

The Mathematical Foundation of Analytic Visualizations

Leland Wilkinson

SYSTAT

University of Illinois at Chicago (Computer Science)

Northwestern University (Statistics)

FODAVA Distinguished Lecture

Georgia Tech

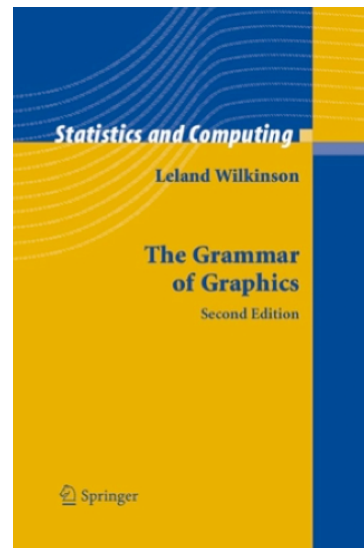
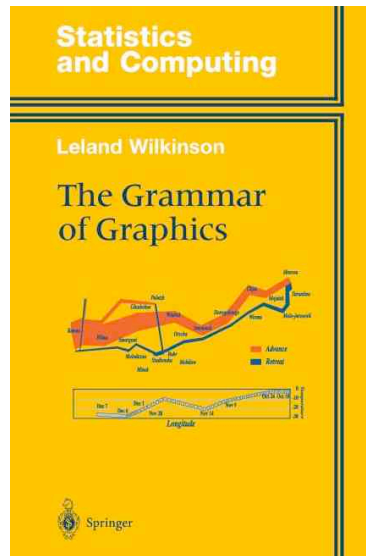
April 2, 2010

Outline

- A formal graphics language for analytics
 - An intelligent analytic system requires a language
 - The Grammar of Graphics (GoG)
- Aspects of a formal graphics language
 - Simple
 - Expressive
 - Coherent
 - Meaningful
- Intelligent graphics systems based on GoG

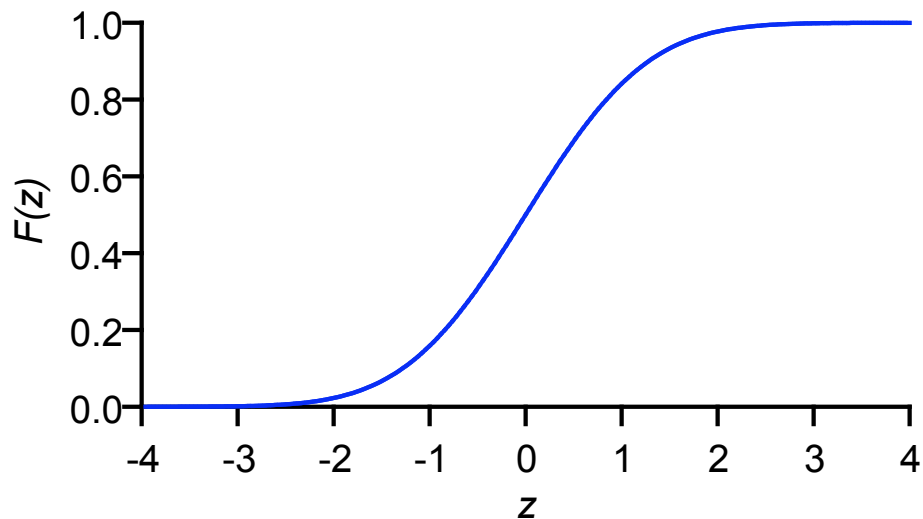
The Grammar of Graphics (1999, 2005)

- Show programmers how to design and implement statistical graphics
- Reveal the mathematical foundation of statistical graphics.



The GoG Rule

- A statistical graphic is a representation of the graph of a function.
- The graph of a function is a subset of the product of its domain and codomain.
- The graphic representing $F(z) = \Phi(z)$ here is blue.
- The rest is annotation.



What is a Graphic?

Graph

$$G = \{(x, f(x)) : x \in \mathbb{R} \text{ and } f(x) = e^{-x^2}\}$$

Frame

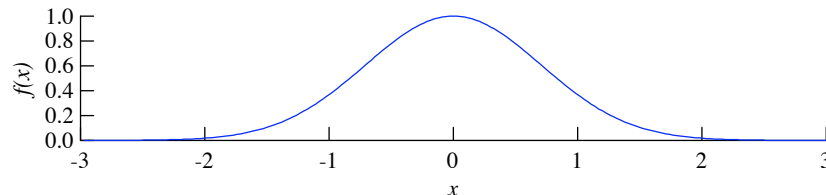
$$F = [-3, 3] \times [0, 1]$$

Aesthetic

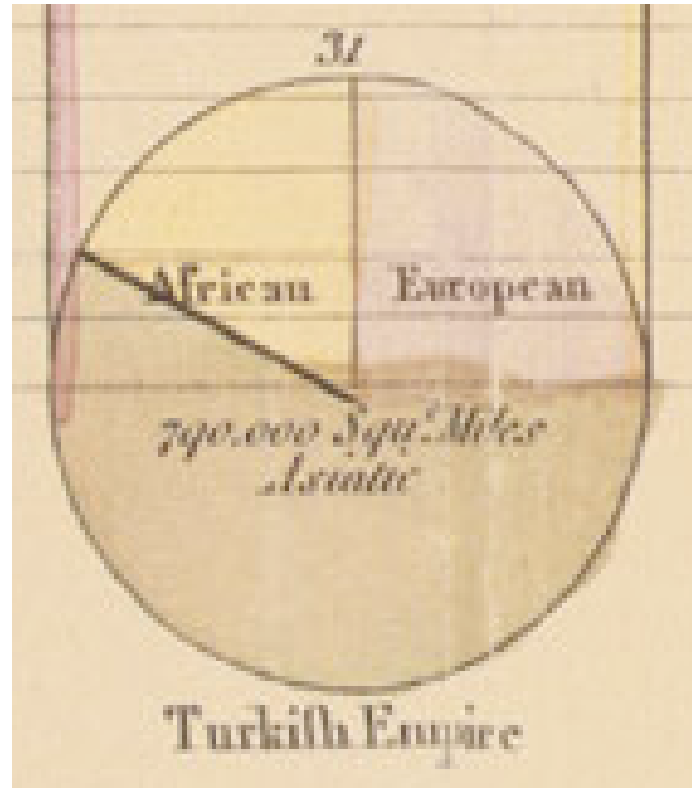
$$A : x \mapsto x_{\text{position}}, f(x) \mapsto y_{\text{position}}$$

Graphic

$$G_A = A(F \cap G)$$

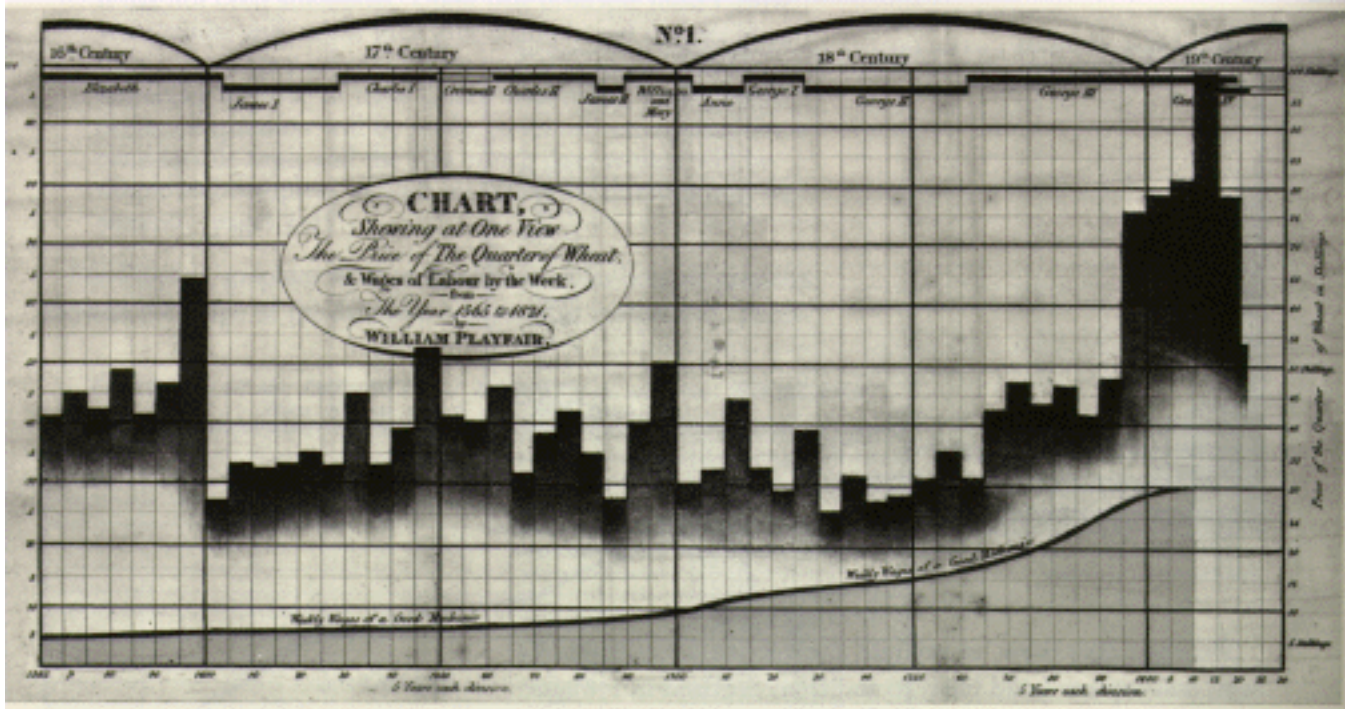


The same rule applies to this graphic



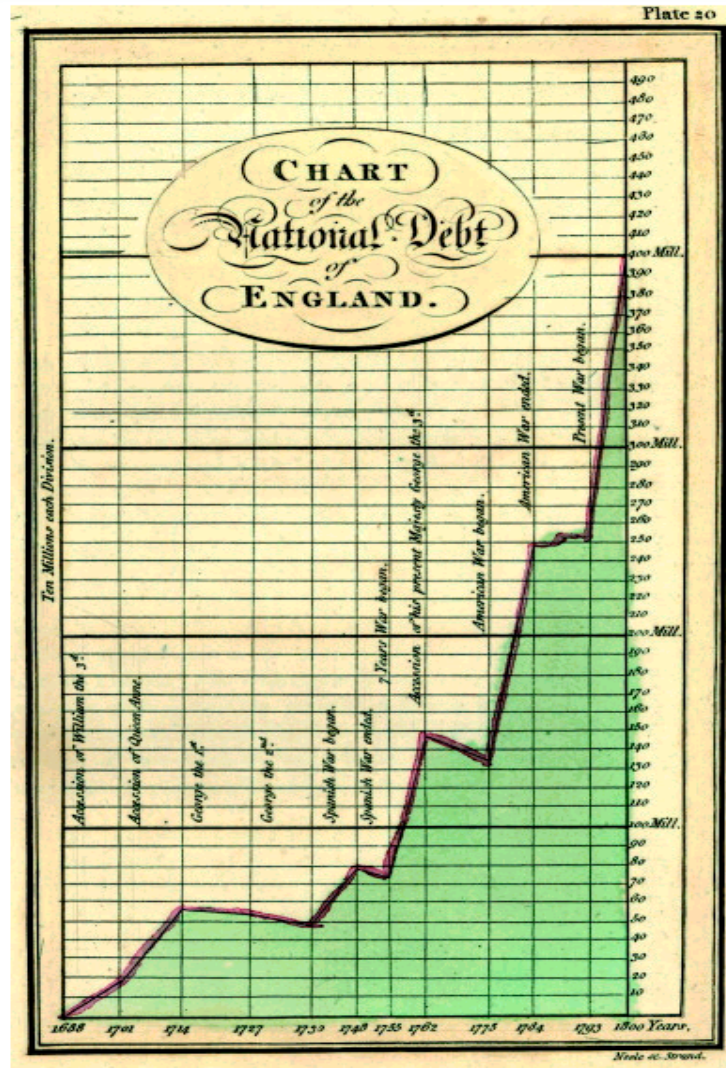
Pie chart, William Playfair

And to this



Composite line and bar chart, William Playfair

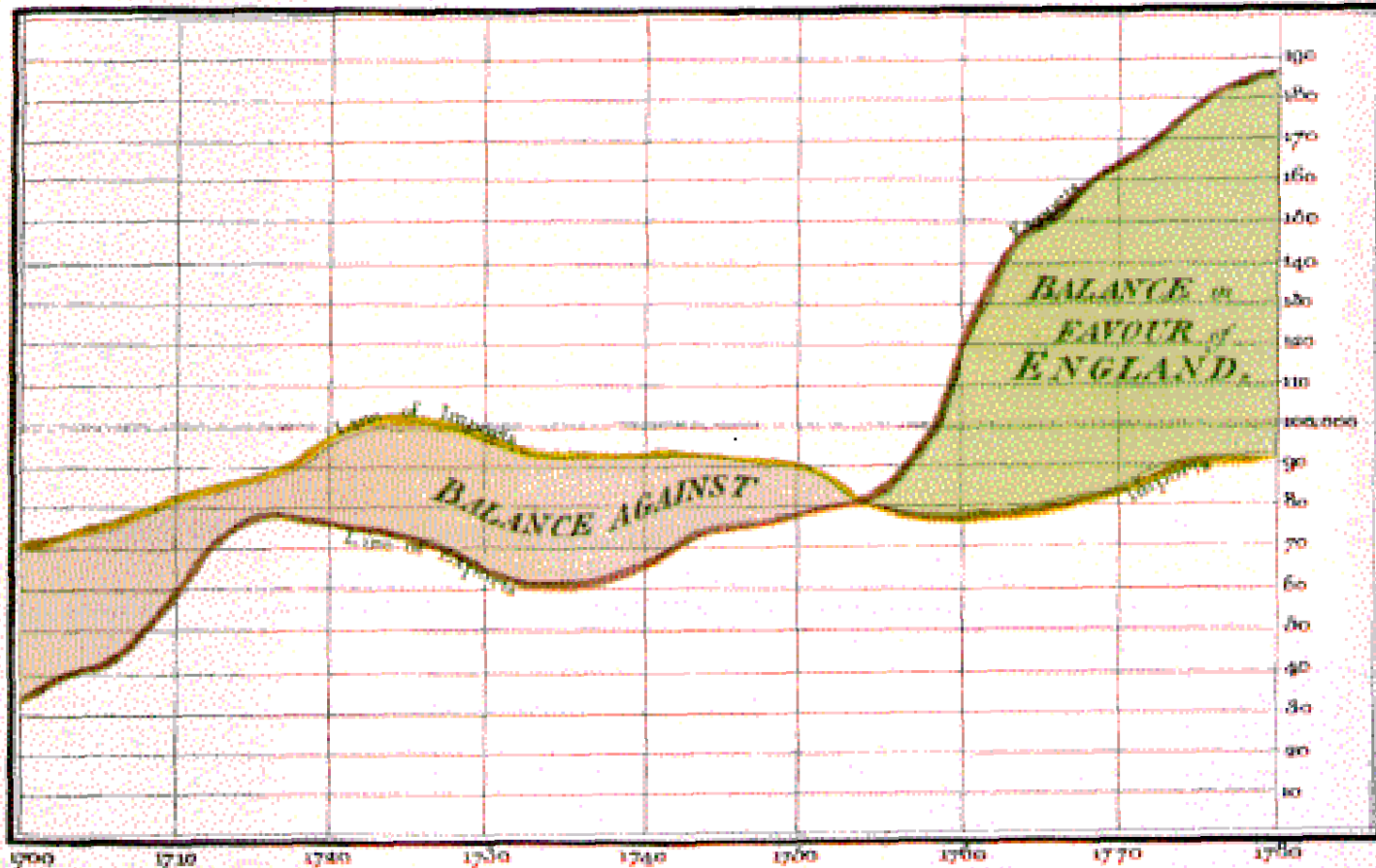
And this



Area chart, William Playfair

And this

Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780



The Bottom line is divided into Years, the Right hand line into £10,000 each.

Published in the first volume, 1791, by Wm Playfair.

Revised and corrected, 1844, by Wm Playfair.

Differenced area chart, William Playfair

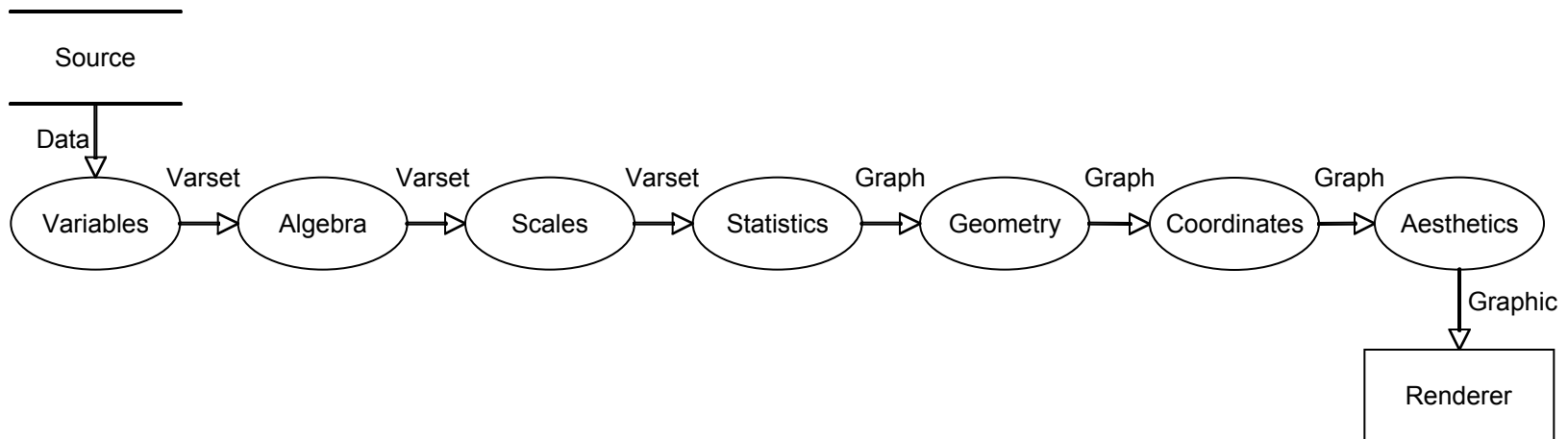
And this



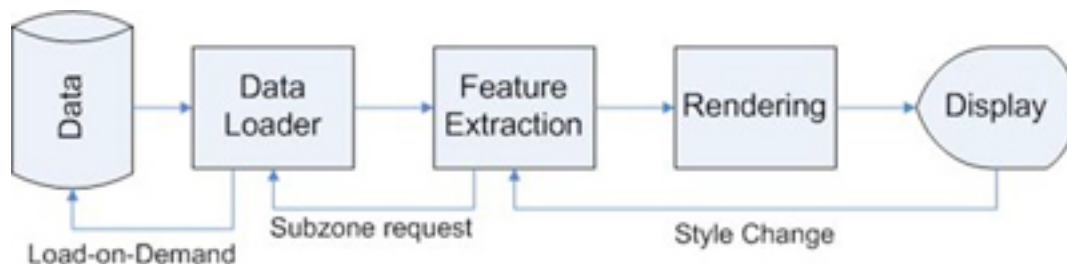
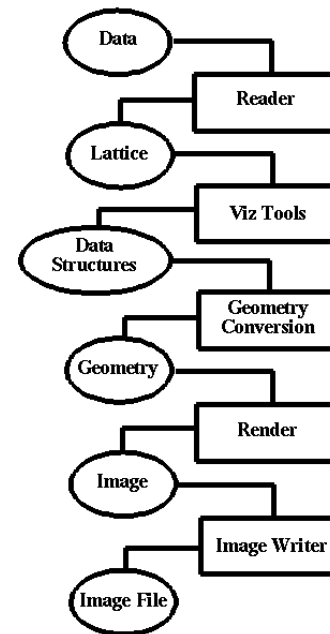
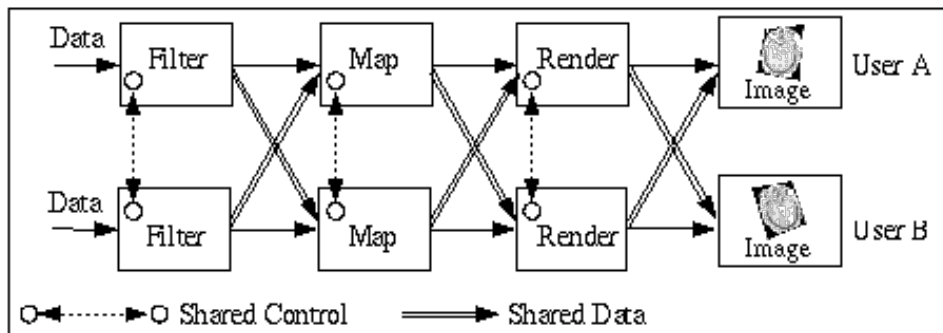
Soletto Map ca. 500 BCE, southern region of Italy's heel, discovered in dig supervised by Thierry van Compernelle, Montpellier University

The GoG Function Chain

- Each ellipse is a class.
- Each class contains member functions.
 - functions are interchangeable within a class
- The chain is a total order.
 - changing this order produces a meaningless graphic
- The chain is invertible.

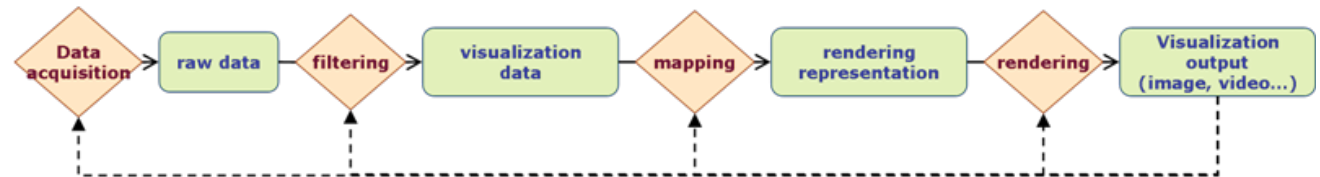
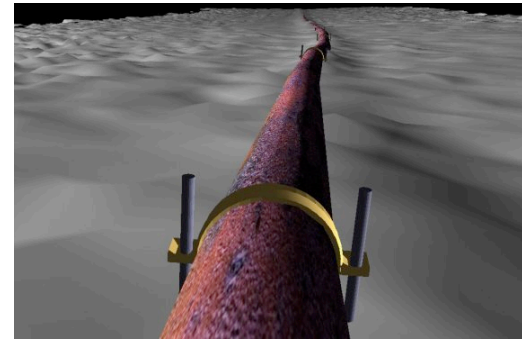
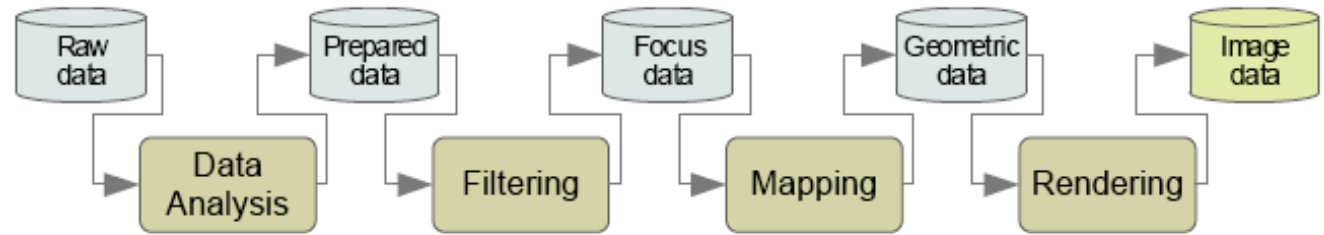
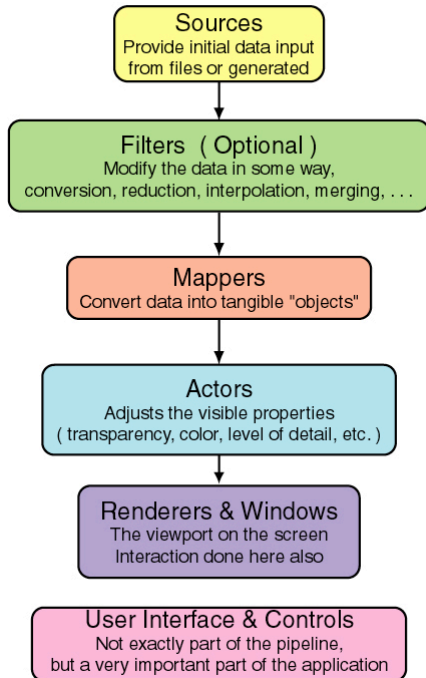


The GoG Chain is NOT a dataflow

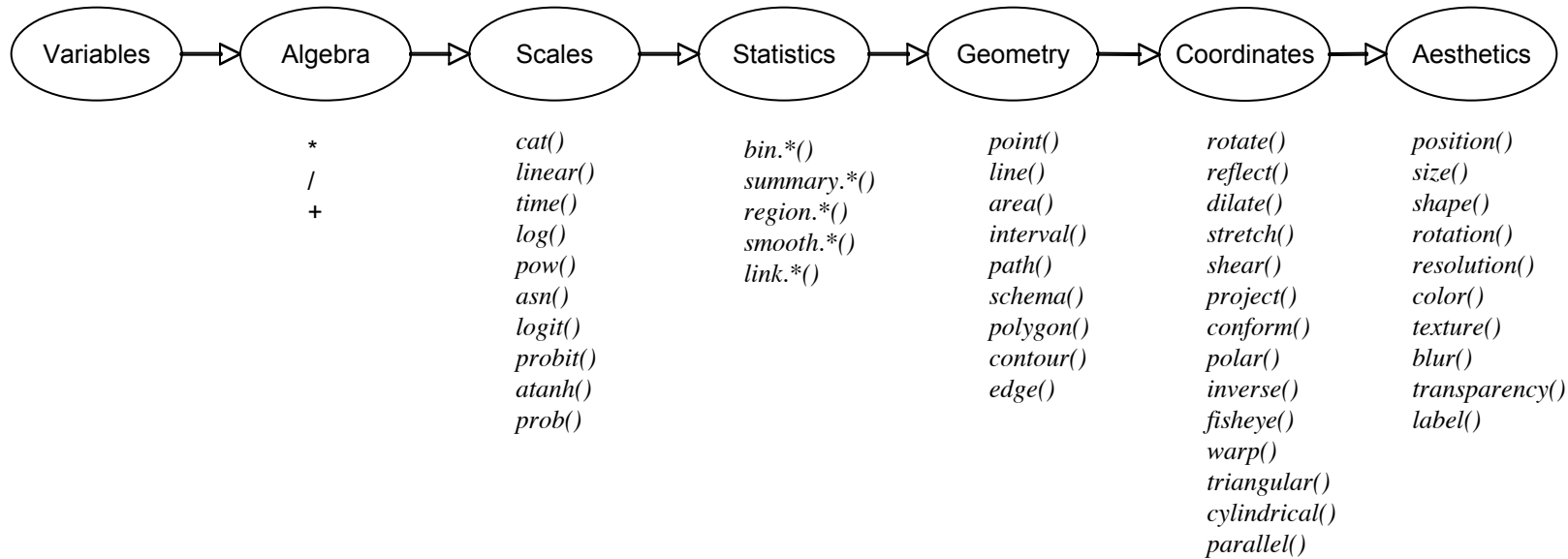


The GoG Chain is NOT a visualization pipeline

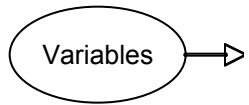
VTK Visualization Pipeline



GoG Member Functions



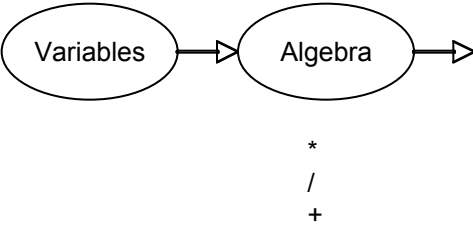
Variables



Variables

- Variables is a class that builds a data view.
- Variables receives a Dataset and outputs a Varset.
- A Dataset is a set of data.
- A Varset is a set of variables.
- A Variable is a function from a set of objects O to a set of values V (a many-to-one mapping).

Algebra



Algebra

- Algebra has three operators

- cross (*)

$$\begin{array}{|c|} \hline x \\ \hline y \\ \hline z \\ \hline \end{array} * \begin{array}{|c|} \hline a \\ \hline a \\ \hline b \\ \hline \end{array} = \begin{array}{|c|c|} \hline x & a \\ \hline y & a \\ \hline z & b \\ \hline \end{array}$$

- nest (/)

$$\begin{array}{|c|} \hline x \\ \hline y \\ \hline z \\ \hline \end{array} / \begin{array}{|c|} \hline a \\ \hline a \\ \hline b \\ \hline \end{array} = \begin{array}{|c|c|} \hline x & a \\ \hline y & a \\ \hline z & b \\ \hline \end{array}$$

- blend (+)

$$\begin{array}{|c|} \hline x \\ \hline y \\ \hline z \\ \hline \end{array} + \begin{array}{|c|} \hline a \\ \hline a \\ \hline b \\ \hline \end{array} = \begin{array}{|c|} \hline x \\ \hline y \\ \hline z \\ \hline a \\ \hline a \\ \hline b \\ \hline \end{array}$$

Algebra

Males	Impala
	Mustang
	Jaguar
Females	Impala
	Mustang
	Jaguar

animal * gender

- Male Jaguar and Female Jaguar are comparable

Algebra

Cars	Impala
	Mustang
	Jaguar
Animals	Impala
	Mustang
	Jaguar
Teams	Impala
	Mustang
	Jaguar

animal / class

- Jaguar Car and Jaguar Animal are not comparable

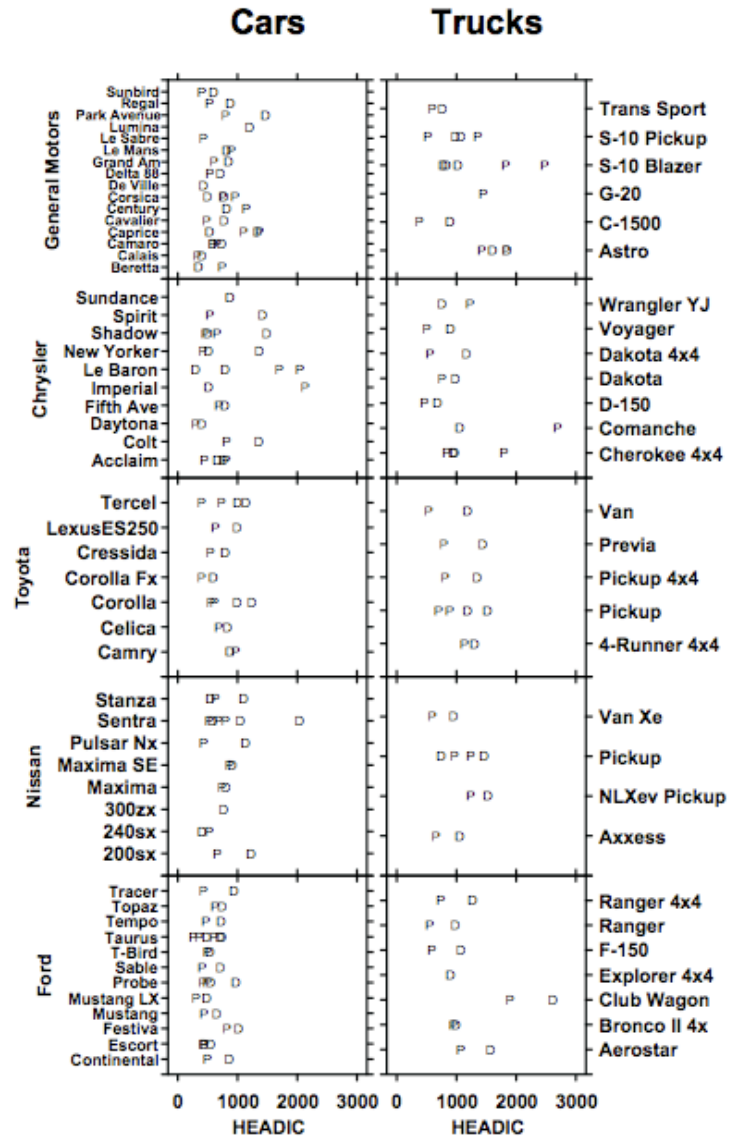
Algebra

Cow
Chicken
Pig
Horse
Snake
Impala
Lion
Giraffe
Jaguar

barnyard + zoo

- Barnyard and Zoo must be from a common class or measurement scale for blend to be legal.

Algebra



Algebra

The GoG algebra expression for this graphic is

$$H * T / (M * V) * O$$

The GoG algebra expression corresponds to the design model

$$H = C + M + V + O + T(MV) + MV + MO + VO + OT(MV) + MVO$$

where the symbols are

H: Head Injury Index

C: constant term (grand mean)

M: Manufacturer

V: Vehicle (car or truck)

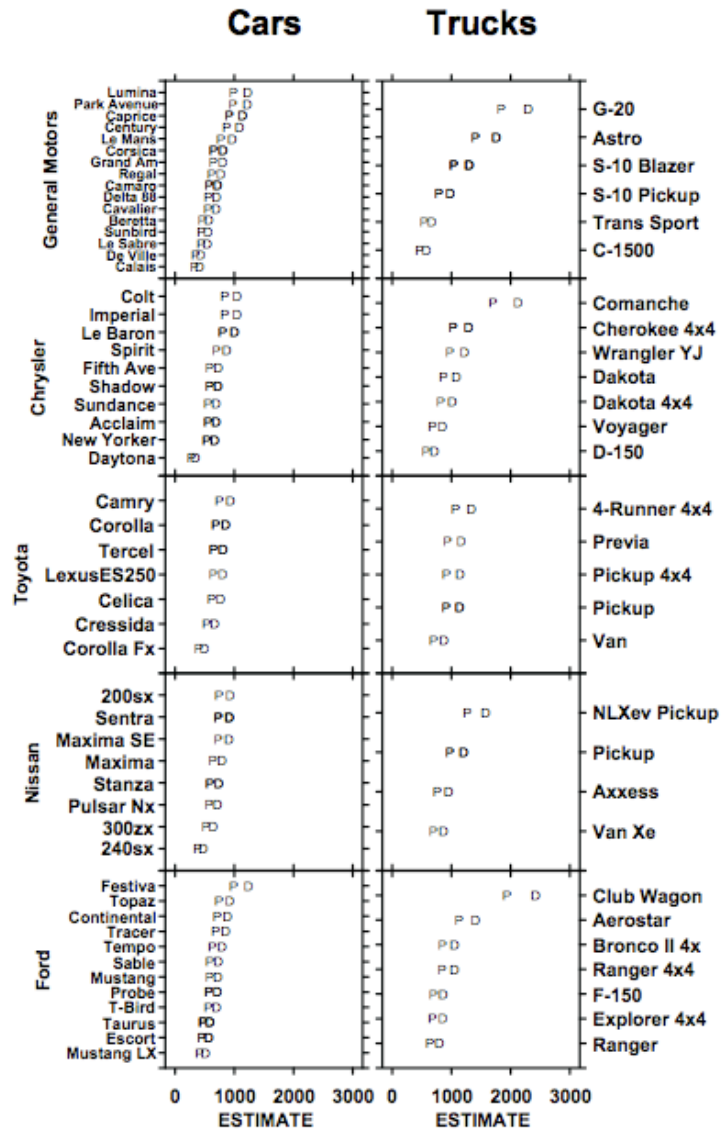
O: Occupant (driver or passenger)

T: Model

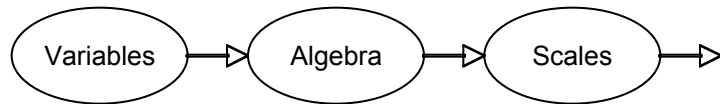
The smallest plausible subset model is:

$$H = C + V + O + T(MV)$$

Algebra



Scales



*
/
+

cat()
linear()
time()
log()
pow()
asn()
logit()
probit()
atanh()
prob()

Scales

- Scales map Varsets to Dimensions.
- A Dimension is a dimension of \mathbf{R}^n .
- Scales use one-to-one mappings within Frames
 - `identity()`, `log()`, `permutation()`, ...
- More than one Varset can map to a Dimension.
- But we need rules for a meaningful mapping.
- Representational measurement is not sufficient.
 - nominal, ordinal, interval, ratio (Stevens)
 - too general for blend (+) operator
 - we can't blend area with weight
 - we can't blend speed with acceleration
 - we can't blend a density and a distribution function

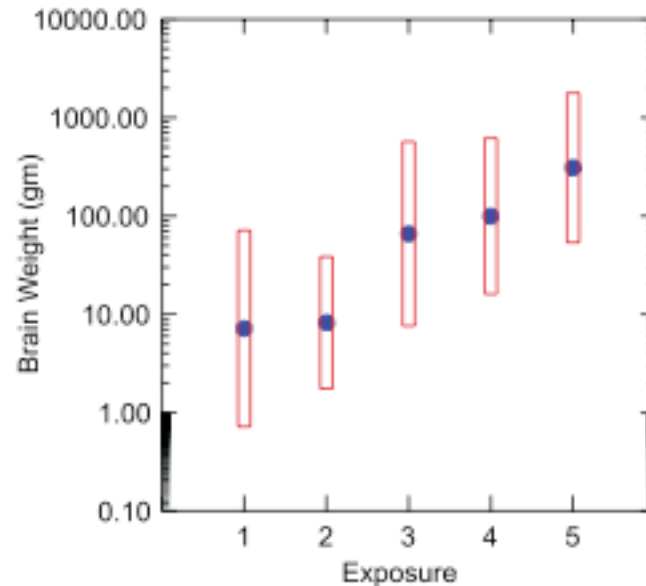
Scales

- We need measurement units.
 - more restrictive than axiomatic measurement scales

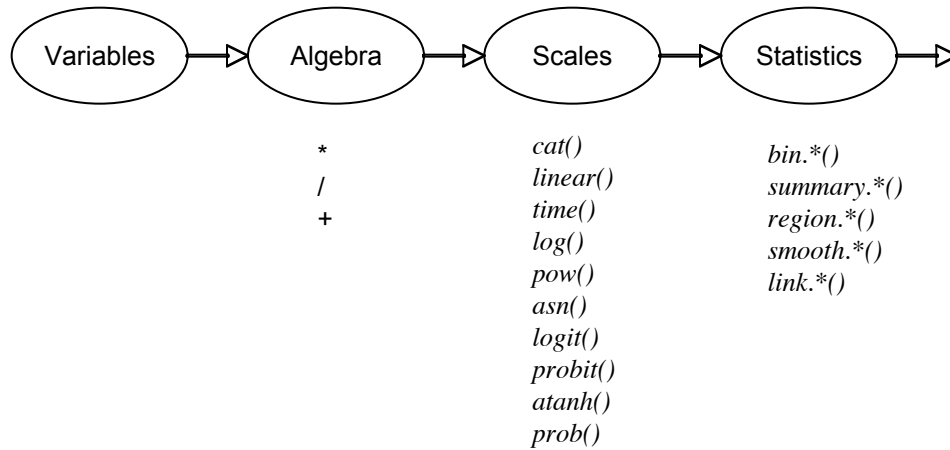
<i>Length</i>	<i>Mass</i>	<i>Time</i>	<i>Current</i>	<i>Temperature</i>
<i>meter</i> <i>point</i> <i>pica</i> <i>inch</i> <i>foot</i> <i>yard</i> <i>mile</i> <i>furlong</i> <i>fathom</i>	<i>kilogram</i> <i>gram</i> <i>grain</i> <i>slug</i> <i>carat</i>	<i>second</i> <i>minute</i> <i>hour</i> <i>day</i> <i>week</i> <i>month</i> <i>quarter</i> <i>year</i> <i>century</i>	<i>amp</i>	<i>kelvin</i> <i>rankine</i> <i>celsius</i> <i>fahrenheit</i>
<i>Substance</i>	<i>Luminous Intensity</i>	<i>Angle</i>	<i>Count</i>	<i>Currency</i>
<i>mole</i>	<i>candela</i>	<i>radian</i> <i>degree</i> <i>minute</i> <i>second</i>	<i>unit</i> <i>dozen</i> <i>gross</i>	<i>dollar</i> <i>euro</i> <i>pound</i> <i>yen</i>

Scales

```
SCALE: cat(dim(1))  
SCALE: log(dim(2), base(10))  
ELEMENT: point(position(summary.mean(exposure*brainweight)))  
ELEMENT: interval(position(spread.sd(exposure*brainweight),  
color(color.red))
```



Statistics



Statistics

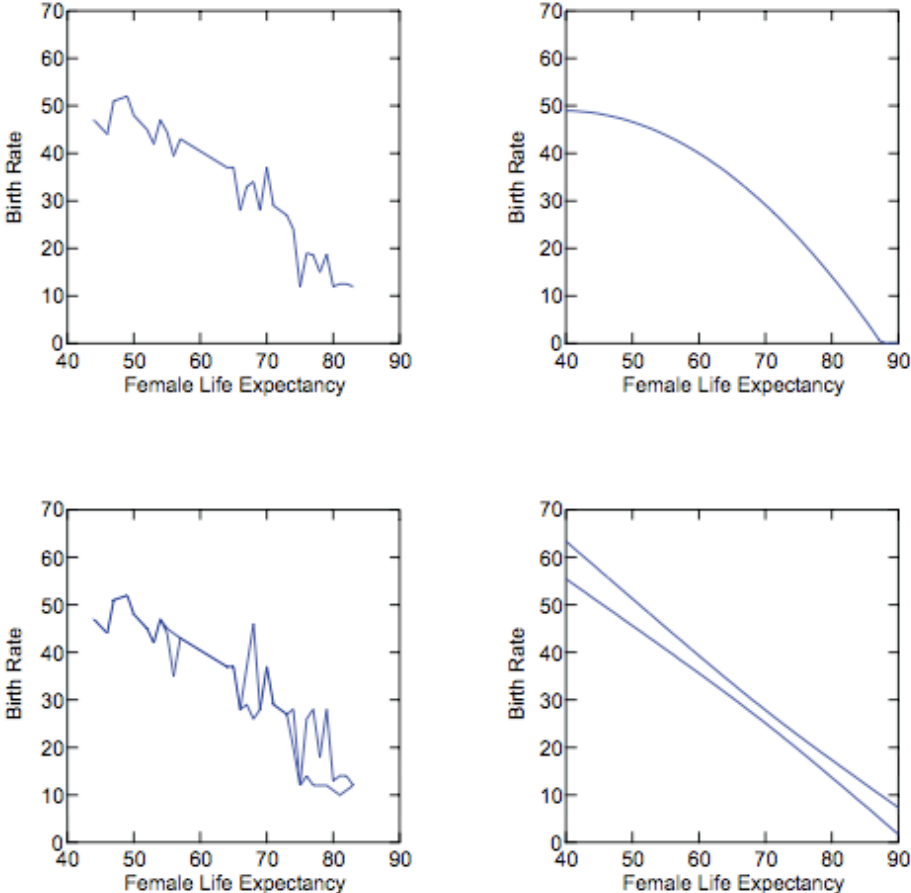


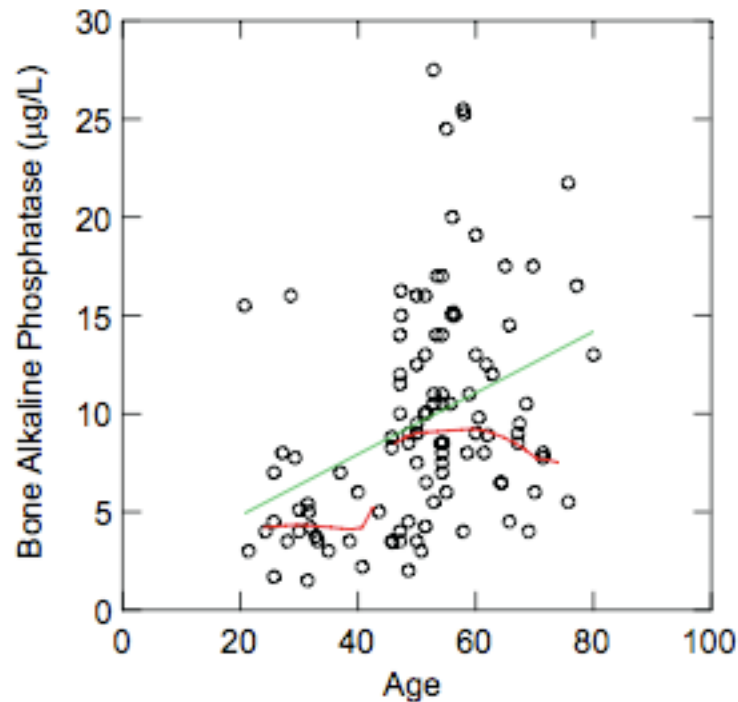
Figure 7.2 Different statistical methods, same graph type

Statistics

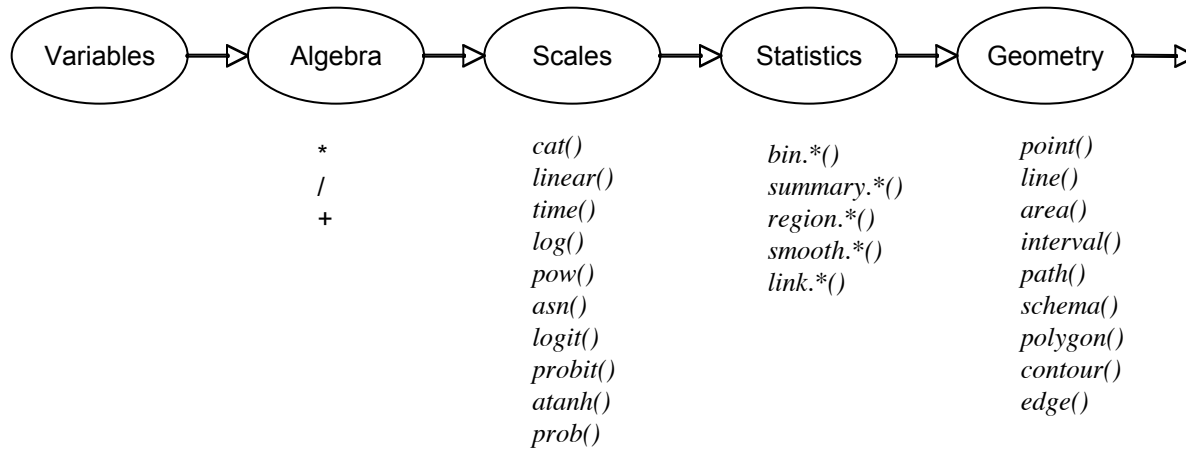
ELEMENT: `line(position(smooth.mode.epanechnikov(age*bone)),
color(color.red))`

ELEMENT: `line(position(smooth.linear(age*bone)), color(color.green))`

ELEMENT: `point(position(age*bone))`



Geometry



Geometry

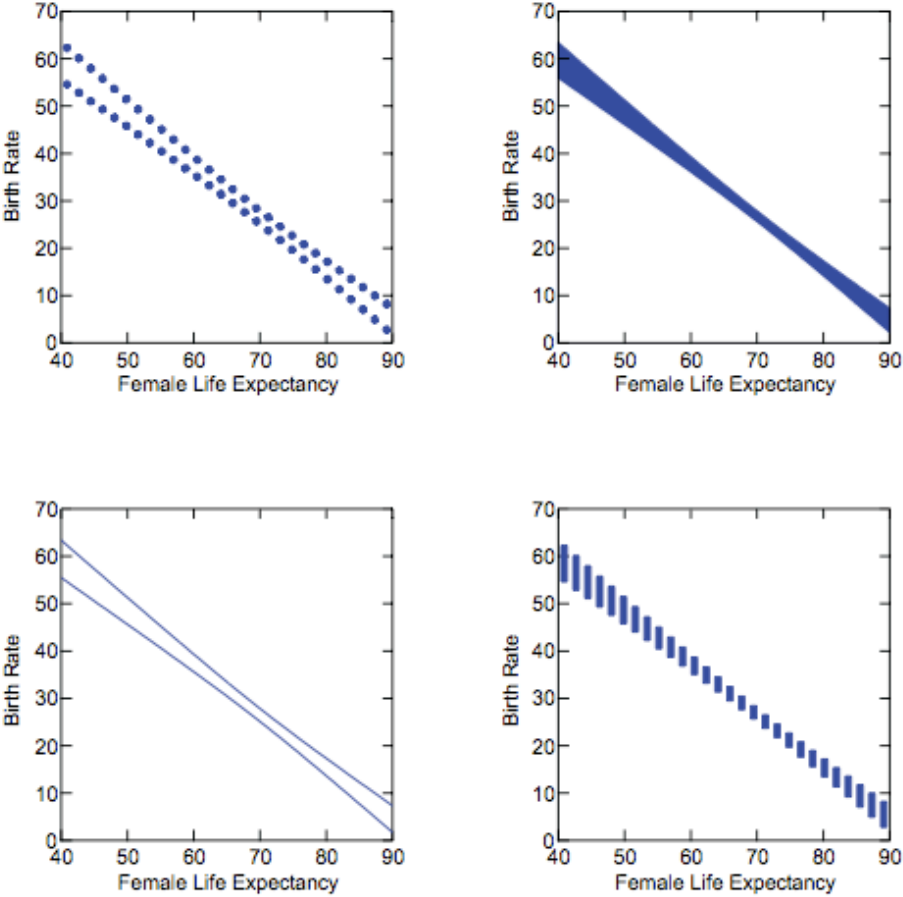


Figure 7.1 Different graph types, same statistical method

Statistics (recall)

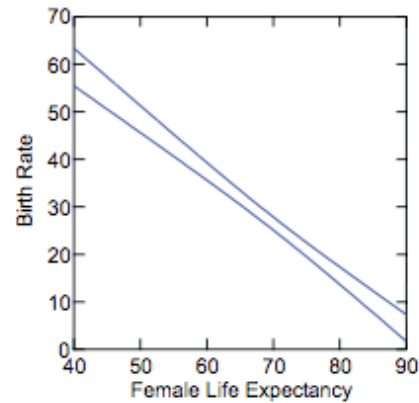
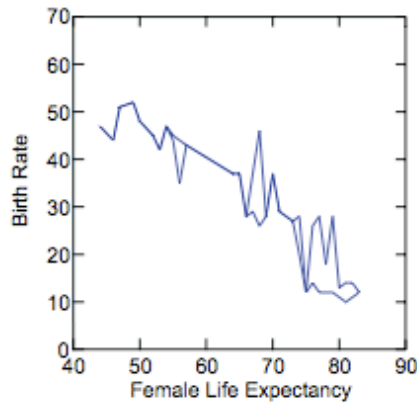
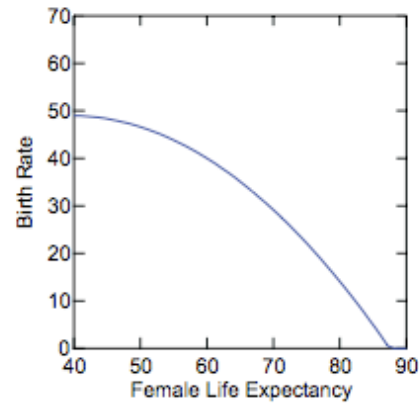
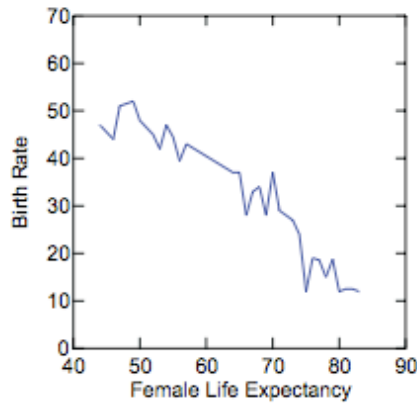
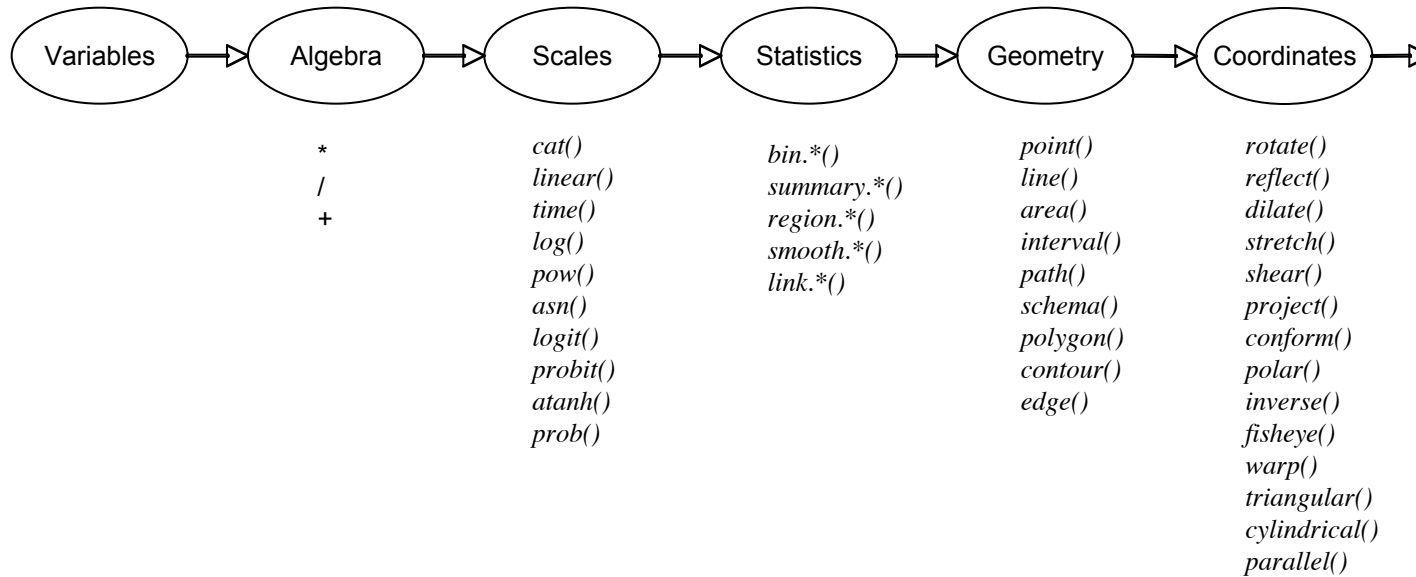


Figure 7.2 Different statistical methods, same graph type

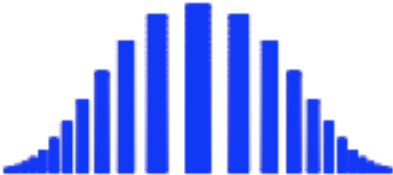
Coordinates



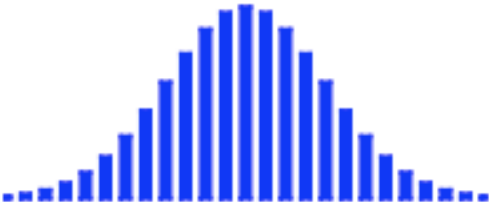
Coordinates



conformal



fisheye



rectangular



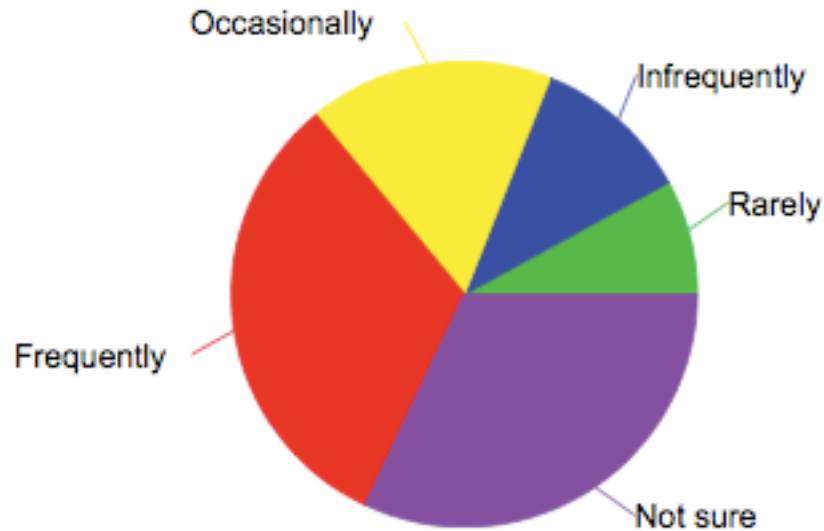
polar



spherical

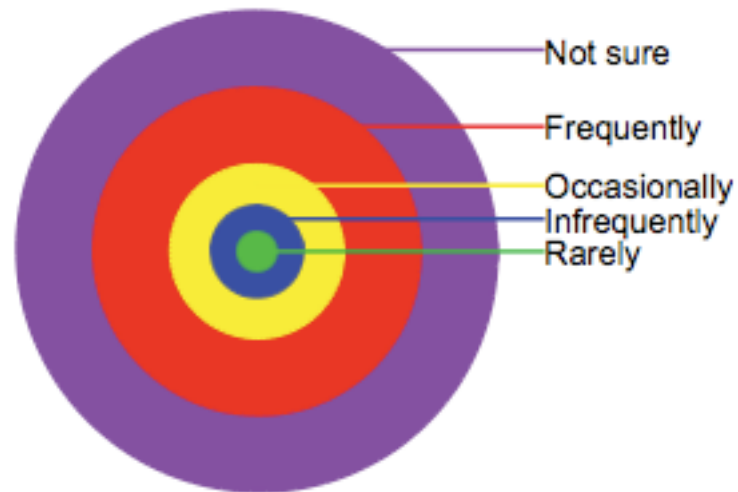
Coordinates

COORD: *polar.theta()*
ELEMENT: *interval.stack(position(summary.proportion(response)),
label(response), color(response))*



Coordinates

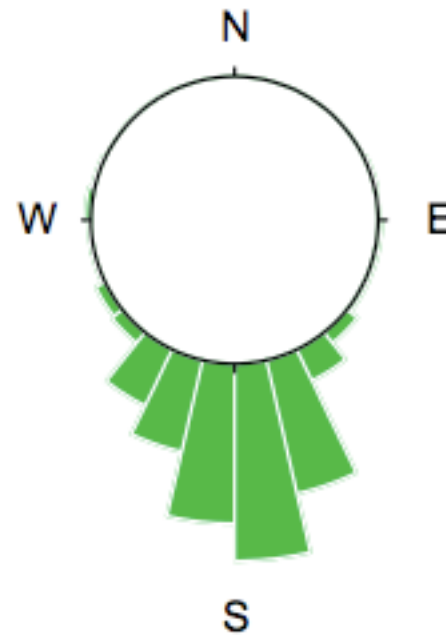
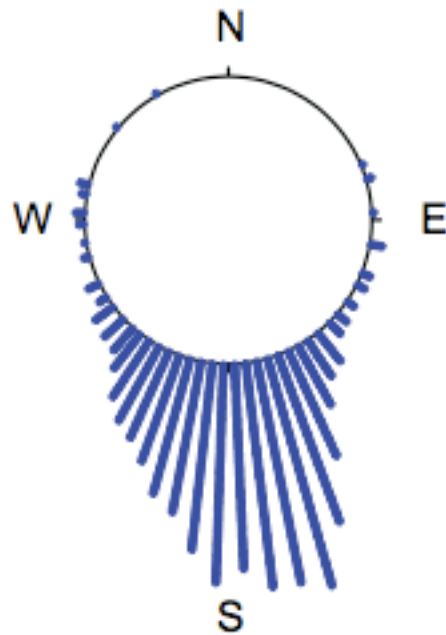
COORD: *polar.rho()*
ELEMENT: *interval.stack(position(summary.proportion(response)),
label(response), color(response))*



Coordinates

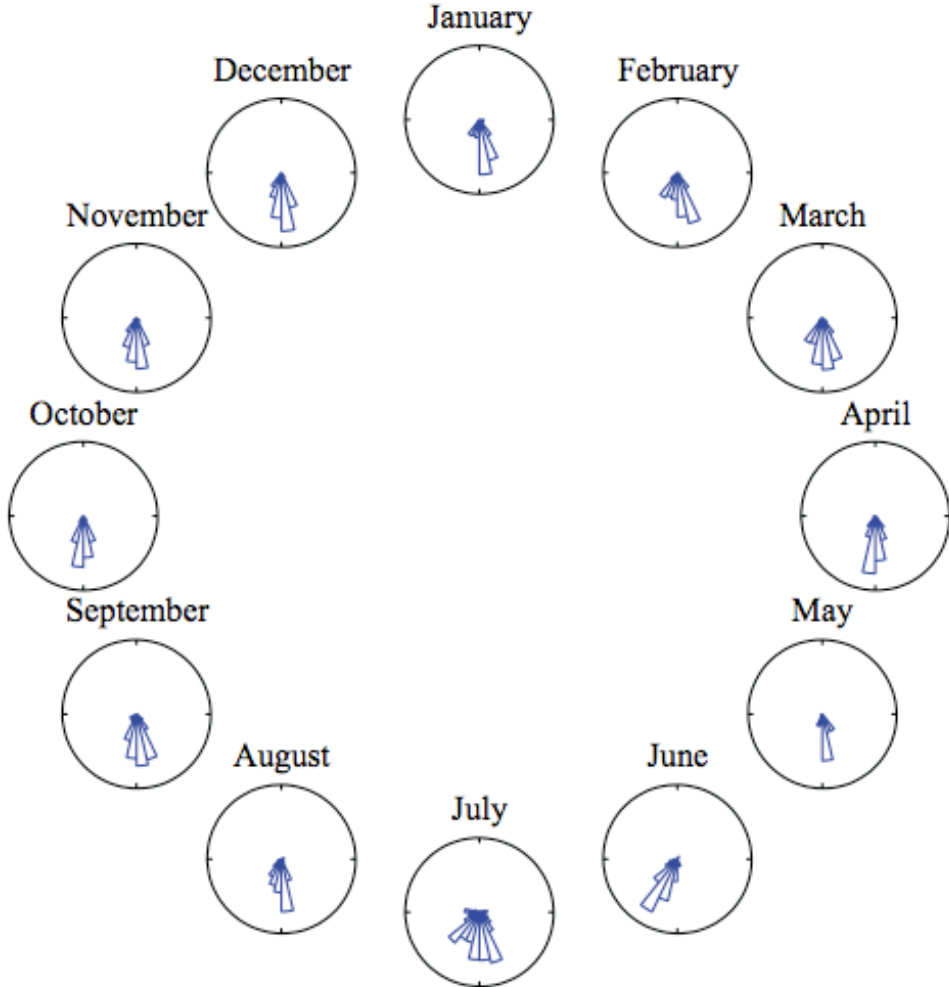
COORD: *polar.rho.plus()*
ELEMENT: *point.dodge(position(bin.dot(direction)))*

COORD: *polar.rho.plus()*
ELEMENT: *interval(position(summary.count(bin.rect(direction))))*

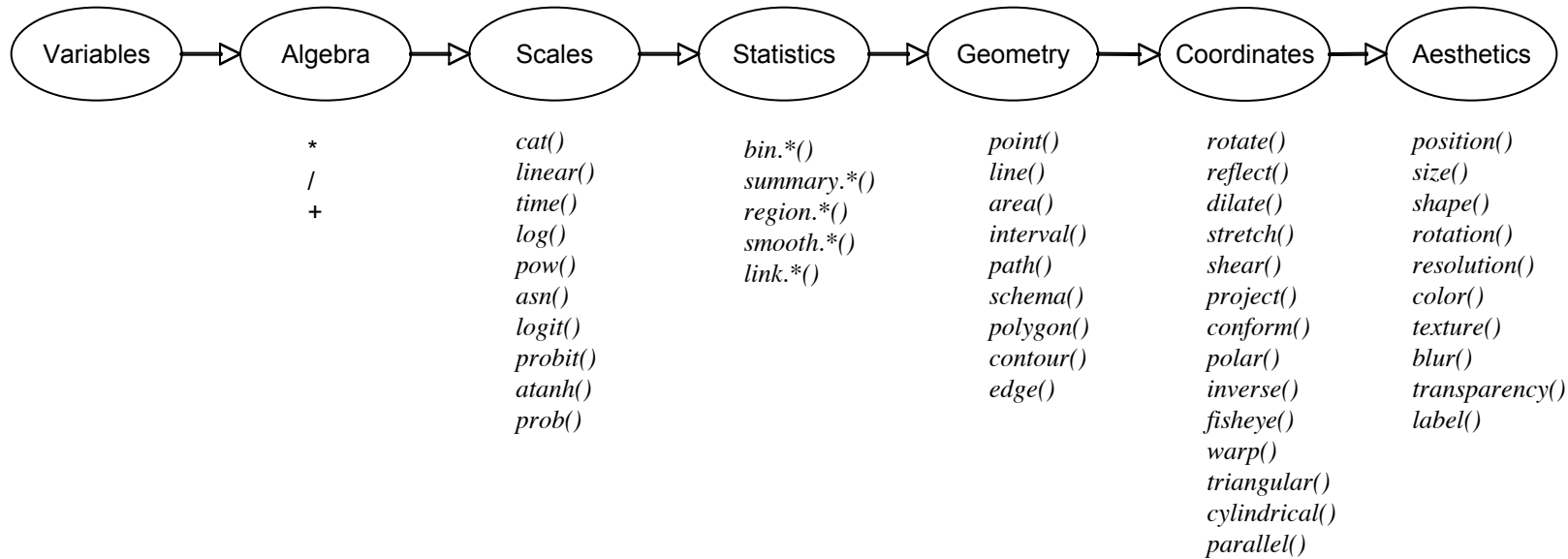


Coordinates

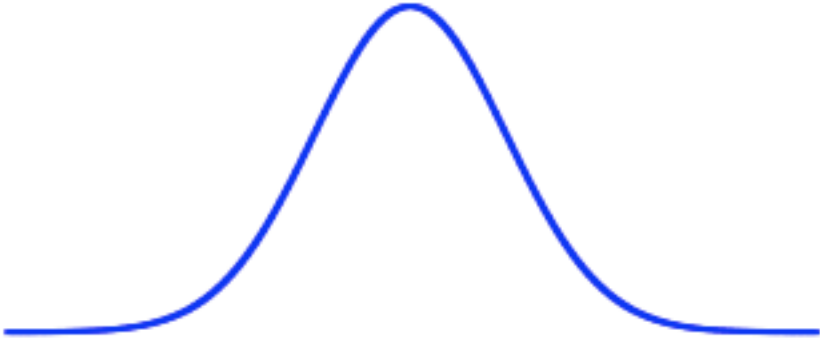
```
COORD: polar.theta(dim(3), polar(dim(1, 2)))  
ELEMENT: interval(position(summary.count(bin.rect(month*direction,  
dim(2))))))
```



Aesthetics



Aesthetics



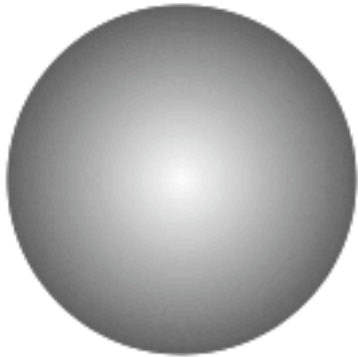
position



blur



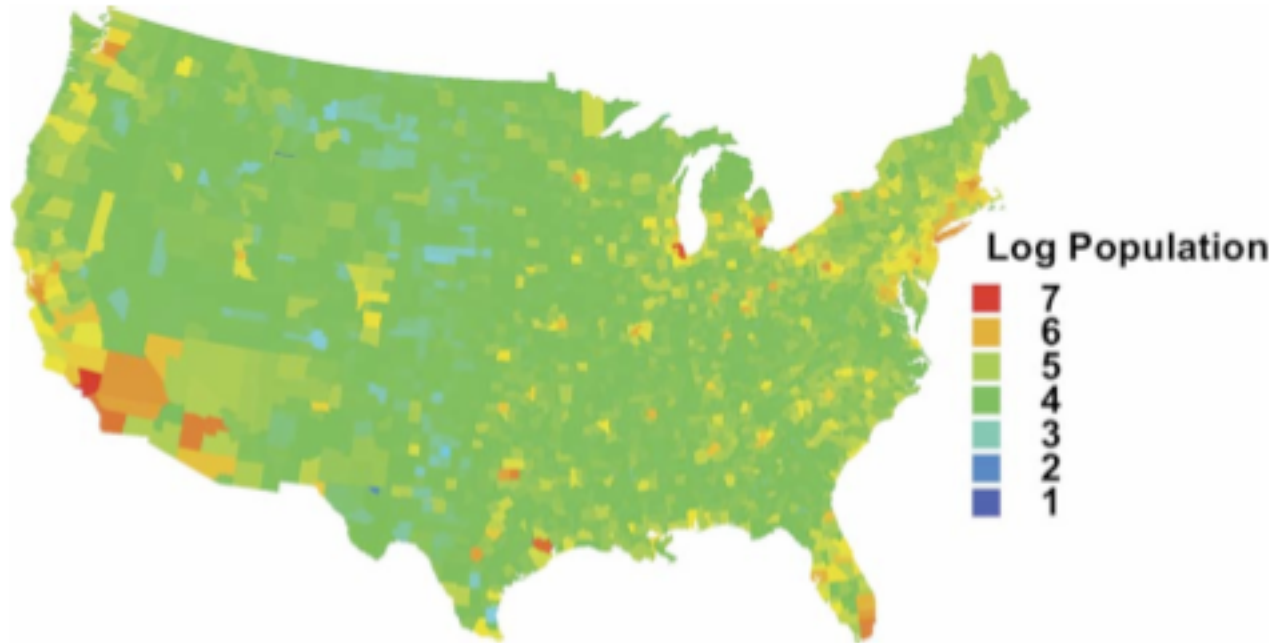
sound



brightness

Aesthetics

```
DATA: longitude, latitude = map(source("US counties"))  
TRANS: lpop = log.10(population)  
COORD: project.stereo(dim(1, 2))  
ELEMENT: polygon(position(longitude*latitude), color.hue(lpop))
```



Four GoG Implementations

SPSS nViZn

Wilkinson, Wills, Rope, ...



Stolte, Hanrahan, Mackinlay, ...

 **ggplot**

Wickham

Protovis

Bostock, Heer

Aspects of a Formal Graphics Language

- These are adjectives we sometimes use to describe formal languages:
 - **Simple** - few primitives, operators, and functions
 - **Expressive** - produces a large result set
 - **Coherent** - based on a grammar
 - **Meaningful** - makes statements that are falsifiable

Simple

- Simplicity in a language means having few primitives, operators, and functions.
- Let's take a look at a famous example and see how simple it is to build it in various graphics languages...

Simple

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

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

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui ont été en Russie; le noir ceux qui en sont sortis. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M.M. Cbiers, de Leger, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps de l'armée Napoléon et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilew et qui rejoignent vers Oescha et Wilkôk, avaient toujours marché avec l'armée.

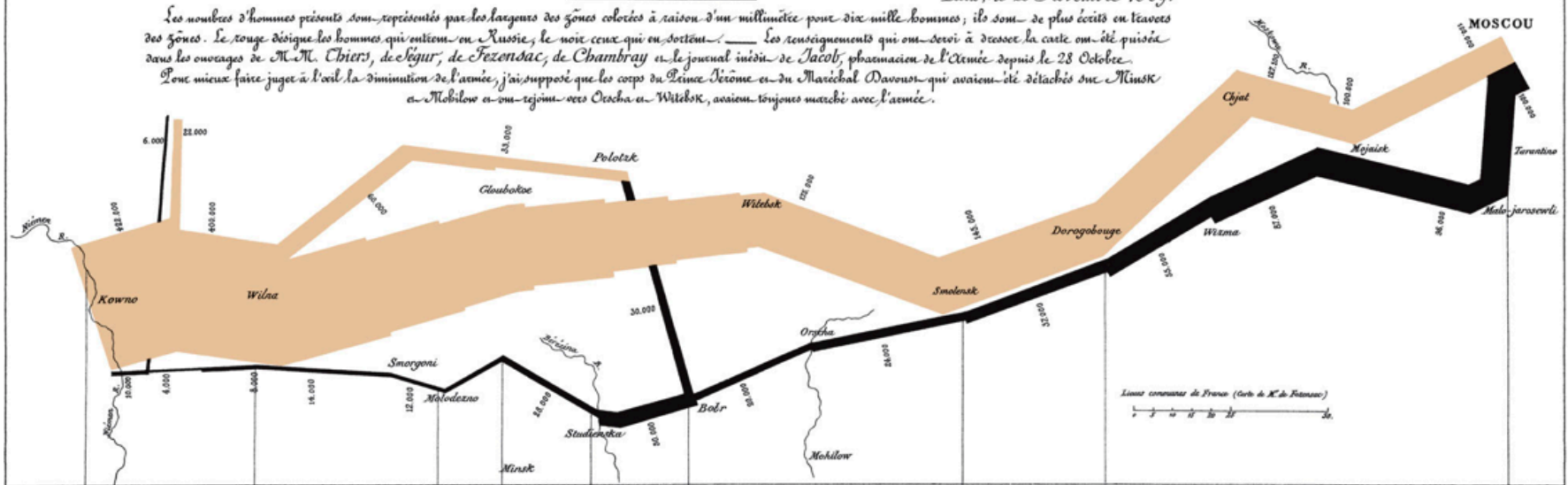
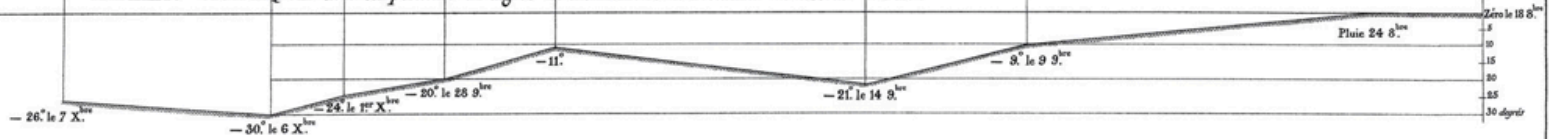


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop le Niémen gelé.



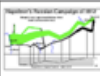
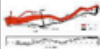


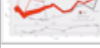






Orig. par Reynier, à Par. 5^e Mars 53 0⁰⁰ à Paris.

Imp. Lit. Reynier et Desobry.

Simple

You can find the following examples at this website:

<http://www.math.yorku.ca/SCS/Gallery/re-minard.html>

Picture	Words
	<p>Mathematica graphic Full size (297x183)</p> <p>In Chapter 4 of <i>Applied Mathematics: Getting Started, Getting it Done</i>, Shaw and Tigg (1993) describe a Mathematica function, <code>RecreateMinardMarch</code> to re-create a close approximation to Minard's graphic.</p> <p>A simpler function, <code>Minard[]</code> in the file <code>Minard.m</code> draws the same graph.</p>
	<p>Grammar of Graphics Full size (561 x 267, 22K)</p> <p>In Chapter 15 of <i>The Grammar of Graphics</i>, Wilkinson (1999) describes the encoding of the information from Minard's graphic according to his graph grammar. This is notable because the specification for the graph is declarative (what the graph consists of) rather than procedural (how to draw it).</p>
	<p>ggplot2: Grammar of Graphics in R March graphic, full size (1139 x 286, 49K), Temperature graphic, full size (1137 x 284, 31K), ZIP archive containing data and R code (56K)</p> <p>Hadley Wickham has implemented most of the ideas in the Grammar of Graphics in the R package <code>ggplot2</code>, available from CRAN servers. In these images, Hadley chose to illustrate the simple use of <code>ggplot2</code> to produce separate graphs of the flowmap of troop strength and the change in temperature on the retreat. See minard.r for the <code>ggplot2</code> calls that produce these images.</p> <p>This example is notable because, like the Grammar of Graphics specification, one can easily see the translation from the code (<code>minard.r</code>) to the graph. More so, because the <code>ggplot2</code> source code is freely available, so one can easily use it or extend it.</p>
	<p>SAS/IML Workshop Full size (1148 x 619, 19KB)</p> <p>SAS/IML Workshop is client-server extension of the SAS Interactive Matrix Language, with an object-oriented flavor. This version, programmed in NapoleonsMarch.jnl, works by defining a data structure and modules suitable to drawing paths on a 2D surface, a generalization of Minard's idea of the flow map.</p> <p>What they were thinking</p>
	<p>English translation Full size (919 x 712, 158K)</p> <p>ODT Maps has recently published Seeing Through Maps: The Power of Images to Shape Our World View by Ward Kaiser and Denis Wood.</p> <p>It contains this version of Minard's image with the caption translated into English and displayed legibly. The scale of temperature is displayed in both F and Celsius and given somewhat more prominence and readability.</p>
	<p>Image Map Full size graphic (541x308) [9K]</p> <p>The WWW and other new technology (CD rom, DVD) allows the easy integration of text, graphics, sound, and other media. Andrew Donoho teaches an Information Design course at the University of Texas at Austin. For the past several years, a course project asked students to design a web presentation based on Minard's graphic. Sunny's page is one example of the graphic linked (as an image map) to pages of text describing the major battles of the campaign.</p> <p>Lori Eichelberger's page is another example, relating the March on Moscow to the memoirs of Francois Bourgogne, a Sergeant in Napoleon's Imperial Guard. Other interesting attempts may be found on the LIS Class Projects page.</p>
	<p>Best online image Full size graphic (1762 x 833) [178K]</p> <p>From John Schneider's Napoleonic Literature page, this is the largest and most detailed re-creation of Minard's graphic I have found on the Web. All the place names and troop strengths are highly legible. The Napoleonic Literature site contains many electronic books and a wealth of information.</p>
	<p>Integrating temperature information Full size graphic (541x308) [9K] SAGE design page</p> <p>This graph, from the Visualization and Intelligent Interfaces group at CMU is one re-design, attempting to link the information about temperature directly to the path and strength of Napoleon's army, emphasizing their interrelations more directly. The (X,Y) coordinates are still map-based (latitude and longitude).</p> <p>"The use of color clearly shows the heat wave during the advance and the steady decline in temperature through the retreat. The exception, a spell of temperatures above freezing, is clearly visible when the retreating army is between the cities of Krasnyj and Bobr."</p>
	<p>Animated graphic Full size graphic (536x182) [166K] from Max Boylan's march page</p> <p>Here is a similar graphic, but done as an animated GIF image, from Max Boylan at SFSU. We start with an empty map on June 24, 1812, as the campaign begins at the Nemèn River with 422,000 troops. The animation shows the step-by-step progress of the Grand Army to Moscow and back. The coding of temperature by color is done somewhat better than in the SAGE graphic above. The animation repeats 10 times.</p>
	<p>Making time explicit Full size graphic (541x308) [9K] SAGE design page</p> <p>"The temporal characteristics of the march were weakly conveyed in the original because they were expressed as text labels. Here the relation between date, troop location (longitude only), and temperature is more strongly conveyed because they are all expressed as properties of the rectangles. Battle sites are indicated in text."</p>
	<p>Animated graphic Full size graphic (547 x 346) [388K]</p> <p>This version, by Aaron Walburg and Stephen Hartzog, listed on the University of York History of Statistics page uses a real map of the region from Poland to Moscow as the background. It overlays on this an animated sequence showing the size of Napoleon's army (with dates), and the retreat, together with the decline of temperature. In many ways, this is my favorite modern rendition.</p>

Simple? - Mathematica



Simple? - Mathematica

```
BeginPackage["Minard`"]
```

```
Minard::usage = "Minard.m is a package for producing thematic maps  
in the style of C. J. Minard"
```

```
ProcessTemp::usage = "ProcessTemp[list] computes the graphics primitives to be associated with a list of elements  
of the form {x, y, temp}"
```

```
stuff::usage = "stuff[list] computes the graphics primitives to be associated with a list of elements of the form  
{x, y, strength}"
```

```
ProcessStrength::usage = "ProcessStrength[list] computes the graphics primitives to be associated with a list of  
elements of the form {{x1, y1, strength1},..., {xn, yn, strengthn}}"
```

```
ProcessRivers::usage = "ProcessRivers[list] computes the graphics primitives to be associated with a list of river  
coordinates"
```

```
ProcessPoints::usage = "ProcessPoints[list] computes the graphics primitives to be associated with a list of city  
or other point coordinates"
```

```
ProcessTitle::usage = "ProcessTitle[location, title] computes the graphics primitives to be associated with the  
title"
```

```
ProcessBoxes::usage = "ProcessBoxes[list] computes the graphics primitives to be associated with a list of box  
coordinates"
```

```
ProcessText::usage = "ProcessText[list] computes the graphics primitives to be associated with a list of  
coordinates and associated text"
```

```
NapoleonicMarchOnMoscowAndBackAgainPlot::usage =  
"NapoleonicMarchOnMoscowAndBackAgainPlot[strength_list, temp_list, riverdata_list, boxdata_list, titledata_list,  
pointdata_list, textdata_list] shows the thematic chart associated with the strength temperature, river, box,  
title, point and text lists, in the style of C. J. Minard."
```

Simple? - Mathematica

```
Begin["`Private`"]
ProcessTemp[tdata_] :=
Module[{tprimlist, coords = Array[coordarr, 100], numb, k,
temprules, firstx, lastx},
temprules = {0, -10, -20, -30};
firstx = First[tdata][[1]];
lastx = Last[tdata][[1]];
tprimlist = Map[({Thickness[0.001], Line[{{firstx, 0.0916 + 0.00188333*#},
{lastx, 0.0916 + 0.00188333*#}}]}&, temprules];
numb = Length[tdata];
Do[
coords[[k]] = {tdata[[k, 1]], 0.0916 + tdata[[k,3]]*0.00188333}
,{k, 1, numb}];
Do[
tprimlist = Append[tprimlist,
{Thickness[0.001], Line[{tdata[[k, Range[2]]] , coords[[k]]} ]}];
,{k, 1, numb}];
Do[
tprimlist = Append[tprimlist,
{Thickness[0.001], Line[{coords[[k]] , coords[[k+1]]} ]}];
,{k, 1, numb-1}];
tprimlist]

stuff[sd_] := Module[{k, m, primlist, pca, pcb, pcc, pcd,
scale, dir = Array[dirarr, 100],
rotmdir = Array[rotmdirarr, 100],
l = Array[larr, 100],
avevec = Array[avevecarr, 100]},
scale = 10500000;
primlist = { };
k = Length[sd];
```

Simple? - Mathematica

```
Do[(dir[[m]] = sd[[m+1, Range[2] ]] - sd[[m, Range[2] ]];
l[[m]] = Sqrt[dir[[m, 1]]^2 + dir[[m, 2]]^2),
{m, 1, k-1}];
Do[
If[l[[m]] > 0.0001,
rotmdir[[m]] = {-dir[[m, 2]], dir[[m, 1]]}/l[[m]],
rotmdir[[m]] =
{-dir[[m+1, 2]], dir[[m+1, 1]]}/l[[m+1]]],
{m, 1, k-1}];
avevec[[1]] = rotmdir[[1]];
avevec[[k]] = rotmdir[[k-1]];
Do[
avevec[[m]] =
(rotmdir[[m]] + rotmdir[[m-1]])/(1 +
rotmdir[[m,1]]*rotmdir[[m-1,1]]
+ rotmdir[[m,2]]*rotmdir[[m-1,2]]), {m, 2, k-1}];
Do[
pca = sd[[m,Range[2]]] - sd[[m,3]]*avevec[[m]]/scale;
pcb = sd[[m+1,Range[2]]] - sd[[m+1,3]]*avevec[[m+1]]/scale;
pcc = sd[[m+1,Range[2]]] + sd[[m+1,3]]*avevec[[m+1]]/scale;
pcd = sd[[m,Range[2]]] + sd[[m,3]]*avevec[[m]]/scale;
u = {If[sd[[m, 4]] == 1,
RGBColor[0,1,0],
RGBColor[0,0,0]],
Polygon[{pca, pcb, pcc, pcd}]];
primlist = Append[primlist, u],
{m, 1, k-1}];
primlist]
```

Simple? - Mathematica

```
ProcessStrength[data_] :=  
Module[{strengthprims, str, len, v},  
strengthprims = {};  
len = Length[data];  
Do[  
(strengthprims =  
Append[strengthprims, stuff[data[[v]]]]),  
{v, 1, len}];  
strengthprims  
]  
  
ProcessTitle[titledata_] :=  
Text[FontForm[titledata[[1, 3]], {"Helvetica-Bold", 14}],  
{titledata[[1, 1]], titledata[[1, 2]]]  
  
ProcessPoints[pointdata_] :=  
Map[({Point[Drop[#, -1]], Text[FontForm[Last[#], {"Times-Roman", 8}],  
Drop[#, -1]+{0.01, 0.01}, {0, -1}]}&, pointdata]  
  
ProcessText[textdata_] :=  
Map[({Text[FontForm[Last[#], {"Times-Roman", 6}], Drop[#, -1], {0, -1}]}&, textdata]  
  
ProcessRivers[riverdata_] :=  
Map[({RGBColor[0, 0, 1], Thickness[0.001], Line[#]}&, riverdata]  
  
ProcessBoxes[boxdata_] :=  
Map[({RGBColor[0, 0, 0], Thickness[0.002], Line[#]}&, boxdata]
```

Simple? - Mathematica

```
NapoleonicMarchOnMoscowAndBackAgainPlot[sdata_, tdata_, riverdata_, boxdata_, titledata_, pointdata_, textdata_] :=  
Show[Graphics[  
  {ProcessStrength[sdata],  
   ProcessTemp[tdata],  
   ProcessRivers[riverdata],  
   ProcessBoxes[boxdata],  
   ProcessTitle[titledata],  
   ProcessPoints[pointdata],  
   ProcessText[textdata]}  
]]
```

```
Minard[] :=  
NapoleonicMarchOnMoscowAndBackAgainPlot[  
  StrengthData,  
  TempData,  
  RiverData, (* river *)  
  BoxData,  
  TitleData, (* title *)  
  PointData, (* point *)  
  TextData (* text *)  
]
```

```
TempData = {...};
```

```
RiverData = {...};
```

```
armydata = {...};
```

```
StrengthData = {...};
```

```
BoxData = {...};
```

```
TitleData = {...};
```

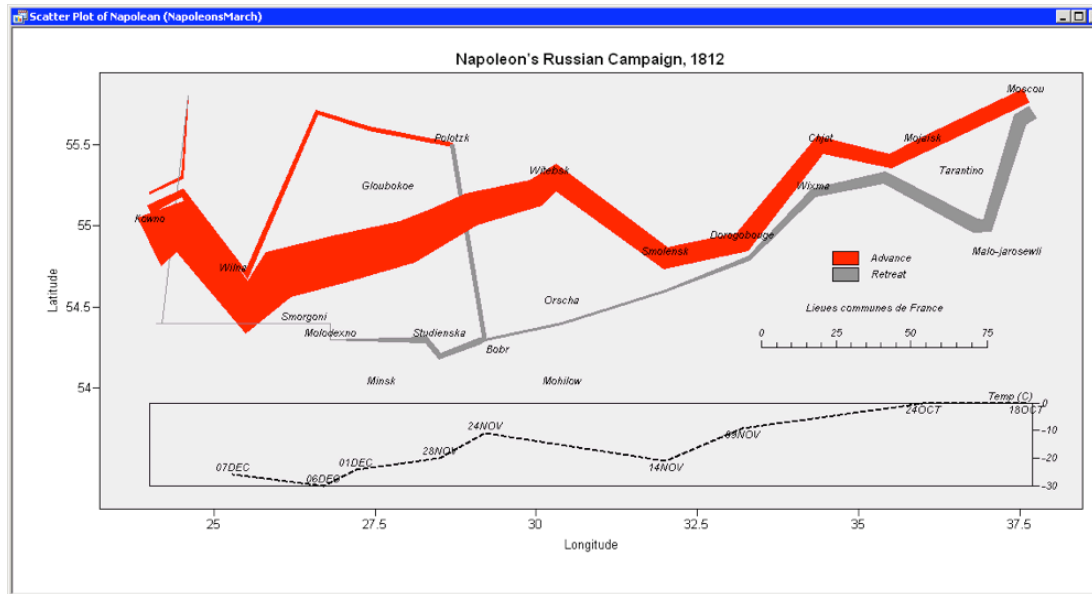
```
PointData = {...};
```

```
TextData = {...};
```

```
End[]
```

```
EndPackage[]
```

Simple? - SAS



Simple? - SAS

```
run DefineData;

declare DataObject dobj;
dobj = DataObject.Create( "Napolean", {"Longitude" "Latitude" "ArmySize" "ArmyDir" "ArmyGroup"}, army );

declare ScatterPlot plot;
plot = ScatterPlot.Create( dobj, "Longitude", "Latitude", false );
plot.SetWindowPosition( 0, 25, 100, 75 );
plot.SetGraphAreaMargins( 0.08, 0.02, 0.08, 0.15 );
plot.SetAxisViewRange( YAXIS, cityLat[><], cityLat[<>] );
plot.ShowObs( false );
plot.SetTitleText( "Napoleon's Russian Campaign, 1812" );
plot.ShowTitle();
plot.DrawUseDataCoordinates();

sizeScale = 6000;
group = unique(armyGroup);
do i = ncol(group) to 1 by -1;
  idx = loc( armyGroup = group[i] );
  /* group's advance */
  forward = loc( armyDir[idx] = 1 );
  plot.DrawSetPenColor( RED );
  plot.DrawSetBrushColor( RED );
  k = idx[forward];
  run PlotContinuousDataPath( armyLon[k], armyLat[k], armySize[k]/sizeScale, plot );

  /* group's retreat */
  backward = loc( armyDir[idx] = -1 );
  plot.DrawSetPenColor( GREY );
  plot.DrawSetBrushColor( GREY );
  k = idx[backward];
  run PlotContinuousDataPath( armyLon[k], armyLat[k], armySize[k]/sizeScale, plot );
end;
```


Simple? - SAS

```
/* label cities for geographical reference */
plot.DrawSetTextSize( 8 );
plot.DrawSetTextStyle( STYLE_ITALIC );
plot.DrawSetTextAlignment( ALIGN_CENTER, ALIGN_BOTTOM );
plot.DrawText( cityLon, cityLat, cityNames );

/* draw reference axis showing scale in miles */
plot.DrawAxis(33.5, 54.25, 37, 54.25, {0 0.33 0.66 1},
    {'0' '25' '50' '75'}, {0.2 0.4 0.6 0.8}, 90 );
plot.DrawSetTextSize( 8 );
plot.DrawText( 35.25, 54.45, "Lieues communes de France" );

/* draw simple legend */
plot.DrawSetTextAlignment( ALIGN_LEFT, ALIGN_CENTER );
plot.DrawSetPenColor( BLACK );
plot.DrawSetBrushColor( RED );
plot.DrawRectangle( 34.6, 54.76, 35, 54.84, true );
plot.DrawText( 35.2, 54.8, "Advance" );
plot.DrawSetBrushColor( GREY );
plot.DrawRectangle( 34.6, 54.66, 35, 54.74, true );
plot.DrawText( 35.2, 54.7, "Retreat" );

/* increase bottom margin and place temperature graphic within */
plot.GetAxisViewRange( YAXIS, ymin, ymax );
plot.SetAxisTickRange( YAXIS, ymin, ymax );
boxMin = ymin - (ymax-ymin)/3; /* decrease by 1/3 */
boxMax = ymin - (ymax-ymin)/20; /* and by 1/20 */

plot.SetAxisViewRange( YAXIS, boxMin, ymax );
plot.GetAxisViewRange( XAXIS, xmin, xmax );

tempMin = -30;
tempMax = 0;
tempInc = 10;
```

Simple? - SAS

```
/* draw bounding box and temperature axis */
plot.DrawRectangle( xmin, boxMin, xmax, boxMax );
plot.DrawNumericAxis( xmax, boxMin, xmax, boxMax, tempMin, tempMax, 4 );

/* scale temperature so we can plot in data coordinates */
temp = (retreatTemp - tempMin)/(tempMax-tempMin);
temp = temp * ( boxMax - boxMin ) + boxMin;

/* draw temperature versus longitude */
plot.DrawSetPenWidth( 2 );
plot.DrawSetPenStyle( DASHED );
plot.DrawLine( retreatLon, temp );

/* label temperature axis */
plot.DrawSetTextAlignment( ALIGN_RIGHT, ALIGN_BOTTOM );
plot.DrawText( xmax, boxMax, "Temp (C)" );

/* label temperature plot with dates of retreat */
k = 1:4;
plot.DrawSetTextAlignment( ALIGN_CENTER, ALIGN_TOP );
plot.DrawText( retreatLon[k], temp[k], putn(retreatDate[k], 'DATE5.' ) );
k = 5:nrow(retreatLon);
plot.DrawSetTextAlignment( ALIGN_CENTER, ALIGN_BOTTOM );
plot.DrawText( retreatLon[k], temp[k], putn(retreatDate[k], 'DATE5.' ) );

plot.ActivateWindow();

start DefineData;
GREY = 0808080x;
/* city names and locations for annotations */
cityNames = {...};
cityPos = {...};
cityLon = cityPos[,1];
cityLat = cityPos[,2];
```

Simple? - SAS

```
/* temperature during retreat */
retreat = {...};
retreatLon = retreat[,1];
retreatTemp = retreat[,2];
retreatDate = retreat[,3];

army = {...};
finish;

start PlotContinuousDataPath( mX, mY, mSize, Plot plot );
  x = rowvec( mX );
  y = rowvec( mY );
  size = rowvec( mSize );

  n = ncol(x);
  if n < 2 then do;
    print "GetContinuousPath: path must contain at least two vertices";
    abort;
  end;
  if n ^= ncol(y) | n ^= ncol(size) then do;
    print "GetContinuousPath: parameter arrays are different lengths";
    abort;
  end;

  /* translate all coordinates into pixel coords */
  plot.GetGraphAreaWidthHeight( winWidth, winHeight );

  /* find pixel coordinates of frame */
  plot.GetGraphAreaMargins( leftGraphMargin, rightGraphMargin,
                           topGraphMargin, bottomGraphMargin );

  plotHeight = (1 - topGraphMargin - bottomGraphMargin) * winHeight;
  plotWidth = (1 - leftGraphMargin - rightGraphMargin) * winWidth;
  plotLeft = leftGraphMargin * winWidth;
  plotRight = plotLeft + plotWidth;
  plotBottom = bottomGraphMargin * winHeight;
  plotTop = plotBottom + plotHeight;
```

Simple? - SAS

```
/* find pixel coordinates of plot area (includes plot margins) */
plot.GetPlotAreaMargins( leftPlotMargin, rightPlotMargin,
                        topPlotMargin, bottomPlotMargin );
drawHeight = (1 - topPlotMargin - bottomPlotMargin) * plotHeight;
drawWidth  = (1 - leftPlotMargin - rightPlotMargin) * plotWidth;
drawLeft   = plotLeft + leftPlotMargin * plotWidth;
drawRight  = drawLeft + drawWidth;
drawBottom = plotBottom + bottomPlotMargin * plotHeight;
drawTop    = drawBottom + drawHeight;

run GetPlotAreaDataCoordinates( plotXMin, plotXMax, plotYMin, plotYMax, plot );

/* translate data into [plot area] pixels */
x = plotLeft + (plotRight-plotLeft)/(plotXMax-plotXMin) * (x - plotXMin);
y = plotBottom + (plotTop-plotBottom)/(plotYMax-plotYMin) * (y - plotYMin);

path = j(2*n, 2, .);

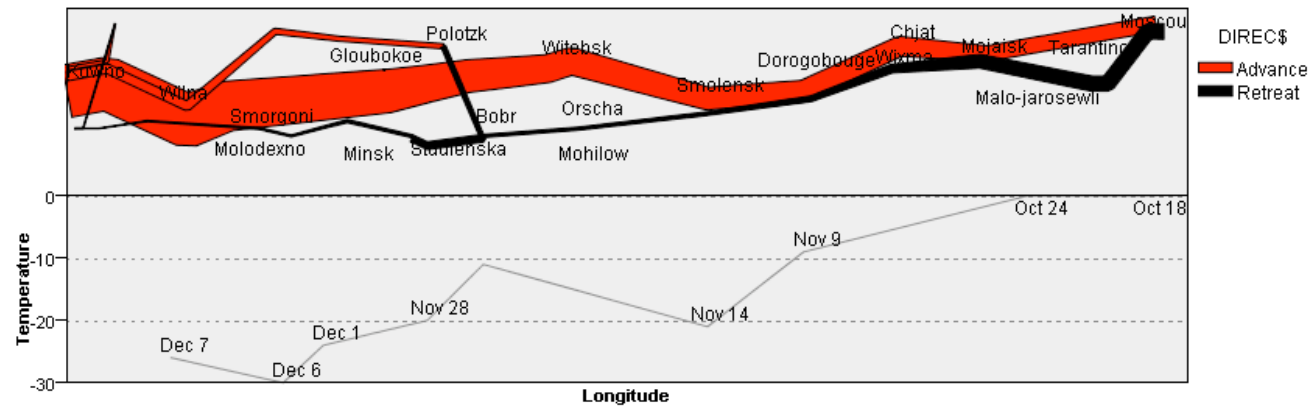
/* find points that define width of first line segment */
u = y[1]-y[2] || x[2]-x[1];
uNorm = u / sqrt(ssq(u)); /* unit length */
wx = size[1]/2 * uNorm[1];
wy = size[1]/2 * uNorm[2];
path[1,1] = x[1] + wx;
path[1,2] = y[1] + wy;
path[2*n,1] = x[1] - wx;
path[2*n,2] = y[1] - wy;
```

Simple? - SAS

```
/* find points that define width of last line segment */
u = y[n-1]-y[n] || x[n]-x[n-1];
uNorm = u / sqrt(ssq(u)); /* unit length */
wx = size[n]/2 * uNorm[1];
wy = size[n]/2 * uNorm[2];
path[n,1] = x[n] + wx;
path[n,2] = y[n] + wy;
path[n+1,1] = x[n] - wx;
path[n+1,2] = y[n] - wy;

/* now convert back to data coordinates */
path[,1] = plotXMin + (plotXMax-plotXMin)/(plotRight-plotLeft) * (path[,1] - plotLeft);
path[,2] = plotYMin + (plotYMax-plotYMin)/(plotTop-plotBottom) * (path[,2] - plotBottom);
plot.DrawUseDataCoordinates();
plot.DrawPolygon( path[,1], path[,2], true );
finish;
```

Simple - nViZn



Simple - nViZn

The blue statements are data definitions. The rest is the graphics program.
This program executes in SPSS to produce the Minard graphic on the last slide.

```
SOURCE: troops = csvSource(file("../data/minard troops.csv"))
SOURCE: cities = csvSource(file("../data/minard cities.csv"))
SOURCE: temperature = csvSource(file("../data/minard temperature.csv"))
DATA: lonp = col(source(troops), name("LON"))
DATA: latp = col(source(troops), name("LAT"))
DATA: surviv = col(source(troops), name("SURVIV"))
DATA: direc = col(source(troops), name("DIREC$"), unit.category())
DATA: division = col(source(troops), name("DIVISION$"), unit.category())
DATA: latc = col(source(cities), name("LAT"))
DATA: lonc = col(source(cities), name("LON"))
DATA: city = col(source(cities), name("CITY$"), unit.category())
DATA: lont = col(source(temperature), name("LON"))
DATA: temp = col(source(temperature), name("TEMP"))
DATA: date = col(source(temperature), name("DATE$"), unit.category())

GRAPH: begin(origin(0.0px, 0.0px), scale(720.0px, 120.0px))
GUIDE: axis(dim(1), null())
GUIDE: axis(dim(2), null())
GUIDE: legend(aesthetic(aesthetic.size), null())
SCALE: cat(aesthetic(aesthetic.color), map(("Advance",color.red), ("Retreat",color.black)))
SCALE: linear(aesthetic(aesthetic.size), aestheticMinimum(size."0px"), aestheticMaximum(size."20px"))
ELEMENT: path(position(lonp*latp), size(surviv), color(direc), split(division))
ELEMENT: point(position(lonc*latc), size(size."0px"), label(city))
GRAPH: end()

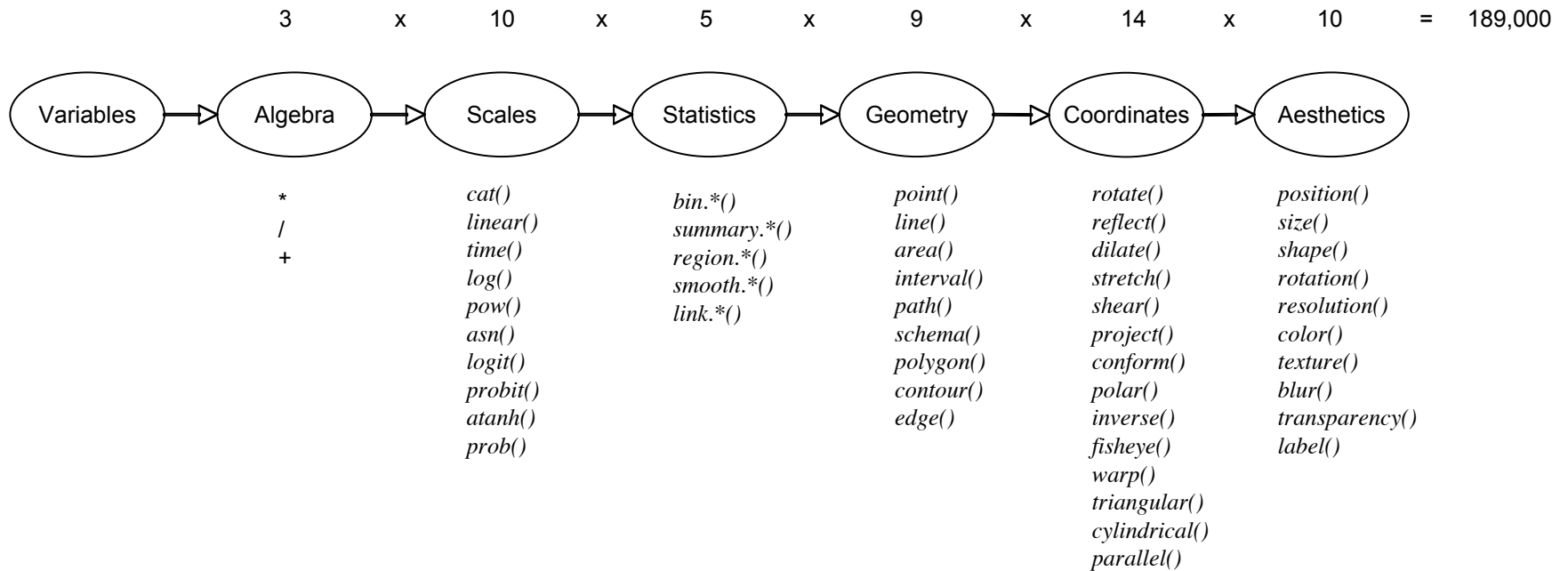
GRAPH: begin(origin(0.0px, 120.0px), scale(720.0px, 120.0px))
GUIDE: axis(dim(1), label("Longitude"), ticks(null()))
GUIDE: axis(dim(2), label("Temperature"), start(0.0), delta(10.0), gridlines())
ELEMENT: path(position(lont*temp), color(color.gray))
ELEMENT: point(position(lont*temp), label(date), size(size."0px"))
GRAPH: end()
```

Expressive

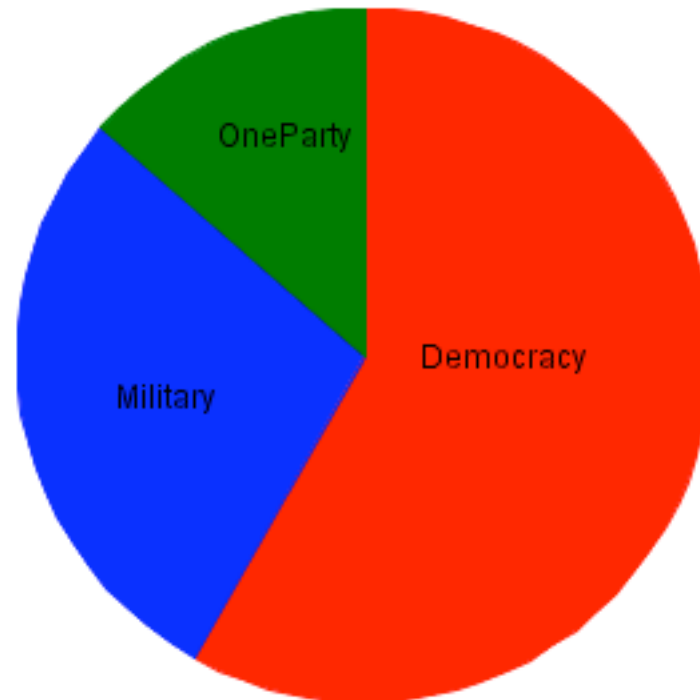
- Expressive means having a large results set.
- The following examples suggest the range of results the GoG model can produce.
- The range of possibilities greatly exceeds the capabilities of graphics production libraries and toolkits such as VTK.

Expressive

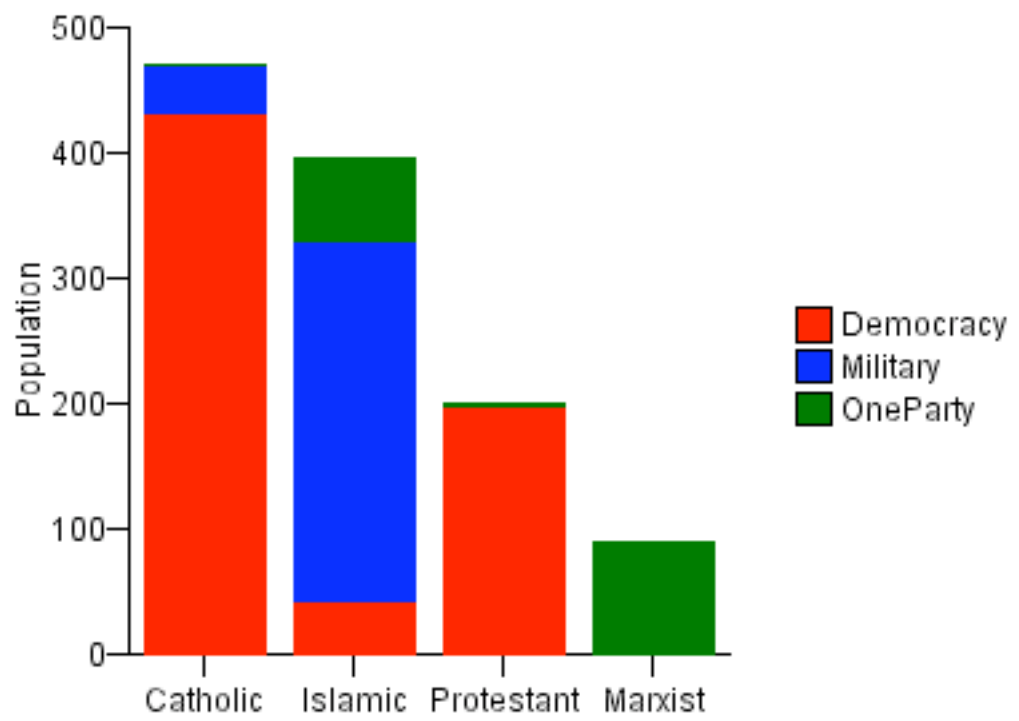
One can get an idea of the magnitude of results by computing a product of methods:



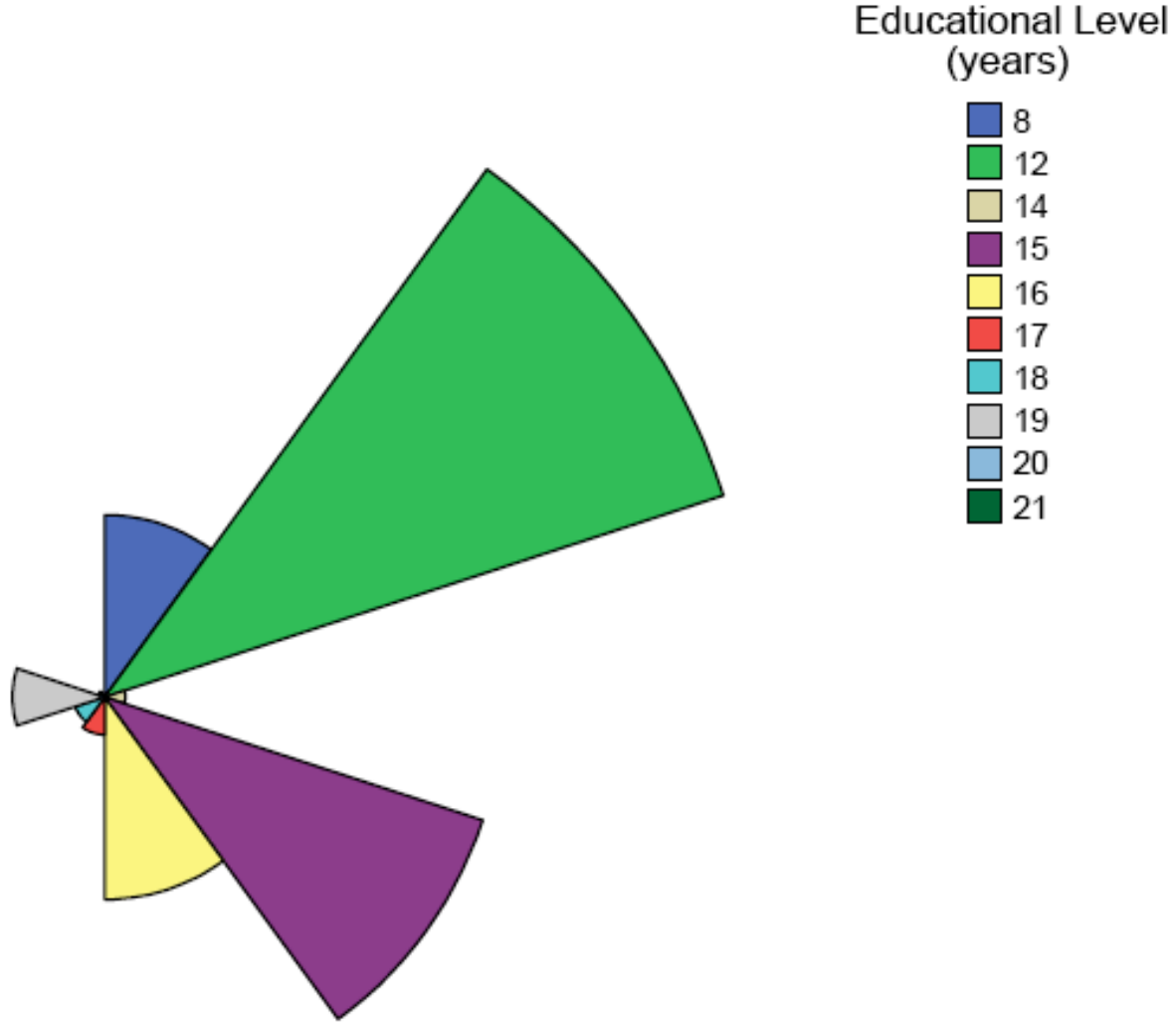
Expressive - some nViZn examples



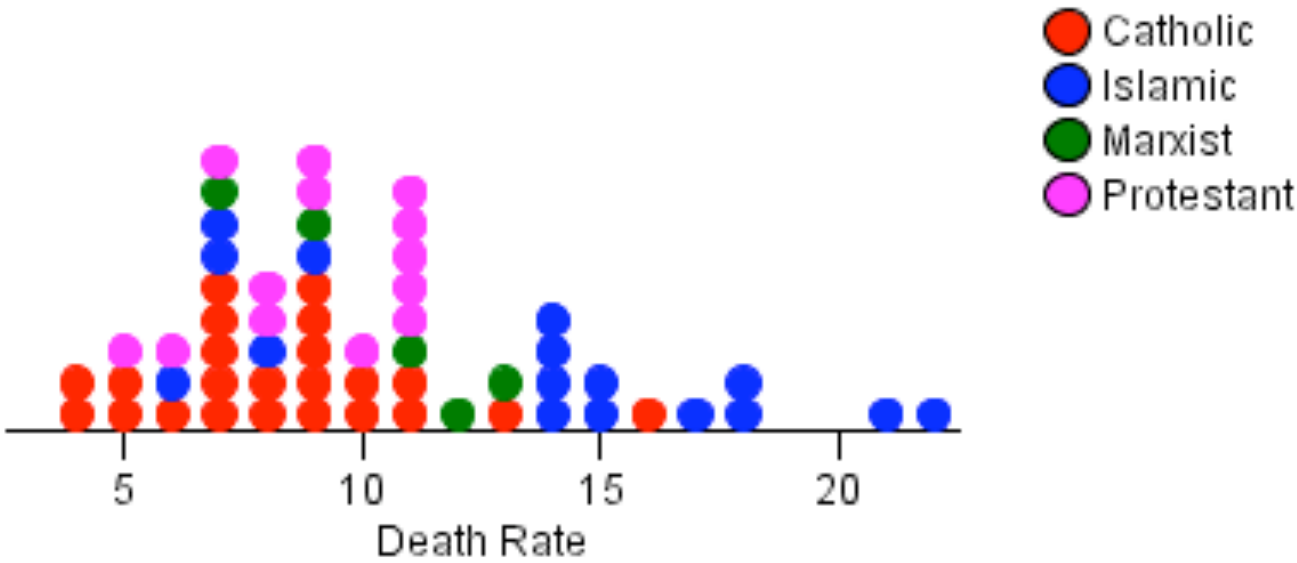
nViZn



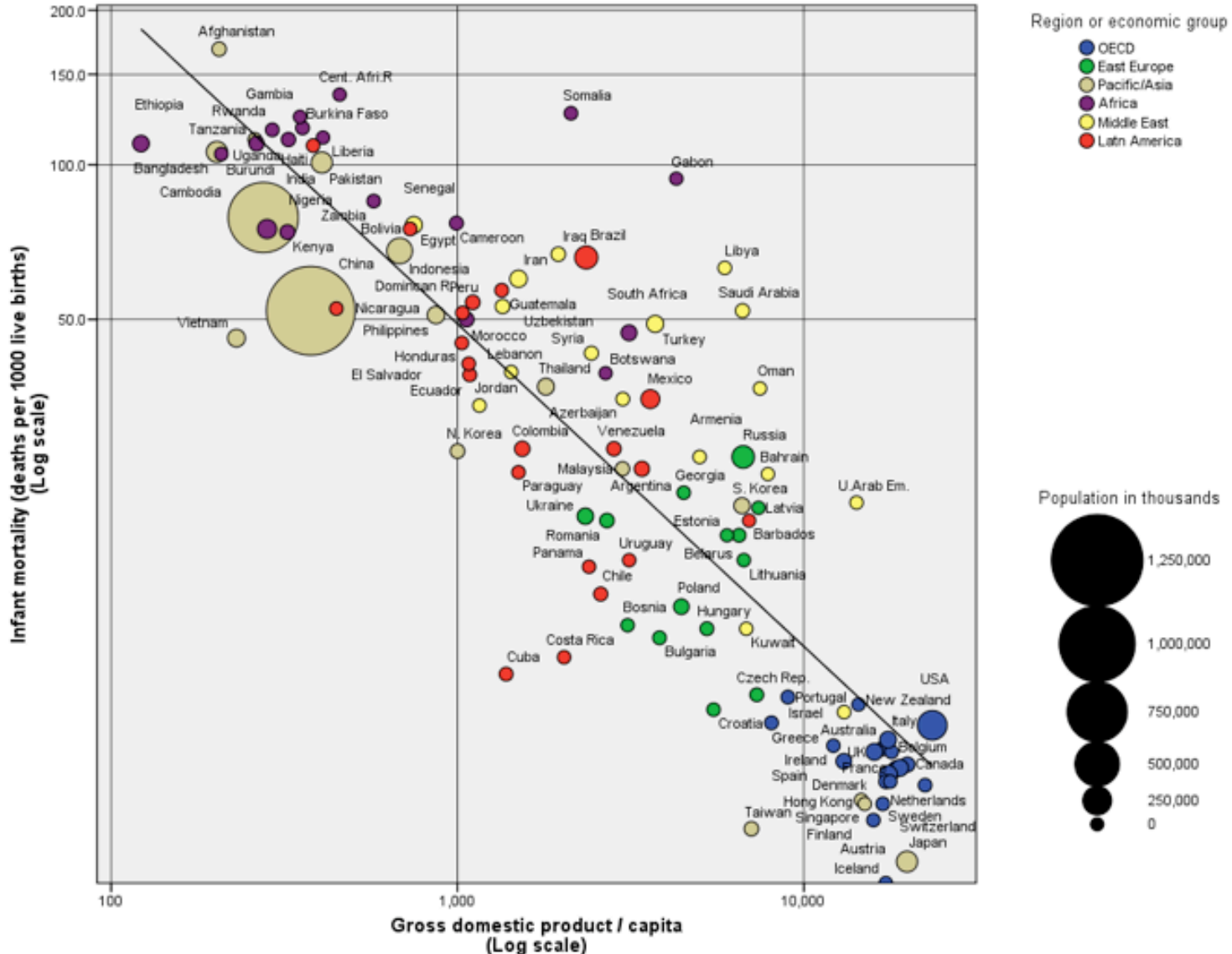
nViZn



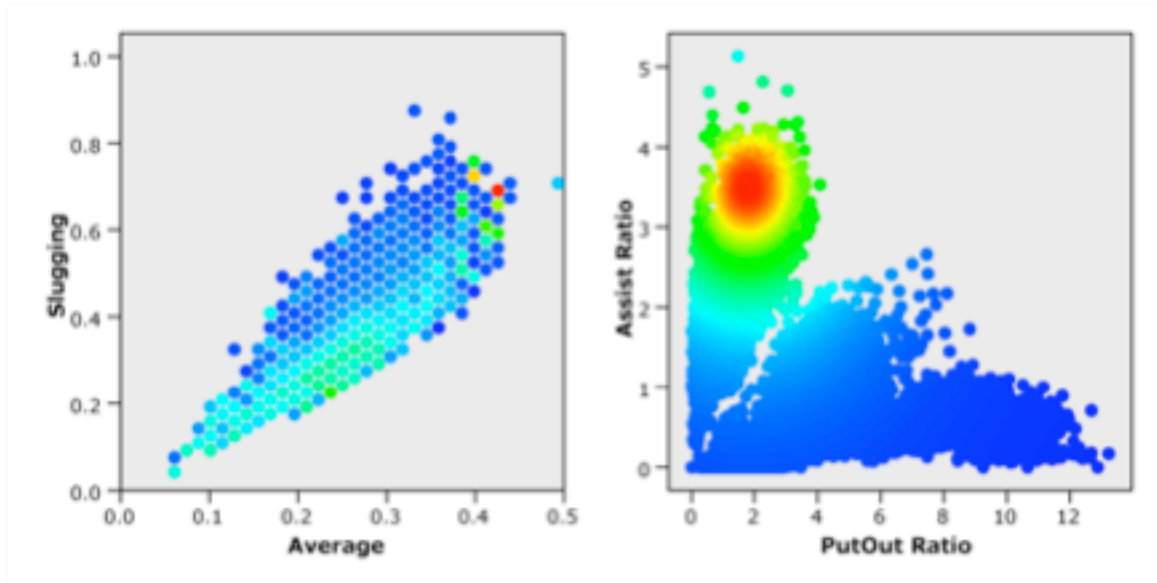
nViZn



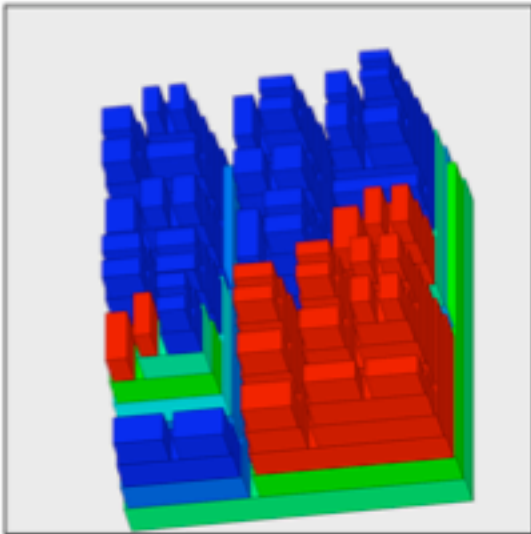
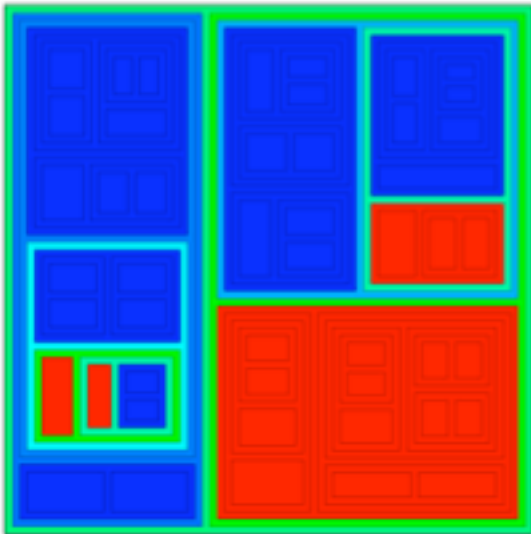
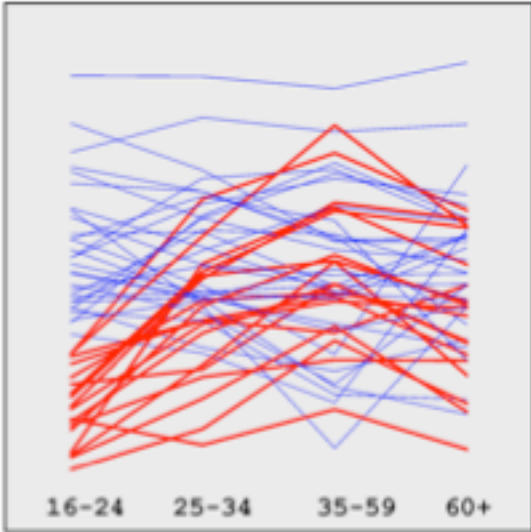
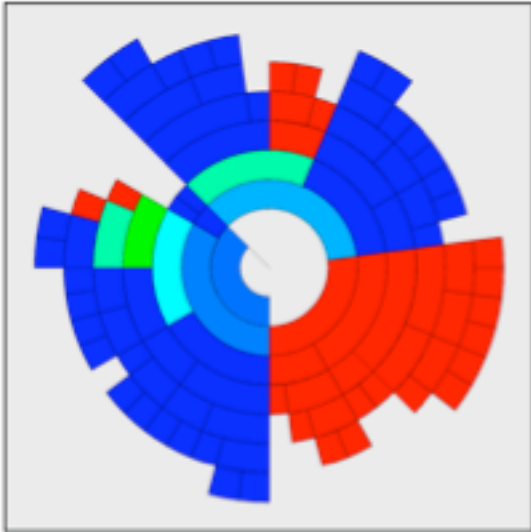
nViZn



nViZn



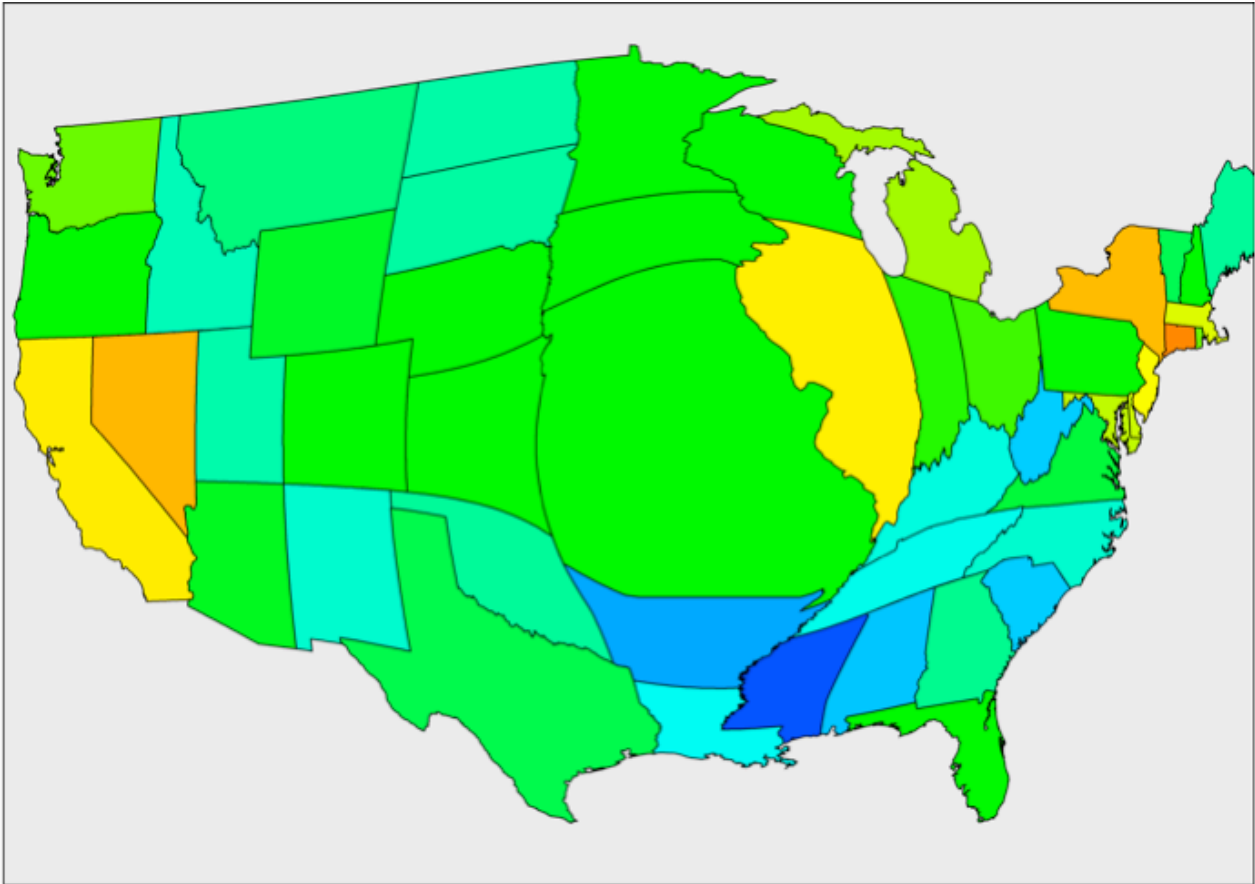
nViZn



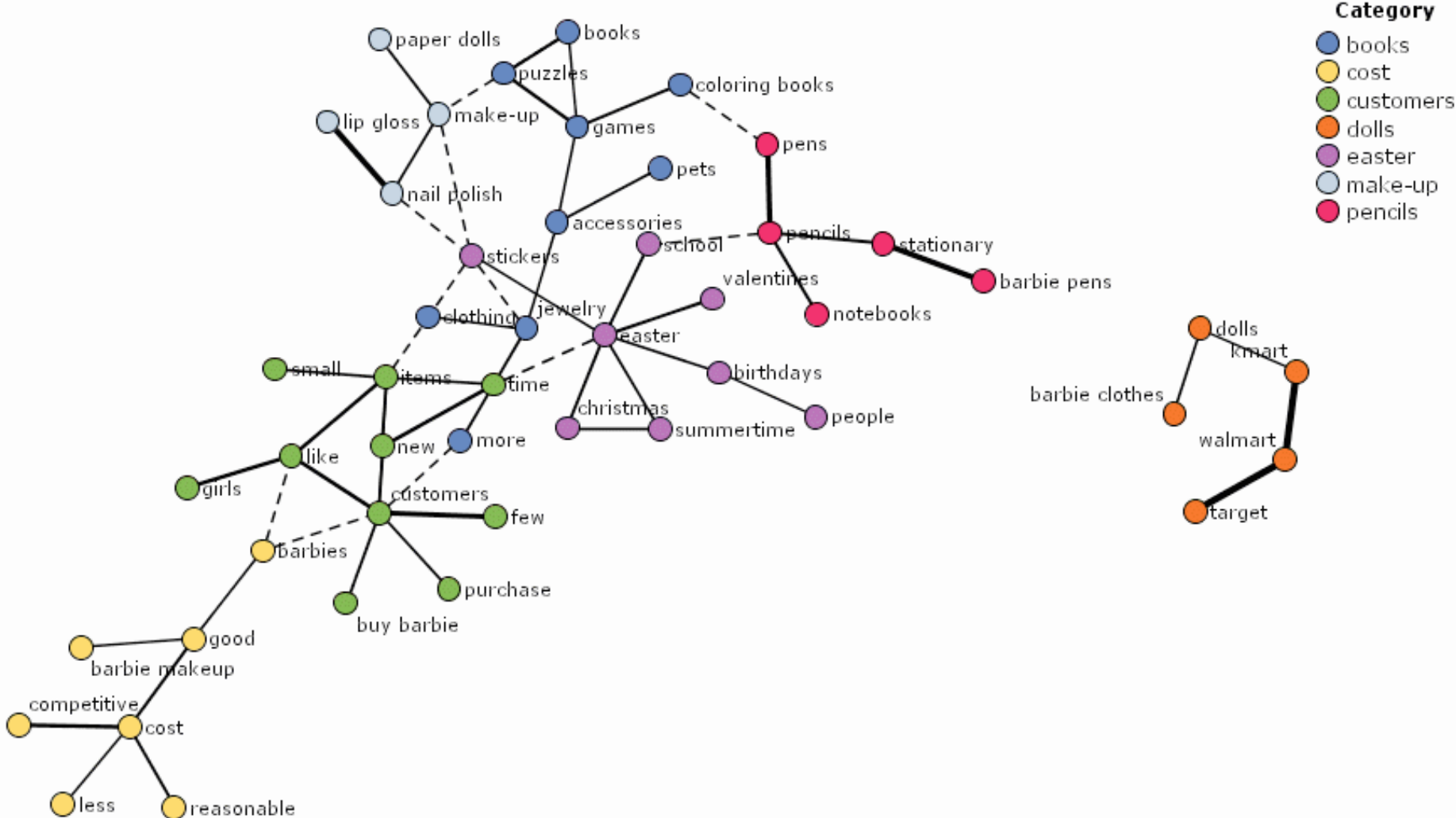
nViZn



nViZn

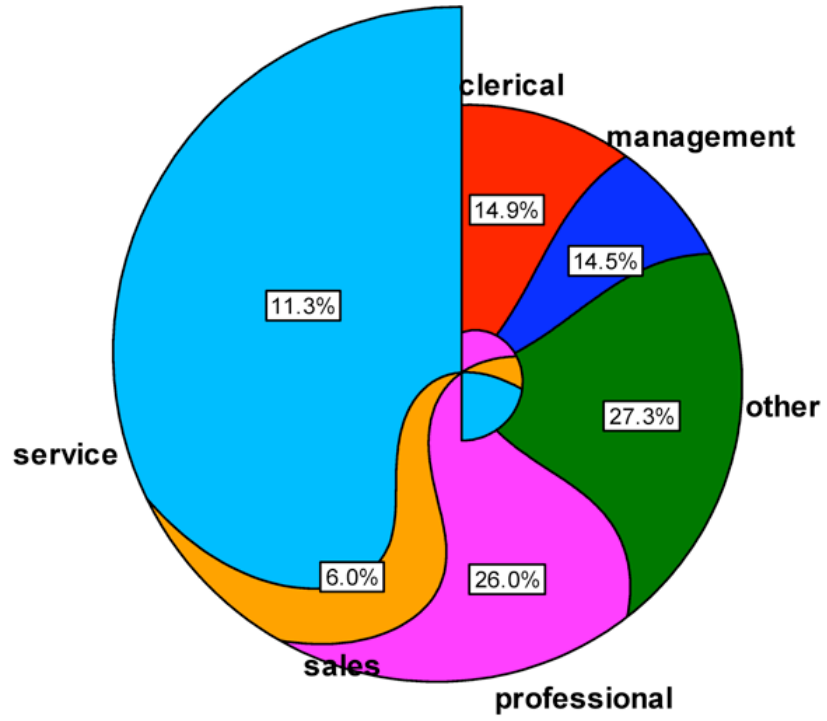


nViZn



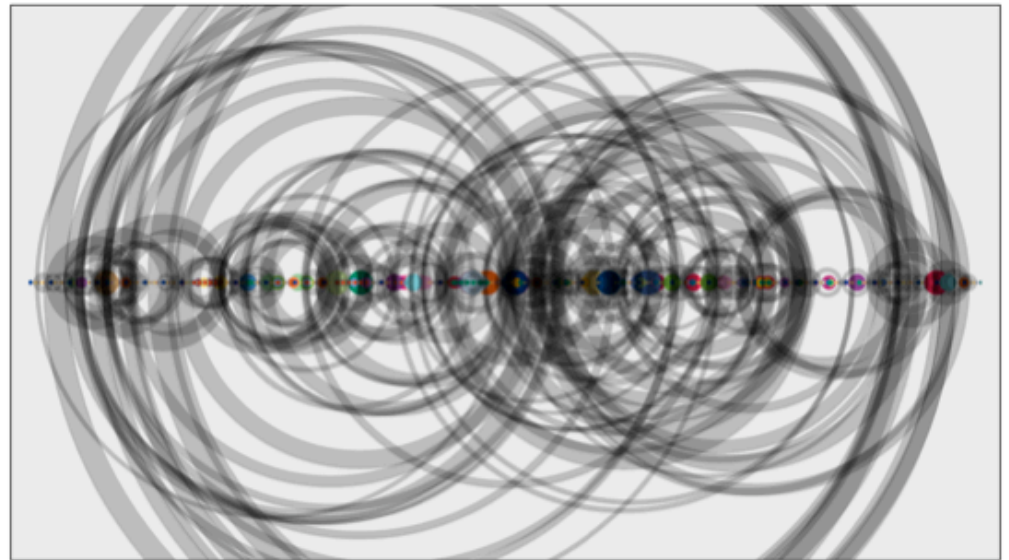
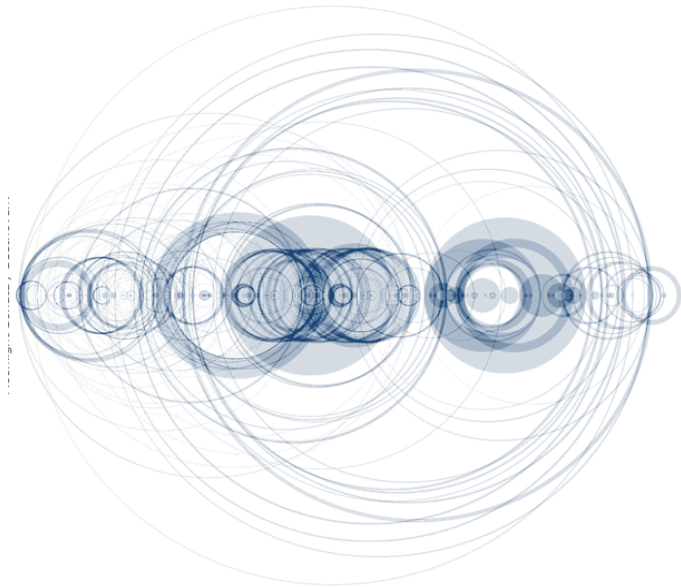
nViZn

Expressive does not mean beautiful or useful. Here is a fisheye transformation that was accidentally applied to a pie chart by a programmer. It is meaningful but ugly.



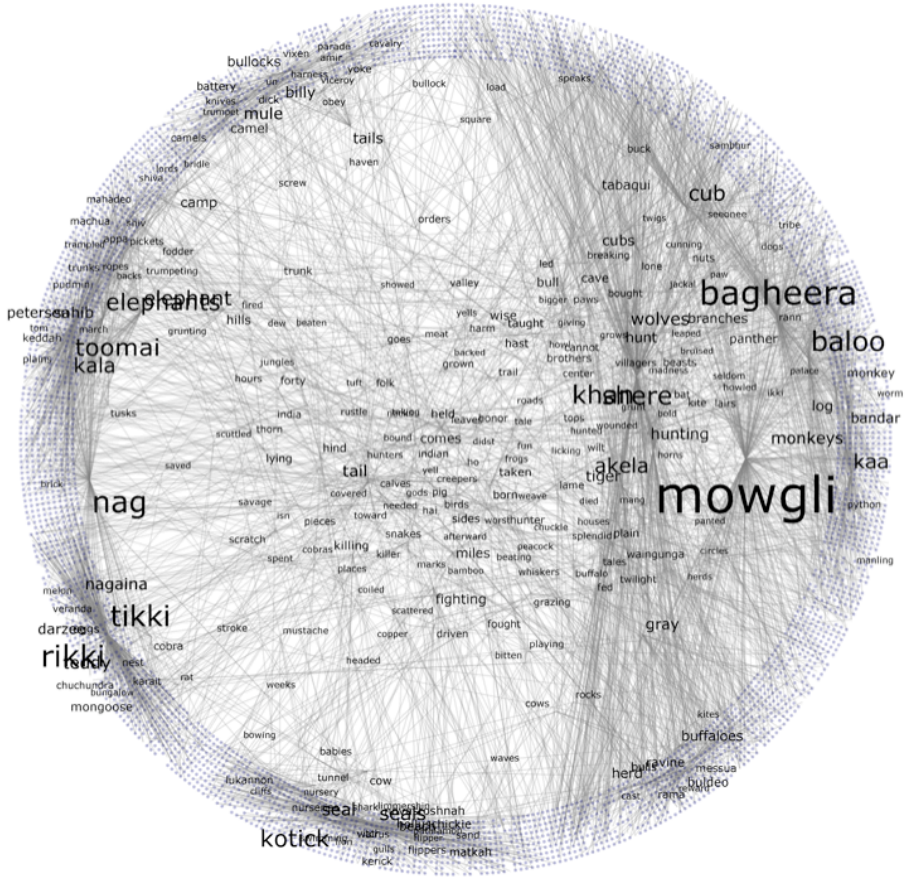
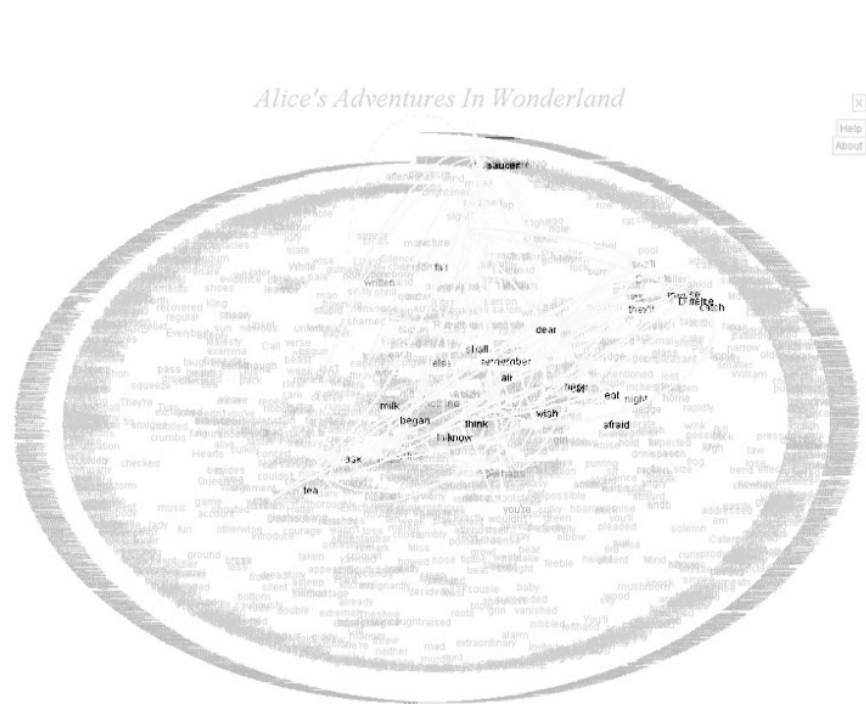
nViZn

a cover tribute to Martin Wattenberg by Graham Wills
left: Wattenberg visualizing music; right: Wills visualizing his own lecture.
Wills programmed his example in nViZn the evening before his talk at the
2008 Joint Statistical Meetings (the day after Wattenberg's talk).



nViZn

a cover tribute to Brad Paley by Graham Wills (drawn by nViZn)
left: Paley visualizing Alice in Wonderland; right: Wills visualizing Jungle Book.
Wills programmed his example in nViZn the evening before his talk at the 2008 Joint Statistical Meetings.



Coherent

- Coherence comes from a grammar that does not depend on surface appearance.
- Taxonomies of visualization are frequently based on surface appearance.
- Taxonomies can be harmful because they mislead programmers, mislead designers, and lead everyone to think the world of graphics is complicated.

Coherent

Taxonomies of charts are harmful.

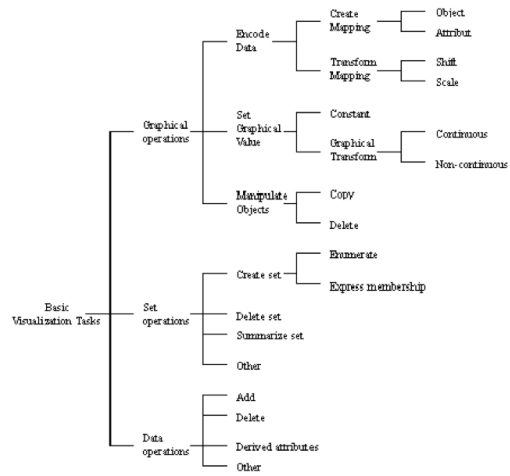
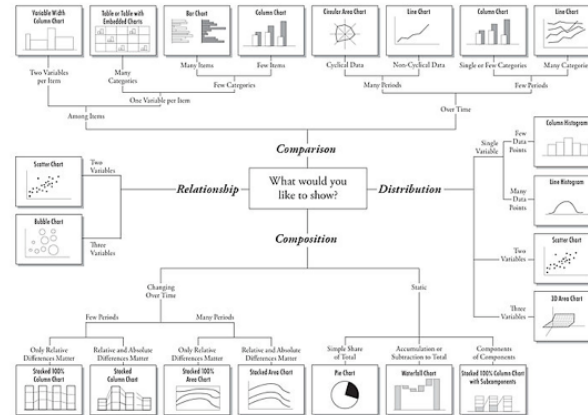
A PERIODIC TABLE OF VISUALIZATION METHODS

Legend:

- Data Visualization:** Visual representation of quantitative data in a structured form (table, map or network map).
- Information Visualization:** The use of computer-based representations of data to aid the user. The means that data is transformed into a display that is related to a user space. The display can be changed by users as they interact with it.
- Concept Visualization:** Methods to represent complex, abstract concepts, ideas, plans, and analysis.
- Strategy Visualization:** Data representation that supports strategic analysis, plans, and optimization of strategic operations.
- Metaphor Visualization:** Visual metaphors provide information primarily as a process and structure information. They also convey an insight about the conceptual structure through the characteristics of the metaphor that is employed.
- Compound Visualization:** The combination of two or different types of representation, such as plan or an image, table or a picture.

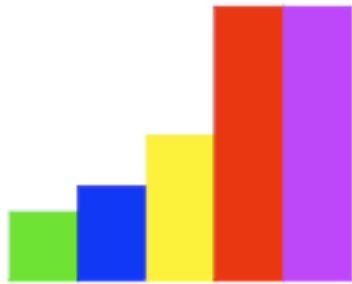
Process Visualization: Cy
Structure Visualization: Hy
Overview: =
Detail AND Overview: ⊕
Divergent thinking: <<
Convergent thinking: >>

Note: Depending on your location and connection speed it can take some time to load a pop-up picture. version 1.5
 © Ralph Langner & Martin J. Spiller www.visual-hierarchy.com



Coherent

- Taxonomies of charts are harmful because...



this

is similar to



this

Coherent

- Taxonomies of charts are harmful because...



this

is similar to

this

Coherent

- Taxonomies of charts are harmful because...



but this

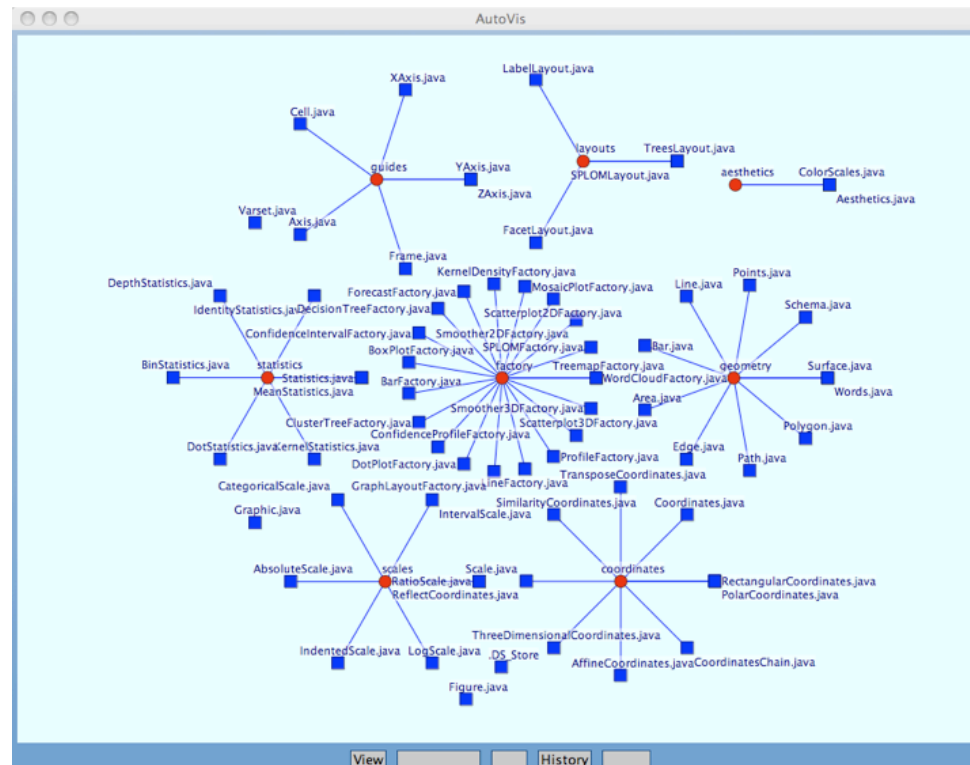
is *not* similar to



this

Coherent

- Here is a graph layout of GoG Java packages.
- The Factory classes in the center contain automation methods for building graphics from the GoG surrounding classes.
- These factory methods all have a single coherent form.
- The next slide shows an example.



Coherent

- A treemap factory:



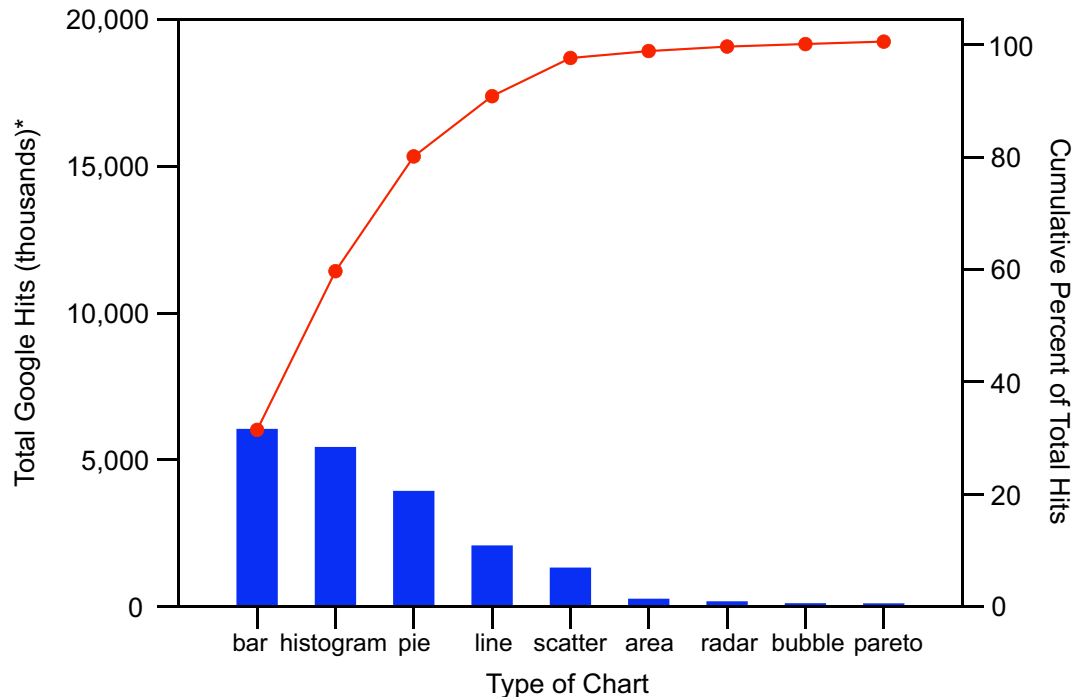
```
public abstract class TreemapFactory {
    public static Graphic build(FlatTable table, GraphicsStyles styles) {
        Renderer go = OutputFactory.makeBitmapGraphicsOutput(styles);
        Graphic graphic = new Graphic(go, styles, null);
        Scale[] scales = new Scale[]{new IntervalScale(table, 0)};
        Varset varset = new Varset(scales, table);
        Aesthetics aesthetics = new Aesthetics(varset);
        aesthetics.setPosition(new int[0]);
        aesthetics.setColor(new int[]{0});
        CoordinatesChain chain = new CoordinatesChain();
        Graphic branch = new Graphic(chain, graphic);
        new Polygon(aesthetics, branch, true);
        new Frame(aesthetics, branch);
        return graphic;
    }
}
```

Meaningful

- A meaningful language is based on identifying statements that are falsifiable.
- Graphics libraries and toolkits do not have this property.
- GoG asserts some popular charts are meaningless because they are ill-formed.
- To call these charts meaningful, defenders must falsify specific assumptions of GoG.
- To do so would lead to characterizing most charts as meaningless.

Meaningful

- GoG tells us the Pareto Chart is ill-formed



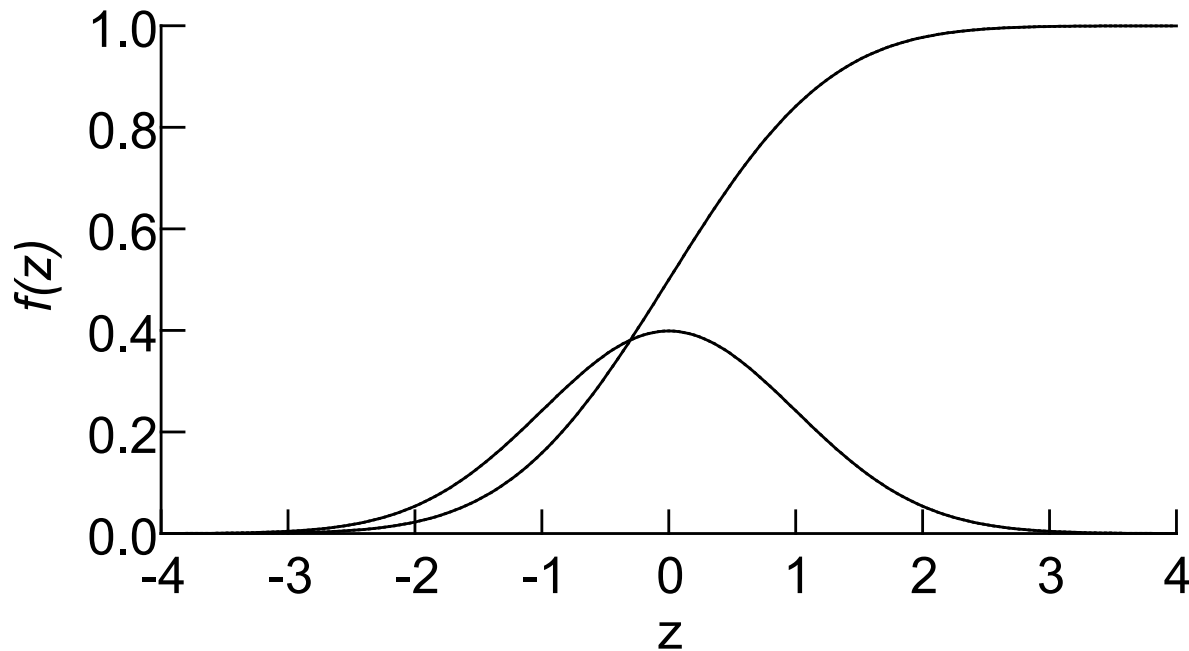
*queries = "<type>" + {" chart", " graph", " plot"} plus {"scatterplot", "histogram"}

This is a Pareto chart of charts.

It is based on total hits in Google searches (using an initial inventory derived from Google Sets).
On this basis, the Pareto chart is the ninth most popular chart in the world today.

Meaningful

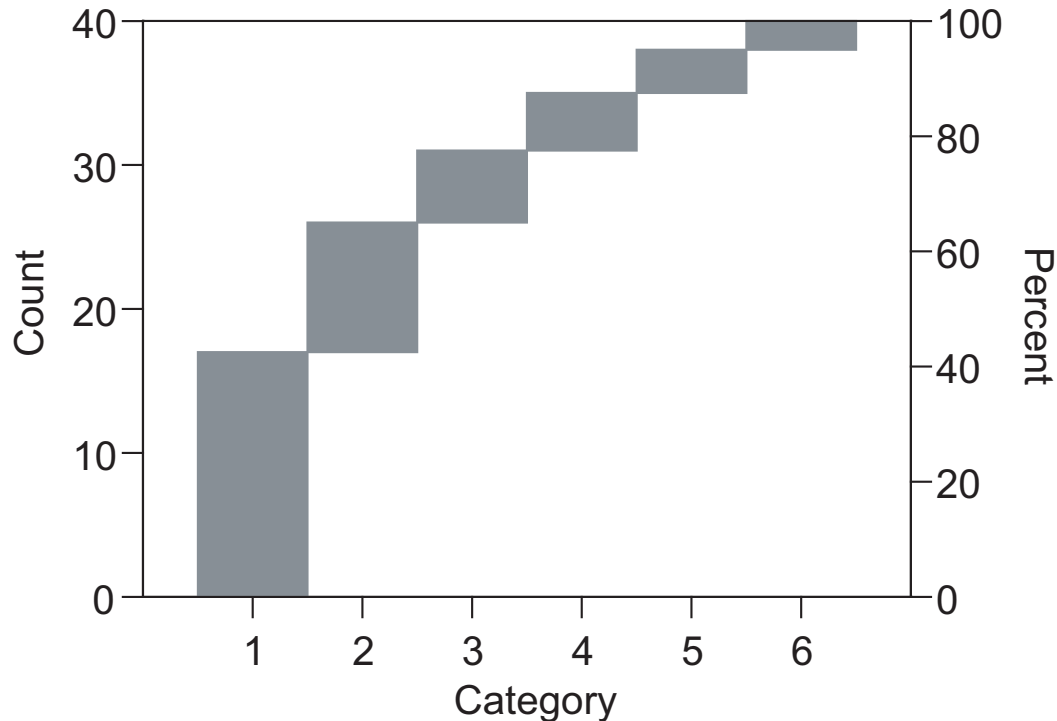
- GoG tells us this chart is ill-formed for the same reason.



The Pareto chart and this graph plot a density and cumulative distribution in the same frame. It is nonsensical to plot a density and cumulative distribution in the same frame. For the similar reasons, it is nonsensical to plot velocity and acceleration in the same frame.

Meaningful

Stacked Pareto Chart (Juran, 1951)



Juran invented the Pareto chart. This was his preferred version, however. It is meaningful because the vertical scale is cumulative. Unfortunately, the world chose the meaningless version. Read more in Wilkinson, *The American Statistician*, 2006.

Intelligent Graphics Systems

- AutoVis
- FASTAT

What is automated visualization?

- The computer looks at raw data and produces visualization(s) without coaching.
- The computer visually highlights aspects of data that are noteworthy.
- The computer presents data visualizations that help evaluate the appropriateness of a model.

What is the point of automated visualization?

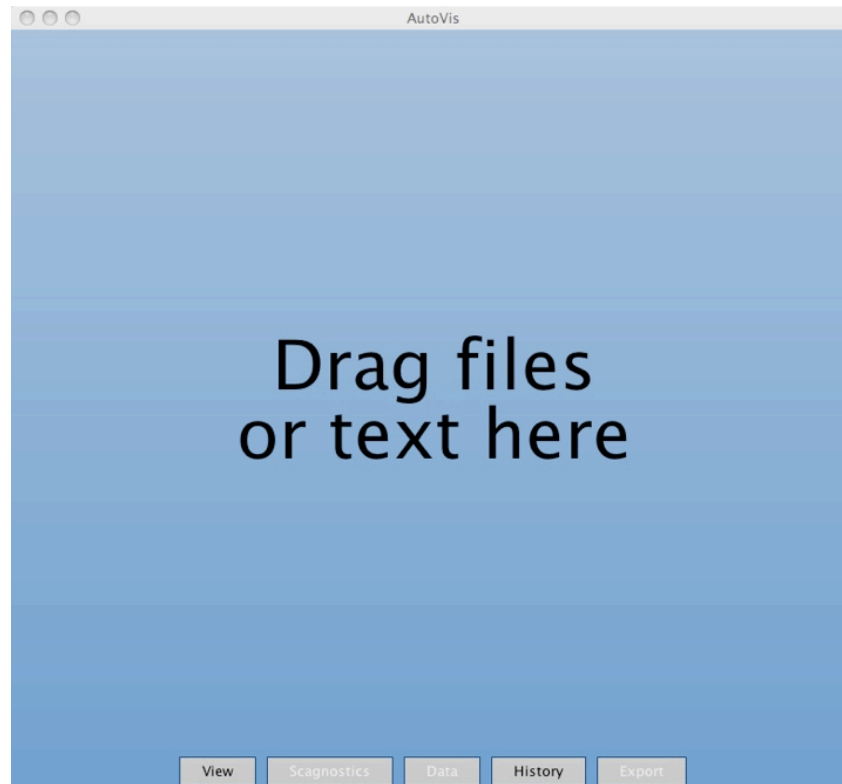
- **Guided discovery**
 - Data cleaning and catching anomalies
 - Checking assumptions
 - Controlling False Discovery Rate
- **Data mining and machine learning**
 - Classification
- **NOT** for deciding whether to buy Dubai real-estate or for submitting a paper to *Science*.

AutoVis

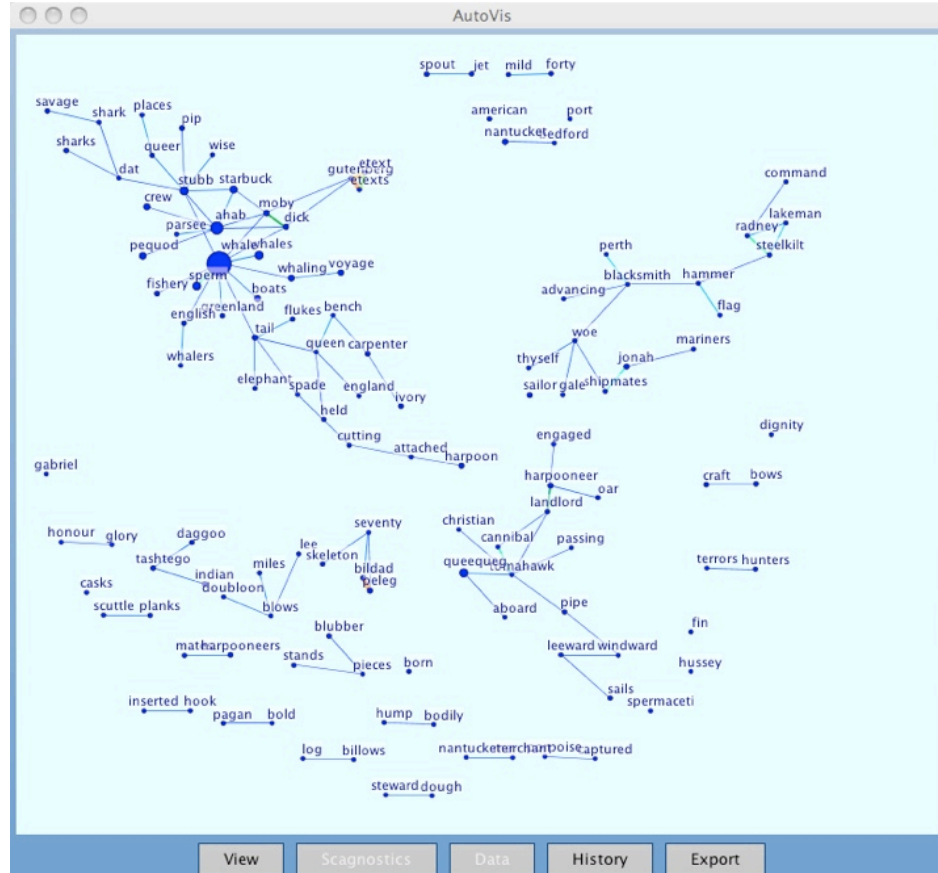
- A data viewer
 - Drag a file into window
 - Graphics appropriate to file structure
 - Uses scagnostics to find interesting views
 - Controls False Discovery Rate to prevent fishing expeditions

AutoVis

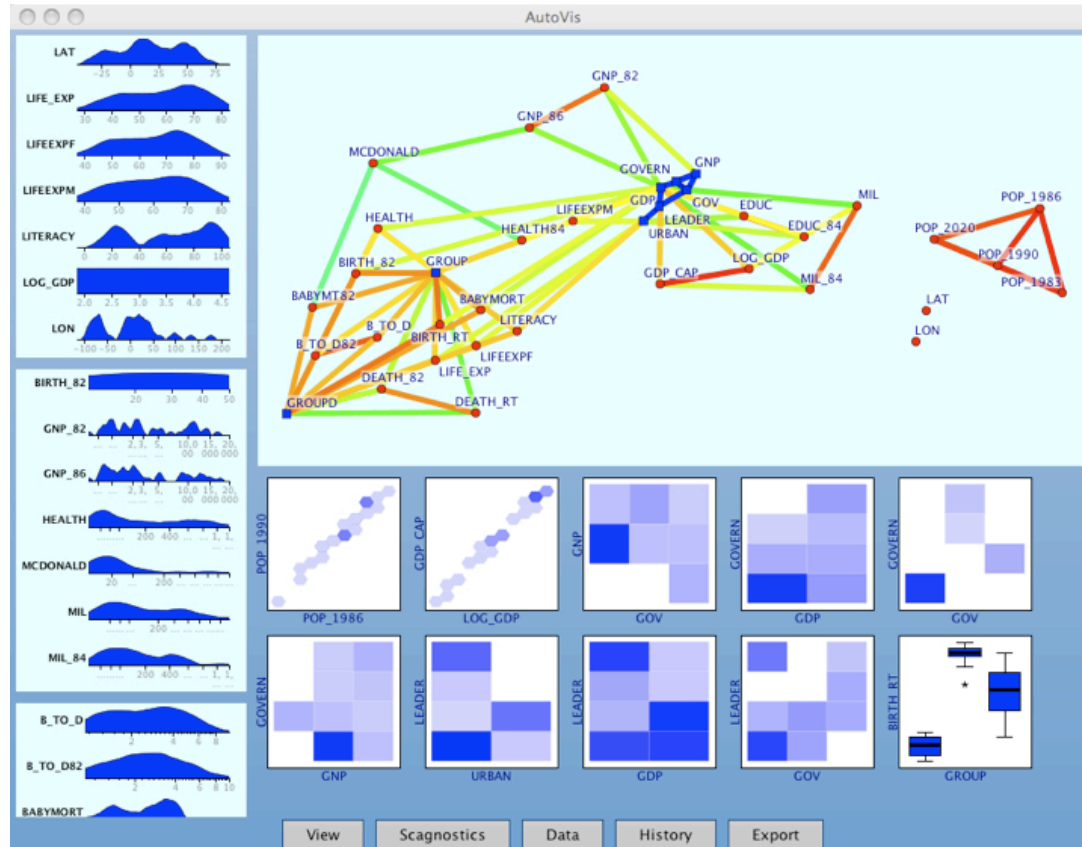
AutoVis has no menus or user manual. Grab a file from your local directory or from the Web and drag it in the window. The following slides show typical results.



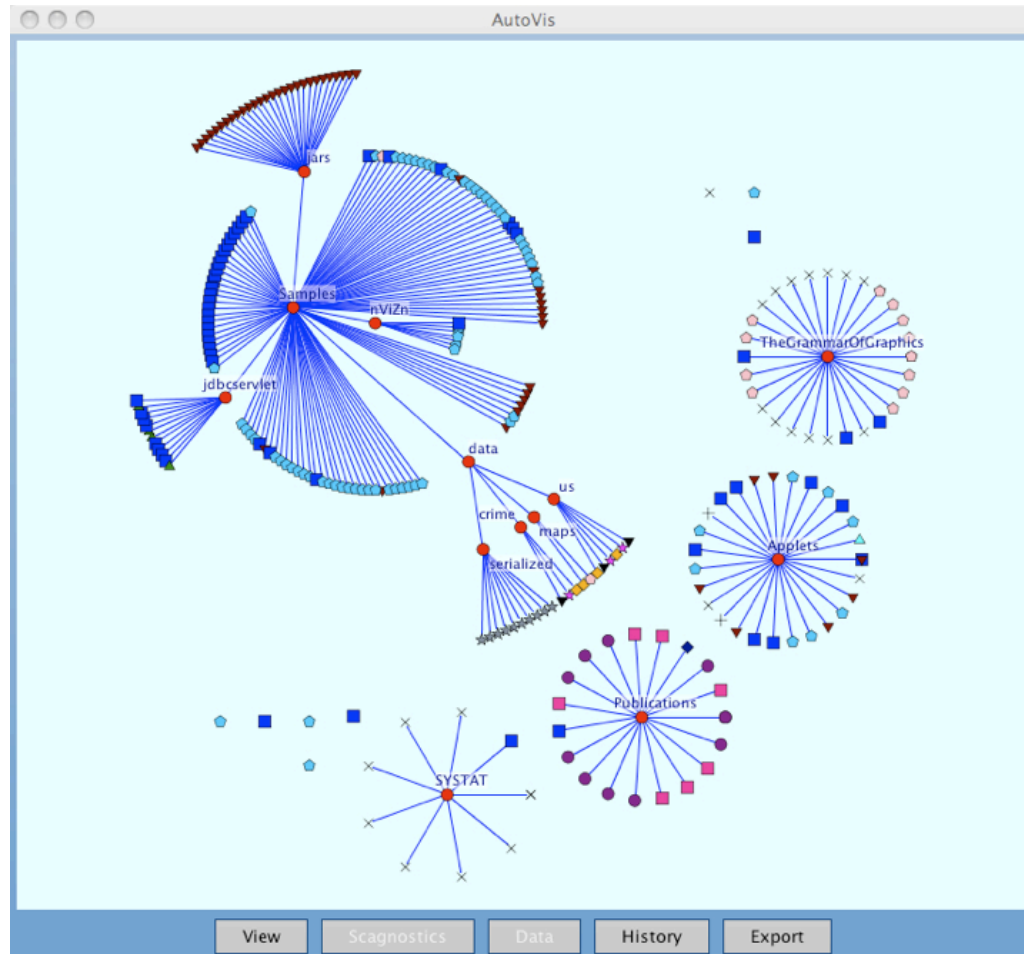
AutoVis



AutoVis



AutoVis



FASTAT

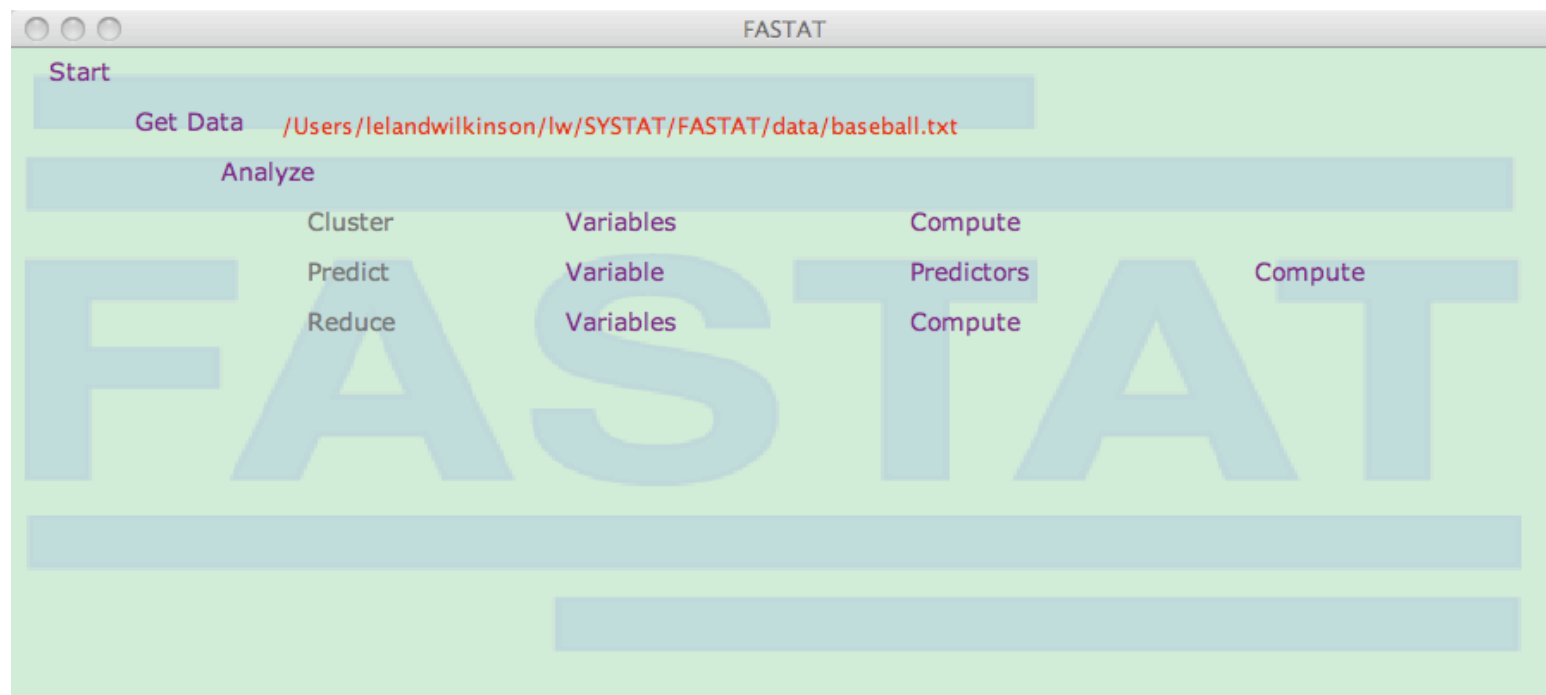
- A second opinion:
 - Look for outliers, distributions, coding errors...
 - Transform to meet assumptions
 - Converse with user
 - Produce annotated, hyperlinked output

FASTAT

FASTAT has no dialogs or user manual. There is no keyboard input. The main window fills up as the user chooses options by touching the purple words, as the following slide shows.

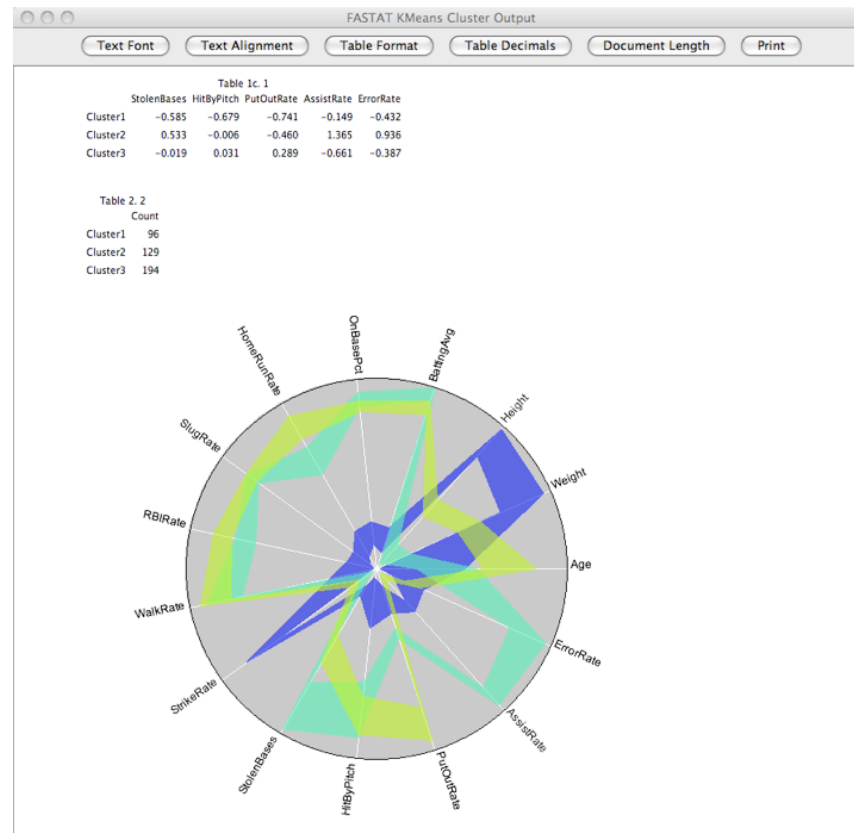


FASTAT

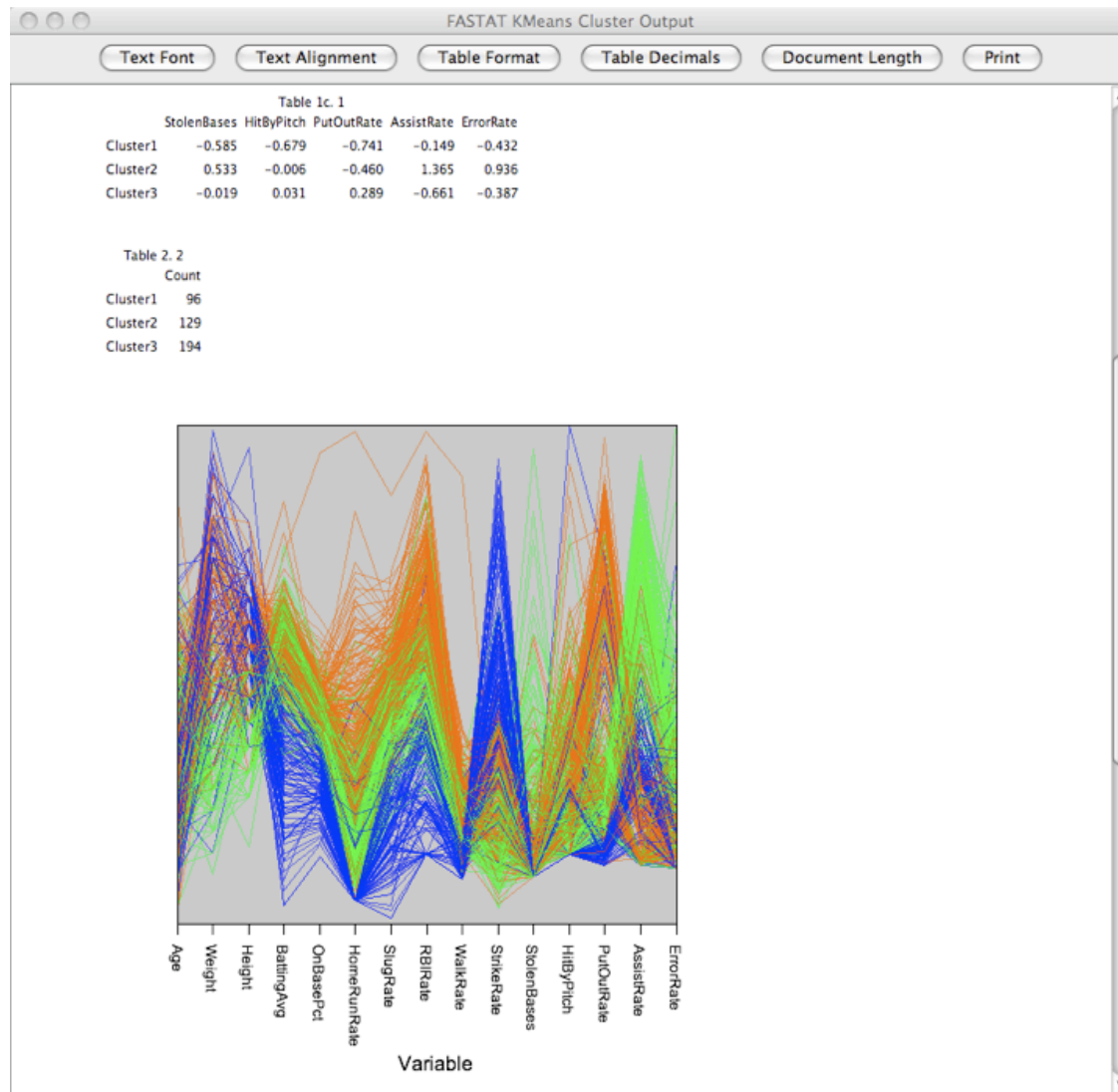


FASTAT

The output is a browser that contains graphics and text. Touching a graphic produces a popup that allows the user to change the form of the graphic in one tap, as the following slide shows.



FASTAT



Conclusion

To build *intelligent* visual analytics, we need to understand the mathematical foundation of analytic visualizations.