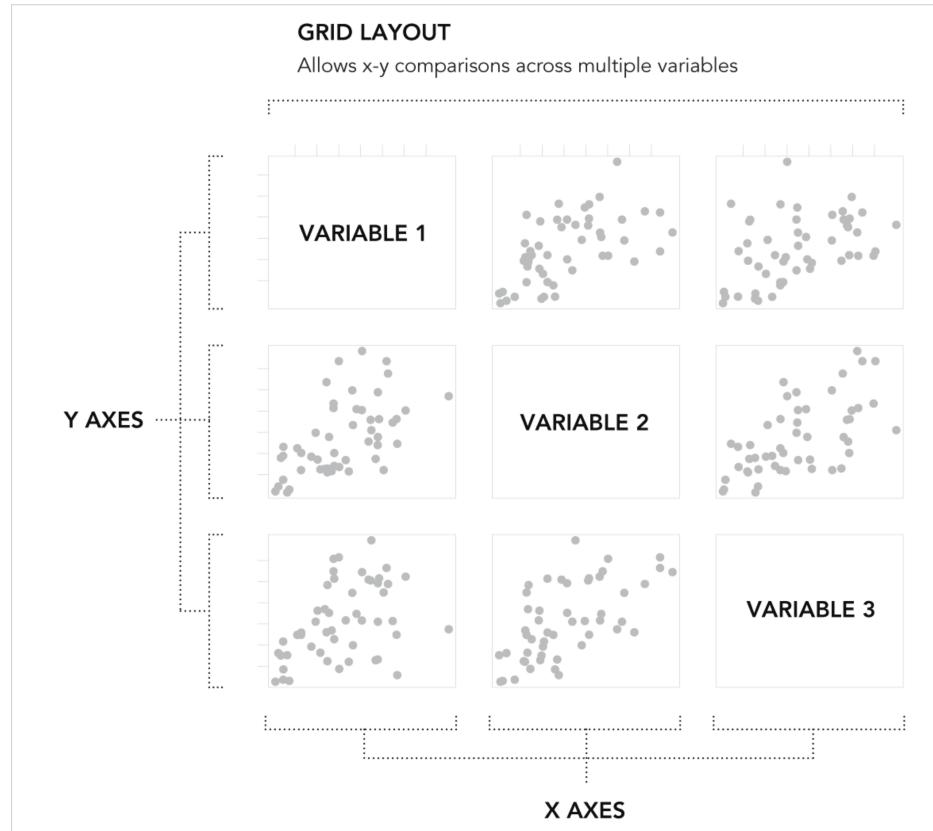


FIGURE 6-8 Scatterplot matrix framework

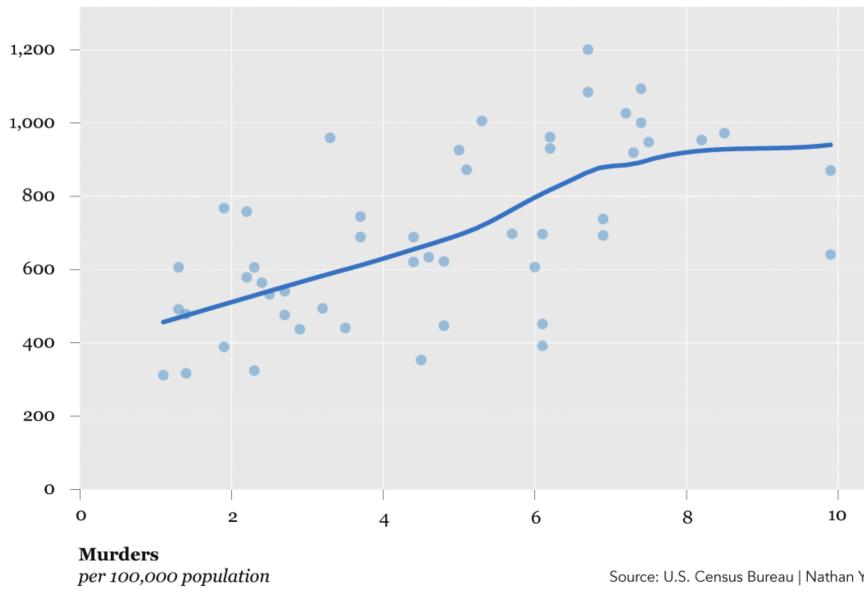
"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

MURDERS VERSUS BURGLARIES IN THE UNITED STATES

States with higher murder rates tend to have higher burglary rates.

Burglaries

per 100,000 population



Source: U.S. Census Bureau | Nathan Yau

FIGURE 6-7 Revised scatterplot on murder versus burglary

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

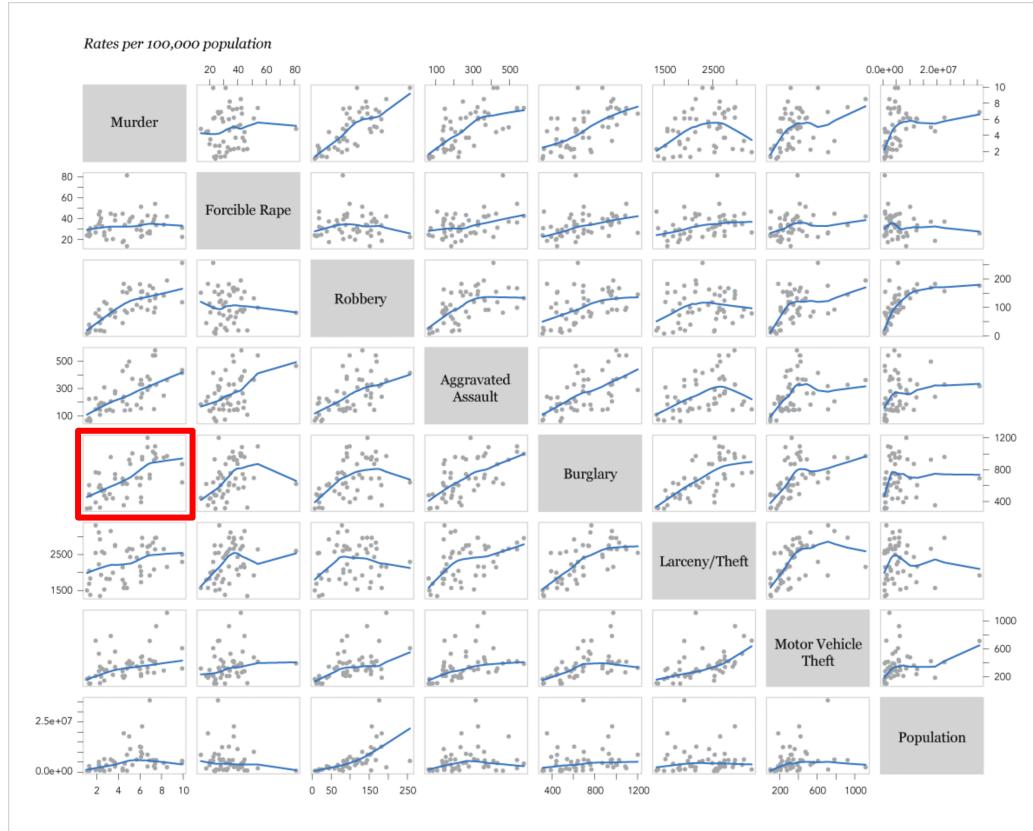
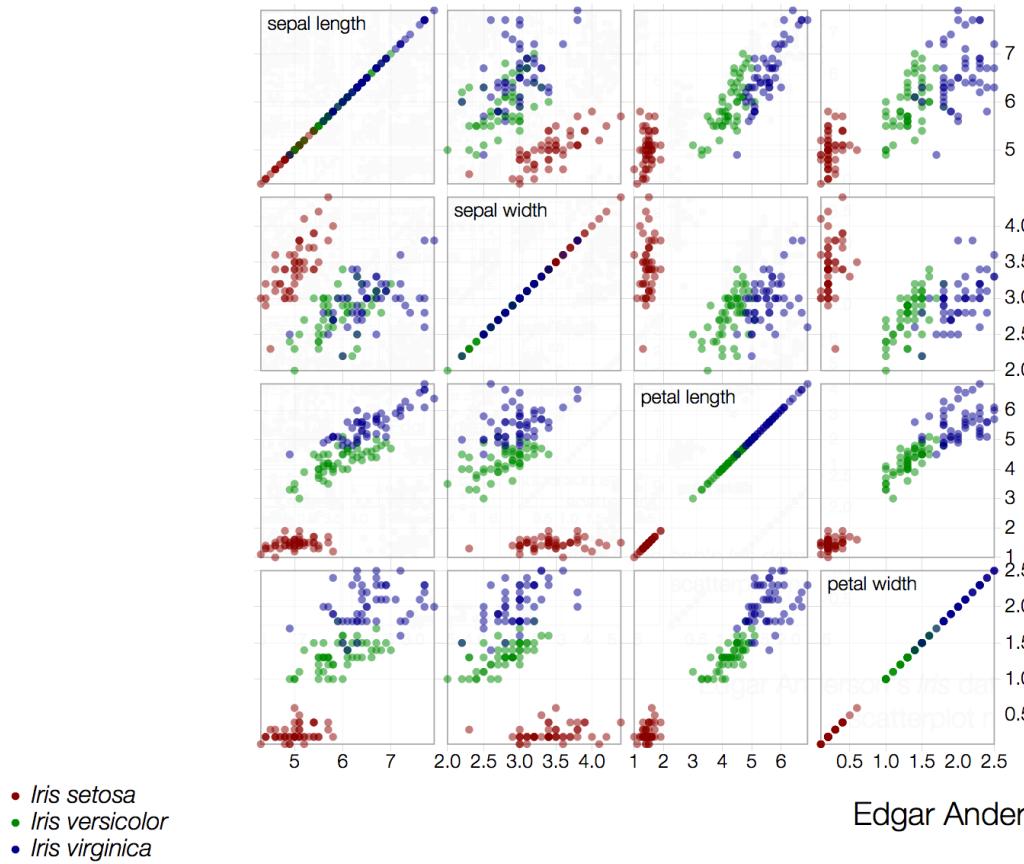


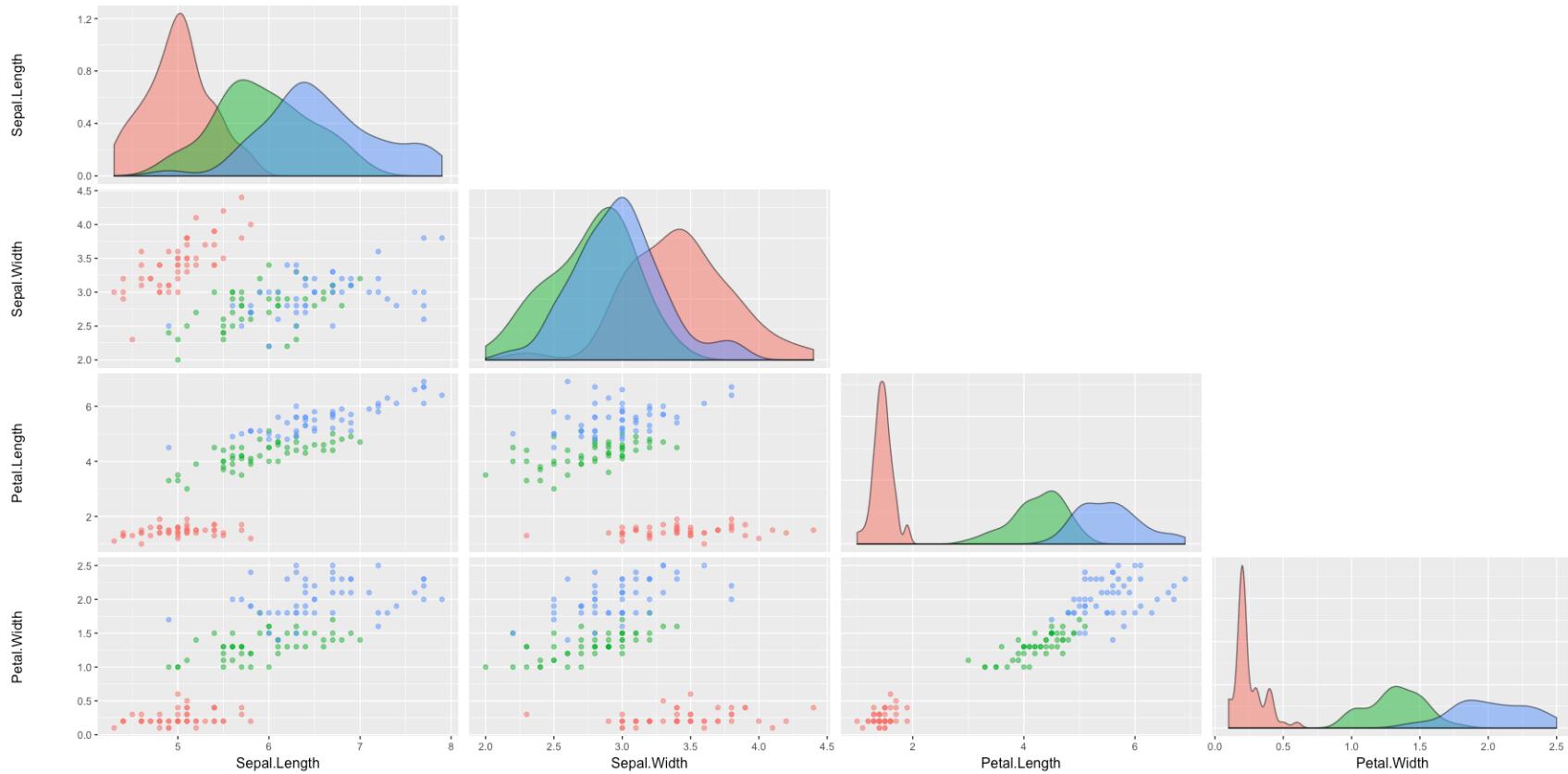
FIGURE 6-9 Scatterplot matrix of crime rates

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

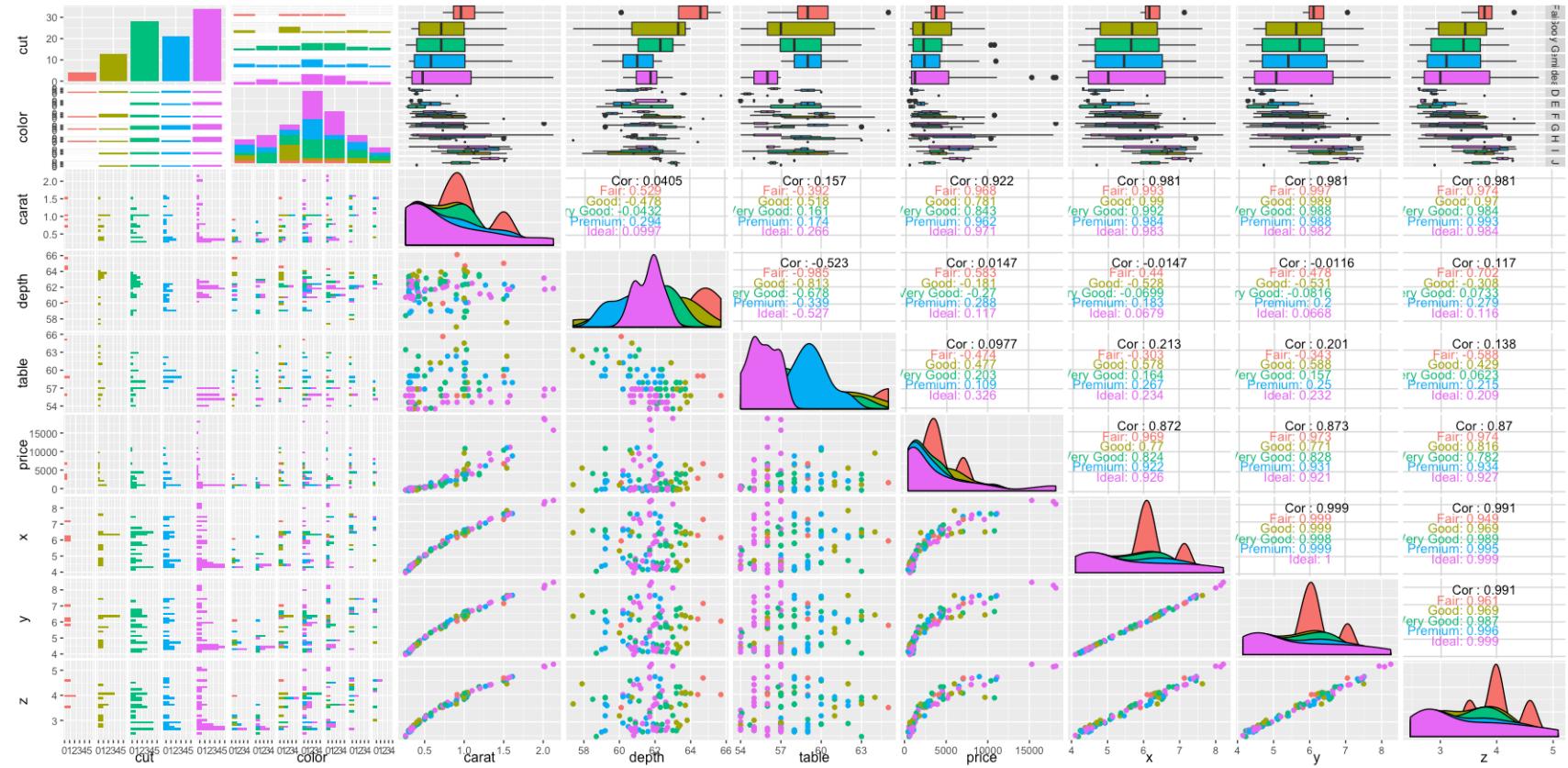


Edgar Anderson's *Iris* data set
scatterplot matrix

<http://mbostock.github.com/d3/talk/20111116/iris-splom.html>



<http://gettinggeneticsdone.blogspot.com/2011/07/scatterplot-matrices-in-r.html>



<https://cran.r-project.org/web/packages/GGally/vignettes/ggpairs.html>



Discussion

- Good for exploration and comparison
 - Can be a little overwhelming at first
- Works with numerical or ordered data
- Form of small-multiples plot
 - Multiple scatterplots

SMALL MULTIPLES

Encode Data

- Group data by a variable to divide it into subsets
 - Usually a categorical variable
- Create a small plot for each subset
- Show all subset plots on same page
 - Encourages comparison

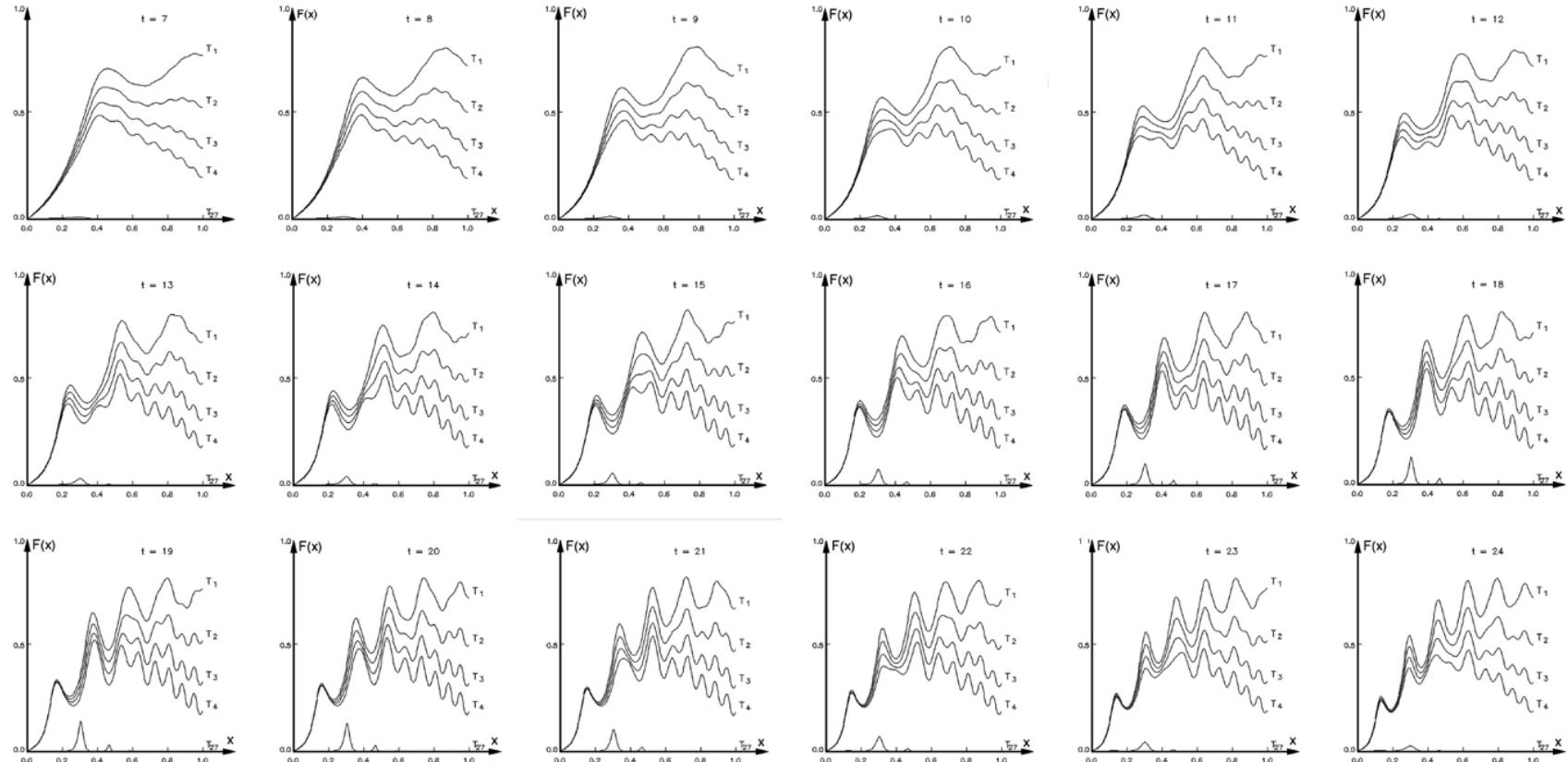
Fresh Originals / Rotten Finales

Out of the 35 selected trilogies, 23 of them had *fresh* originals; however, only 11 trilogies had a *fresh* finale.

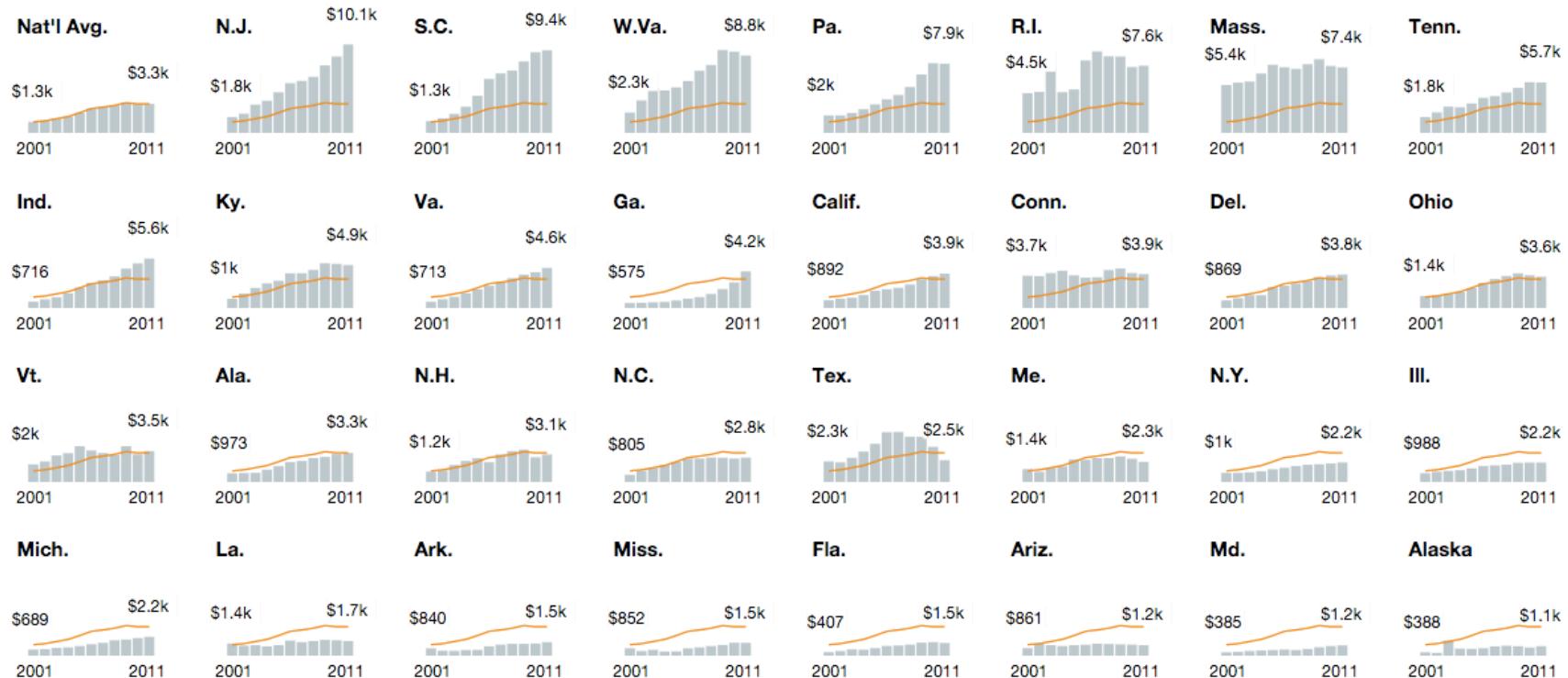
- █ Fresh (at least 60%)
- █ Rotten



"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.



<http://www.sv.vt.edu/classes/ESM4714/methods/CogVizCmp.html>



<http://projects.propublica.org/graphics/ambulances>

Discussion

- Excellent for comparison
 - Depends on how to place the multiples!
- Requires a variable to use for grouping
 - Discrete or categorical data
- Can use with any type of plot
- Harder to tell exact values

PARALLEL COORDINATES

Encode Data

- Create one vertical line for every column
 - Numerical or ordered data
- Plot every row
 - x position is column
 - y position is value for that column
 - line connects values for a single row
- Picture an xy scatterplot, but putting both axis lines vertically

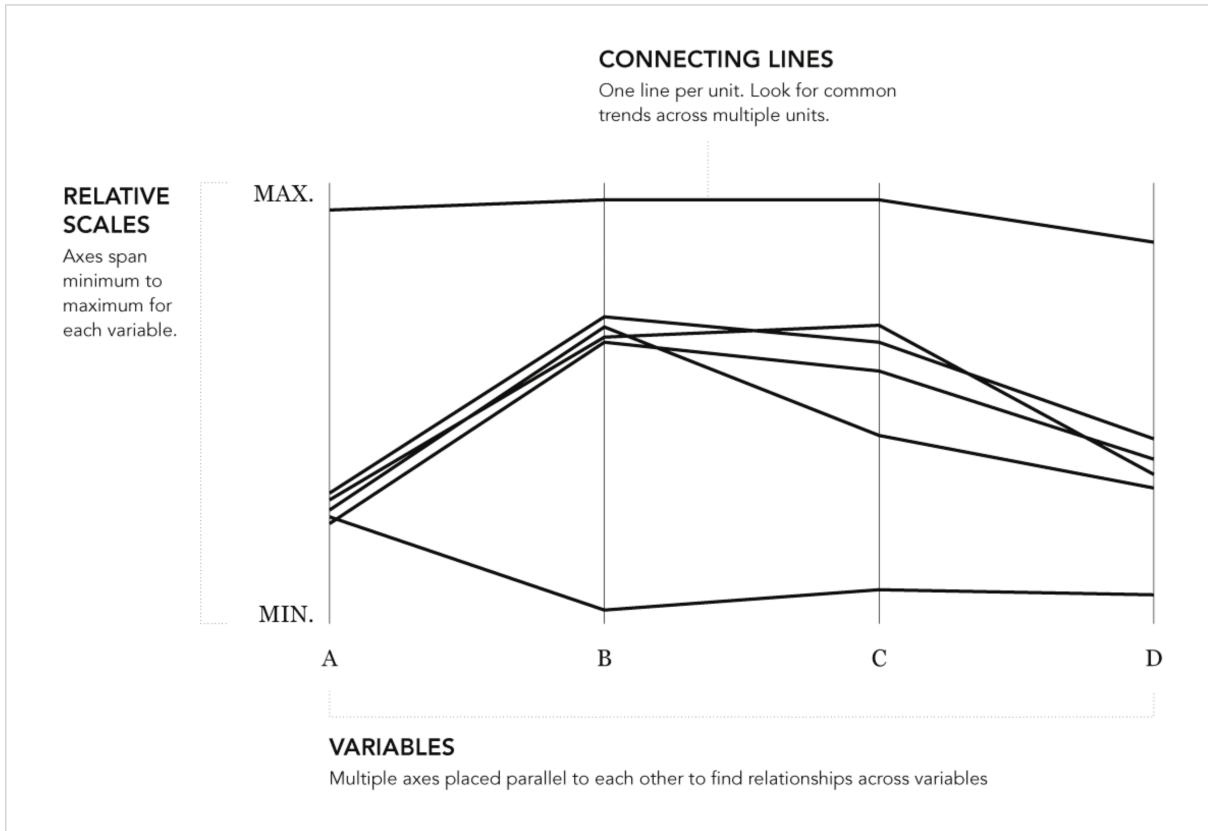


FIGURE 7-20 Parallel coordinates framework

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

EDUCATION IN THE UNITED STATES

States with higher SAT reading scores predictably also have relatively higher math and writing scores. However, this does not necessarily mean that students are better educated in these states. It's more likely an indicator for what percentage of graduates actually took the test.

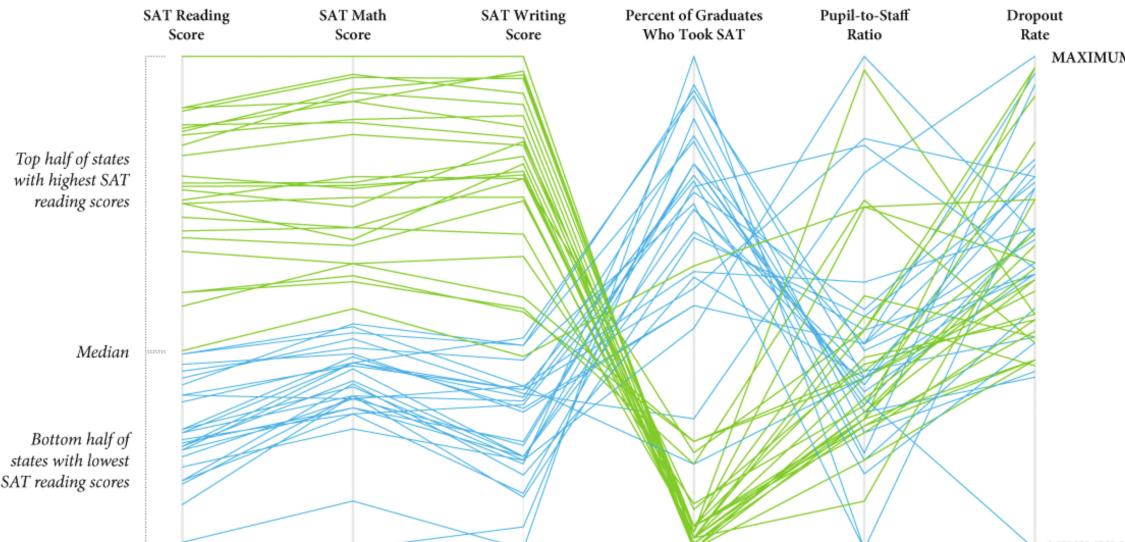
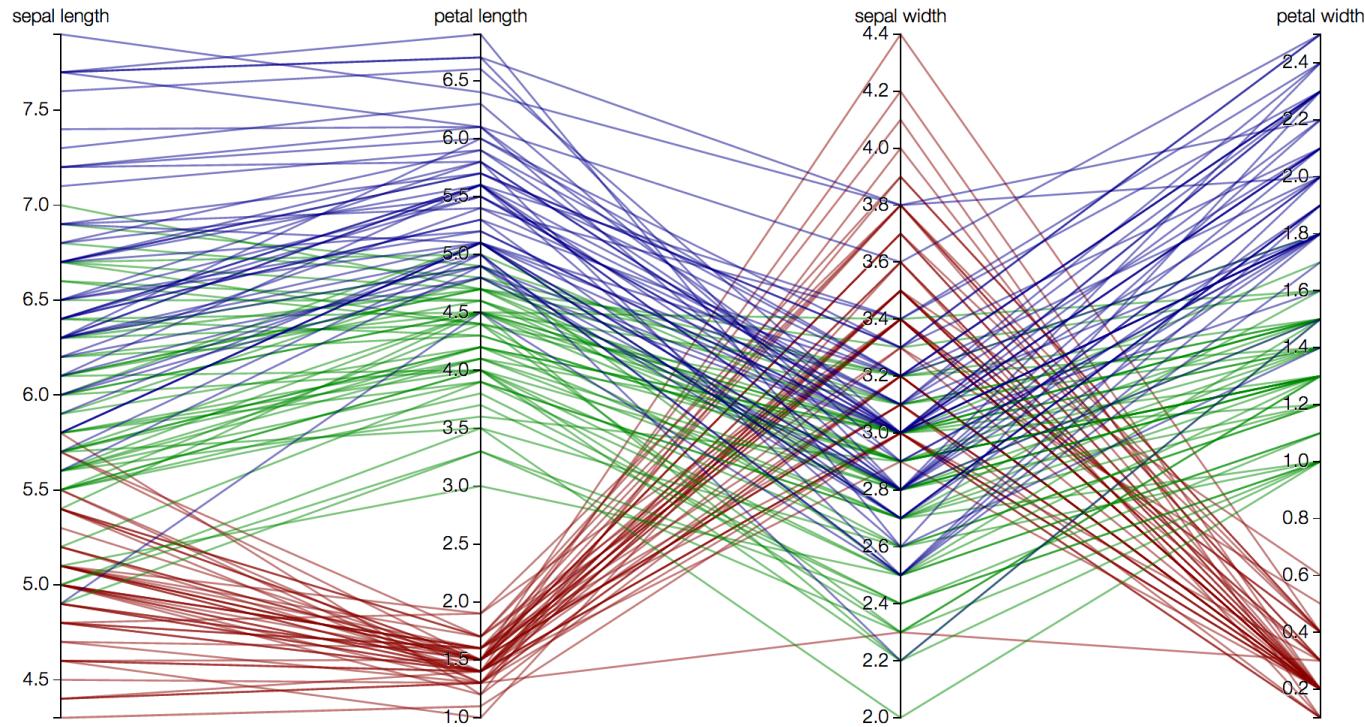


FIGURE 7-26 Standalone parallel coordinates plot on SAT scores

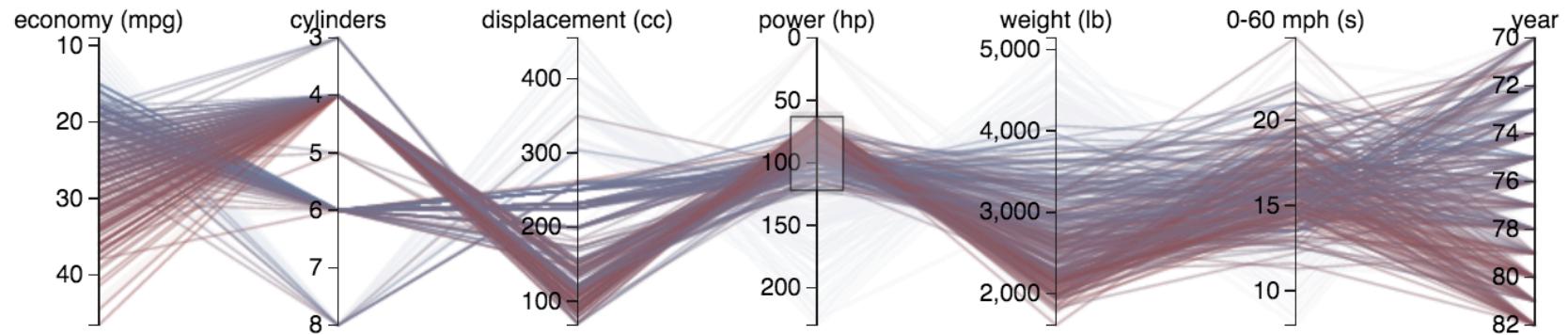
"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.



— *Iris setosa*
— *Iris versicolor*
— *Iris virginica*

Edgar Anderson's *Iris* data set
parallel coordinates

<http://mbostock.github.com/d3/talk/20111116/iris-parallel.html>



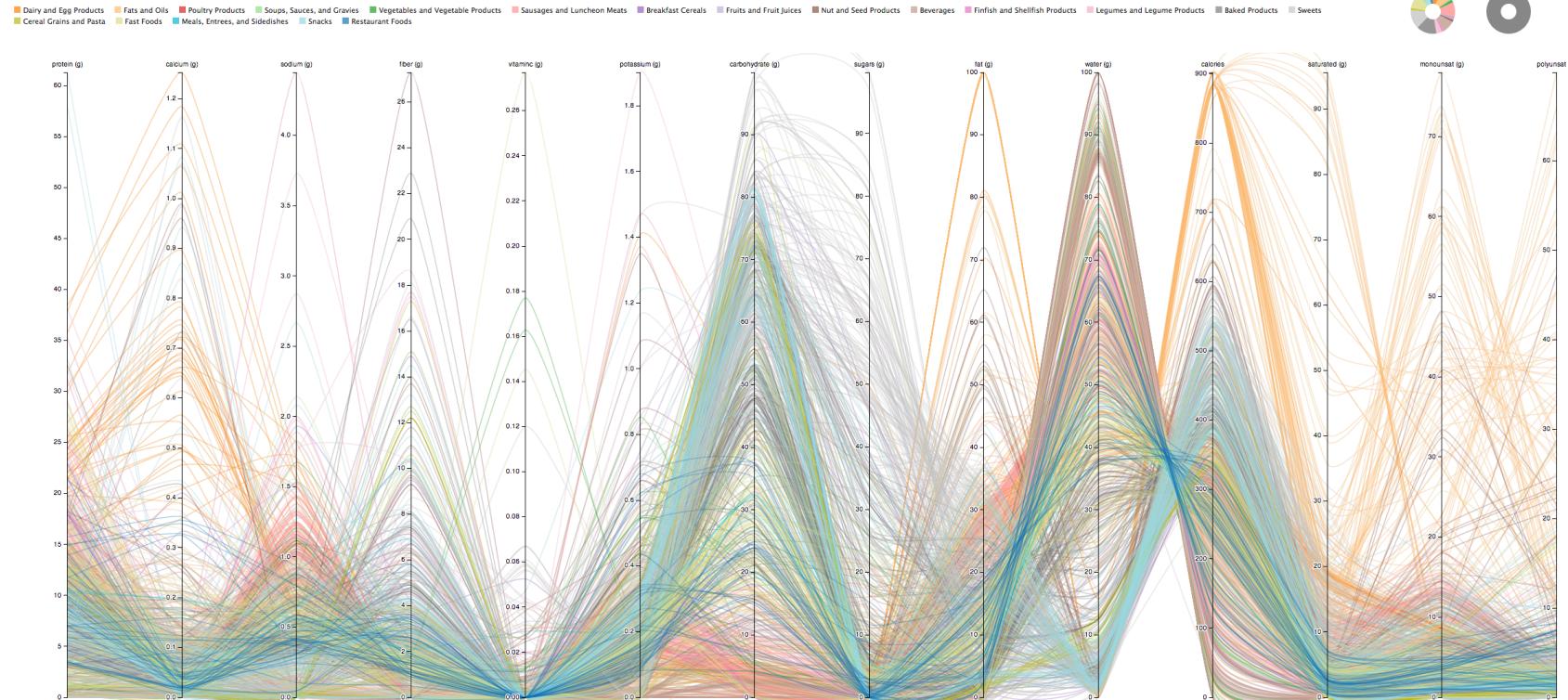
<https://syntagmatic.github.io/parallel-coordinates/examples/brushing.html>

Nutrient Contents – Parallel Coordinates

An interactive visualization of the [USDA Nutrient Database](#). For information on parallel coordinates, read [this tutorial](#).

Hide Ticks | Dark | Shadows | Opacity: 34%

Per 100g of Food



<http://exposedata.com/parallel/>

Visualizing the iOS App Store

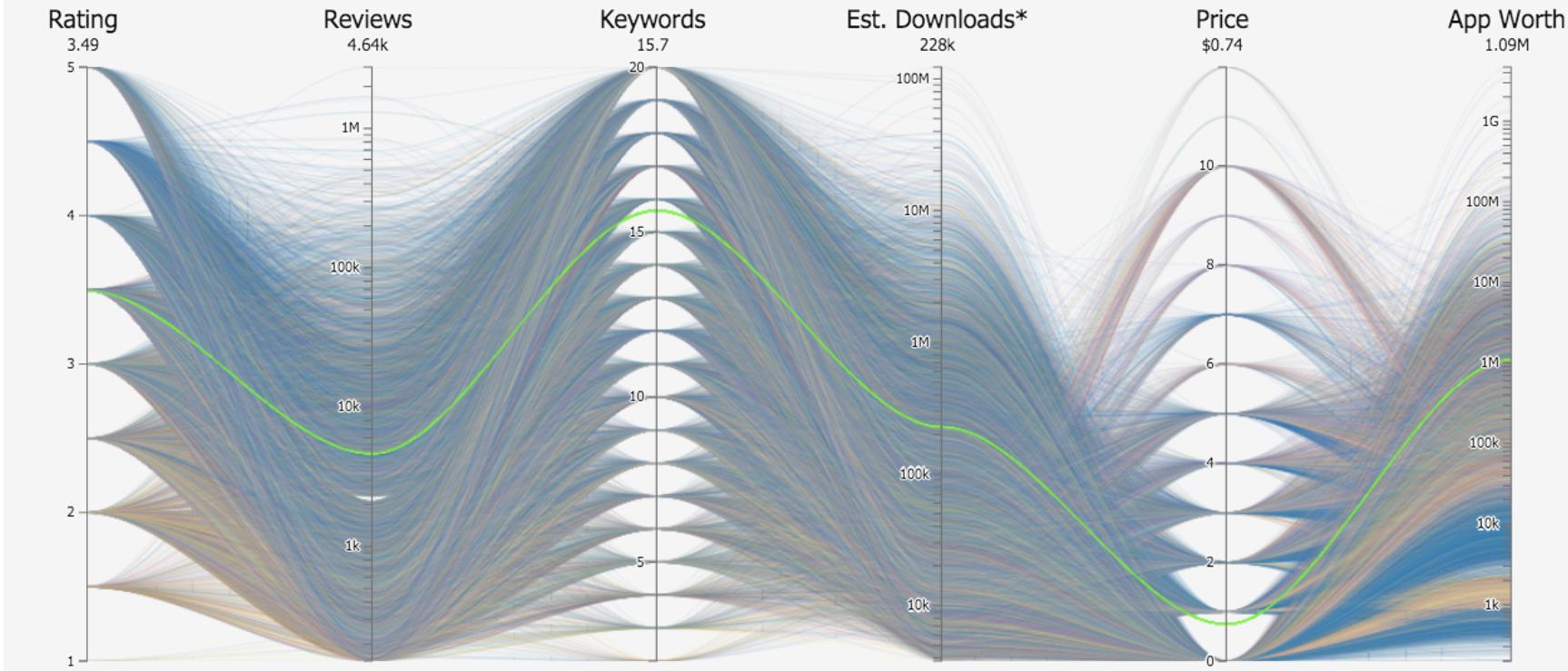
via Parallel Coordinates

Improve the search visibility and downloads of your App with [Sensor Tower](#)

[Follow @SensorTower](#)

Dark Theme

Apps: 42996/42996



<https://sensortower.com/visualizing-the-ios-app-store>

Discussion

- Long startup time
- Almost always requires interactivity
 - Choose column to color by
 - Choose how to sort
 - Highlighting (brushing)
 - Clustering
- Very high density and data ink ratio, low lie factor

OTHER TECHNIQUES

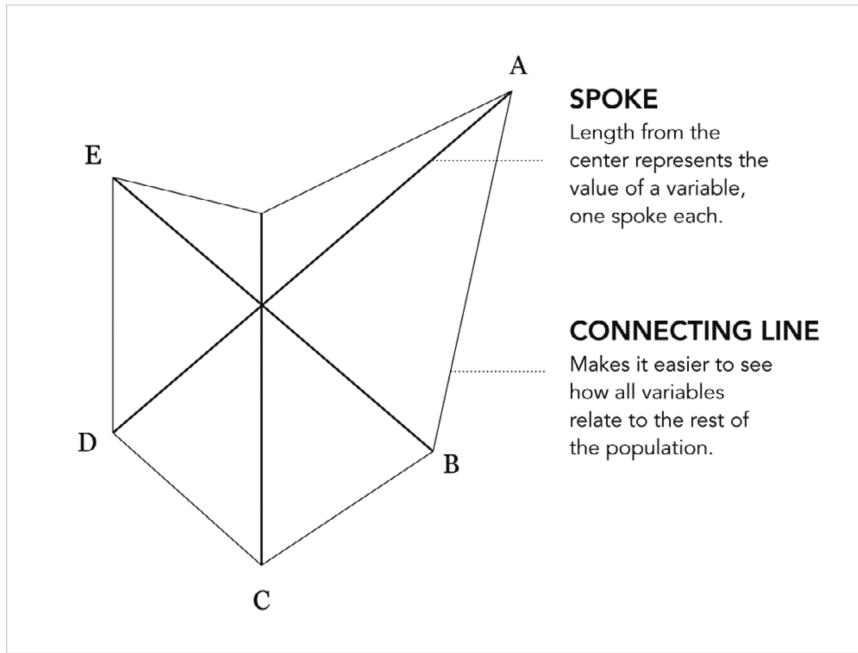


FIGURE 7-14 Star chart framework

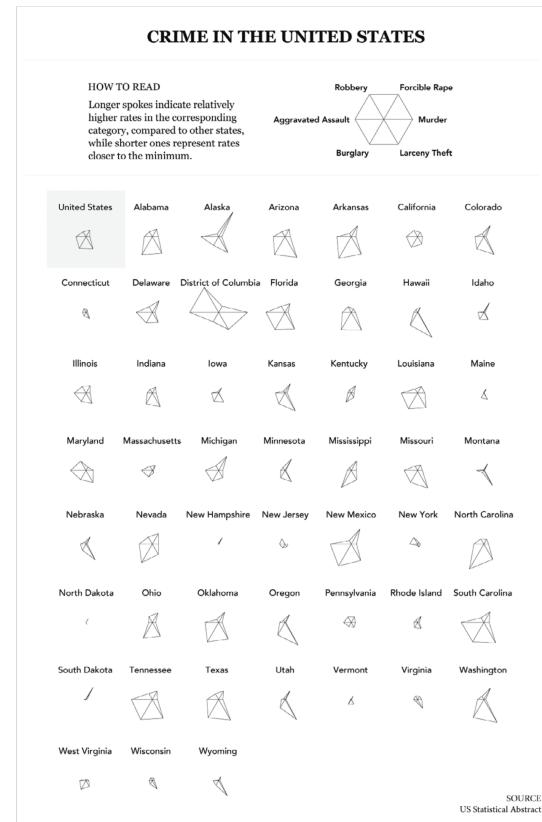
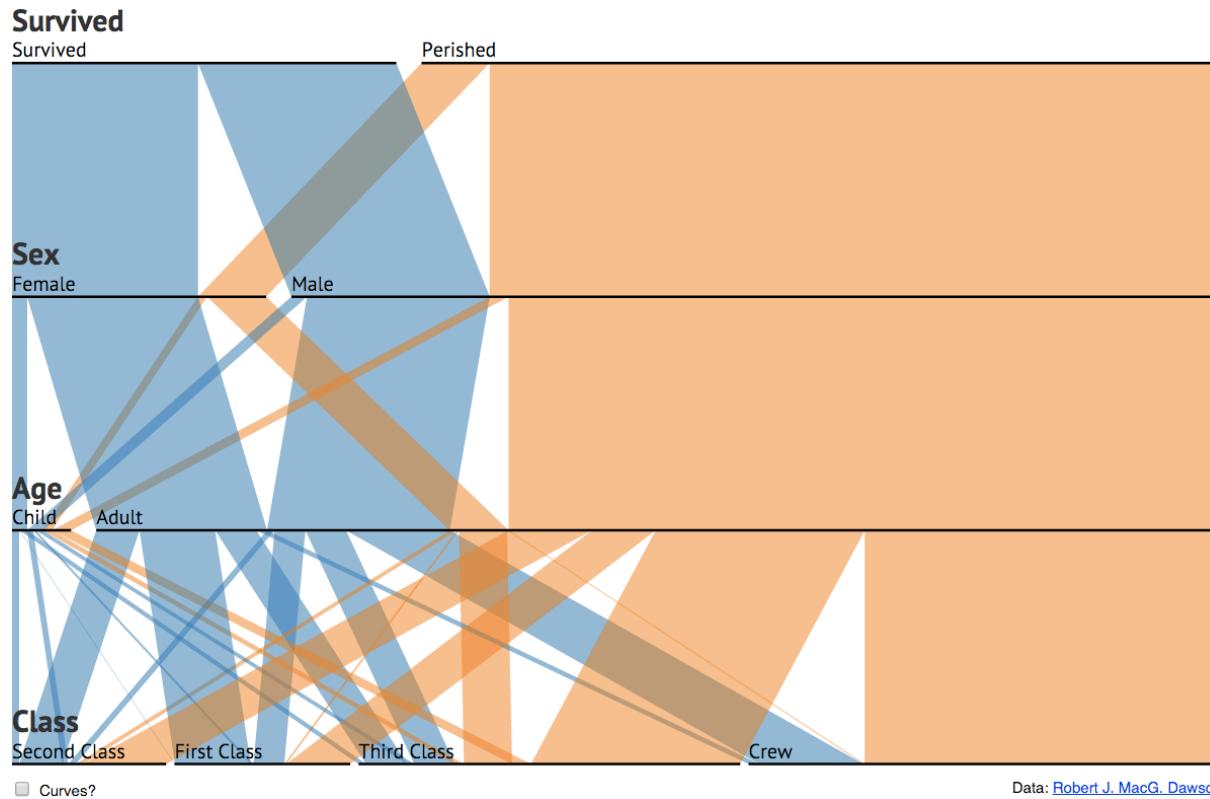


FIGURE 7-19 Series of star charts showing crime by state

"Visualize This: The FlowingData Guide to Design, Visualization, and Statistics" by Nathan Yau, 2011.

Titanic Survivors

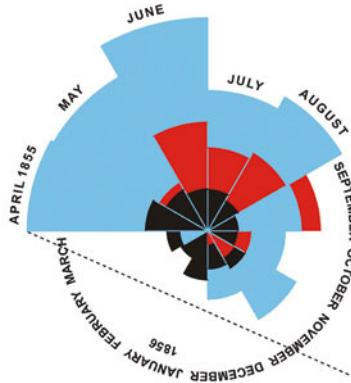


Data: [Robert J. MacG. Dawson](#).

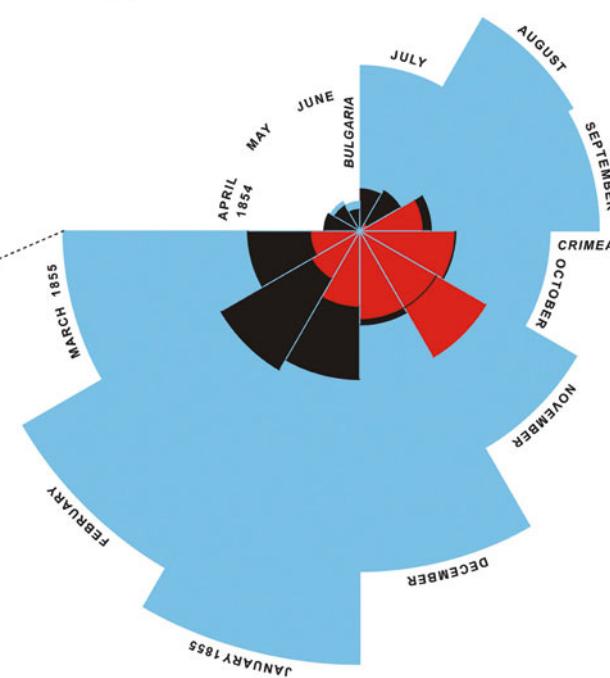
<http://eagereyes.org/parallel-sets> and <https://www.jasondavies.com/parallel-sets/>

**DIAGRAM OF THE CAUSES OF MORTALITY
IN THE ARMY IN THE EAST.**

2.
APRIL 1855 TO MARCH 1856.



1.
APRIL 1854 TO MARCH 1855.



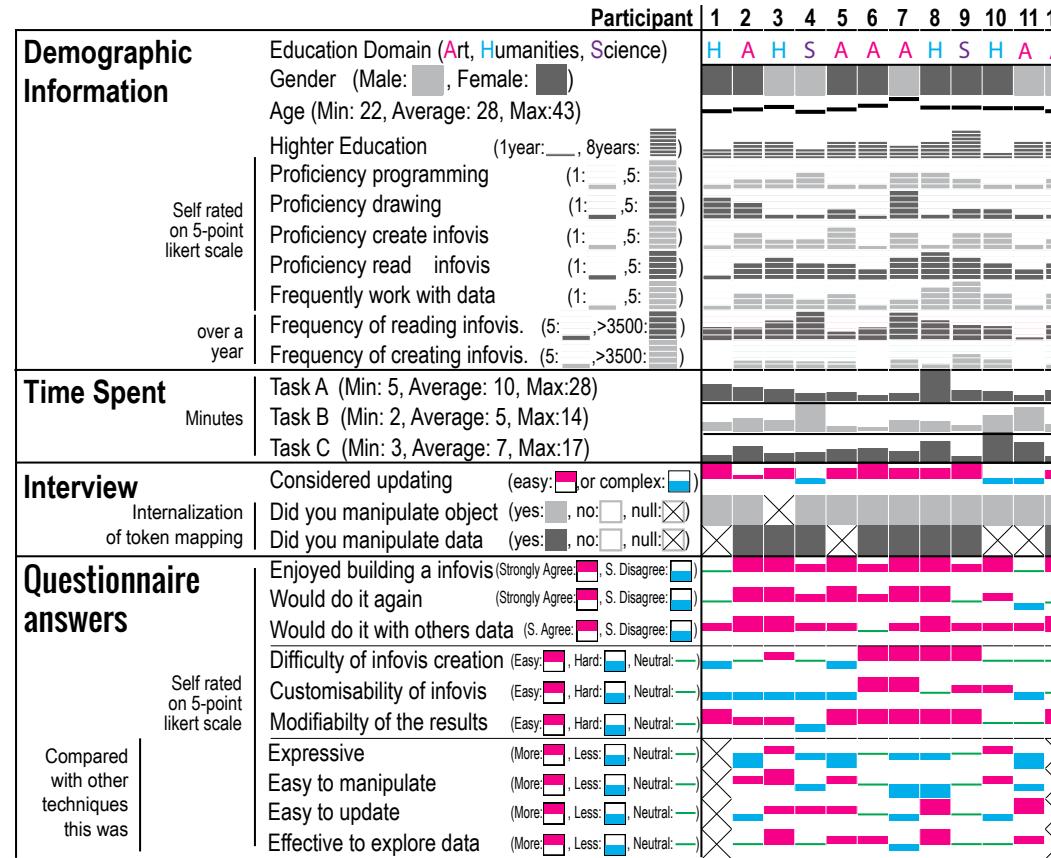
*The Areas of the blue, red, & black wedges are each measured from
the centre as the common vertex*

*The blue wedges measured from the centre of the circle represent area
for area the deaths from Preventive or Mitigable Zymotic Diseases, the
red wedges measured from the centre the deaths from wounds, & the
black wedges measured from the centre the deaths from all other causes
The black line across the red triangle in Nov' 1854 marks the boundary
of the deaths from all other causes during the month*

*In October 1854, & April 1855, the black area coincides with the red,
in January & February 1856, the blue coincides with the black*

*The entire areas may be compared by following the blue, the red & the
black lines enclosing them. ©hugh-small.co.uk*

<http://dd.dynamicdiagrams.com/2008/01/nightingales-rose/>



<http://bertifier.com/> and <http://constructive.gforge.inria.fr/#!questionnaires.md>

CONCLUSION

References

- Nathan Yau, **Visualize This: The FlowingData Guide to Design, Visualization, and Statistics**, Wiley Publishing, 2011.
- Mike Bostock, **Data Driven Documents** (D3.js)
 - <https://github.com/mbostock/d3/wiki/Gallery>
 - <http://bl.ocks.org/mbostock>
- Other references sourced on slides





CHANGE THE WORLD FROM HERE