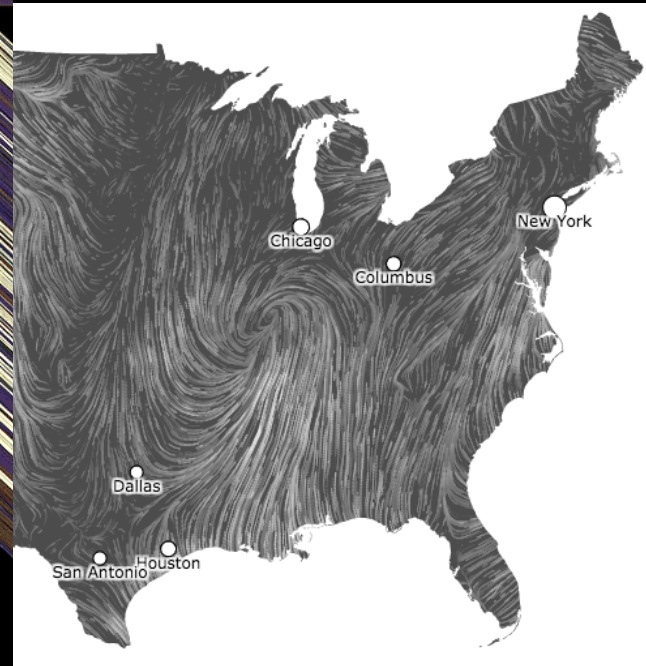
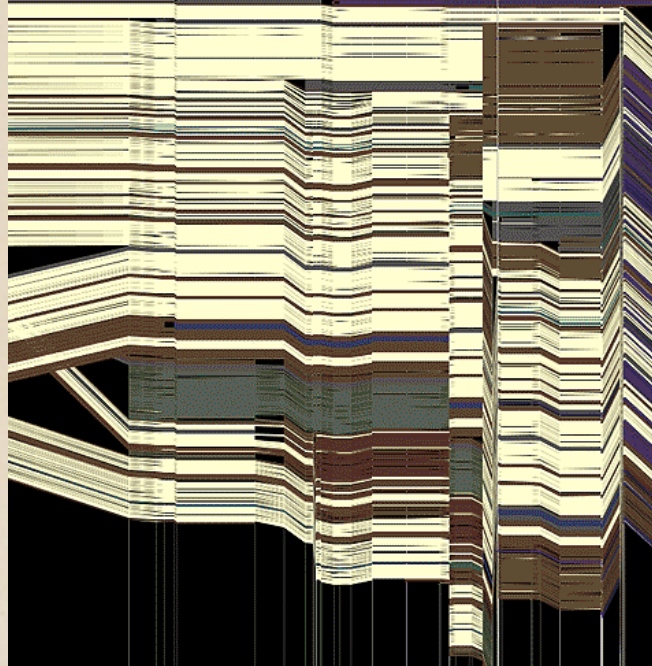
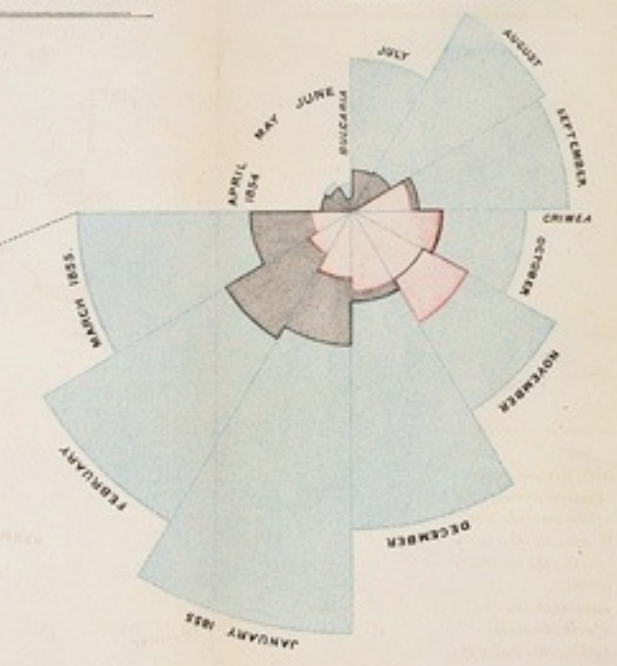


CSE 442 - Data Visualization

Visual Encoding Design



Matthew Conlen University of Washington

Re-Design Exercise

Re-Design Exercise

Task: Analyze and Re-design visualization

Identify data variables (N/O/Q) and encodings

Critique the design: what works, what doesn't

Sketch a re-design to improve communication

Be ready to share your thoughts with the class

Break into groups with those sitting near you

(~4 people per group)

Effectiveness Rankings [Mackinlay 86]

QUANTITATIVE

Position
Length
Angle
Slope
Area (Size)
Volume
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Shape

ORDINAL

Position
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Length
Angle
Slope
Area (Size)
Volume
Shape

NOMINAL

Position
Color Hue
Texture
Connection
Containment
Density (Value)
Color Sat
Shape
Length
Angle
Slope
Area
Volume

Teacher Salaries: Is It Really That Bad?

National and State averages for K-12 Public-School Teachers



UNITED STATES

AVG. SALARY: \$47,674

Avg. vacation days: 63

HOURLY

Hours per week on-site: 36.5
 Public-School Teacher: \$34.06
 Private-School Teacher: \$21.08
 Average Worker: \$25.08
 Police: \$22.64
 Fire: \$17.91

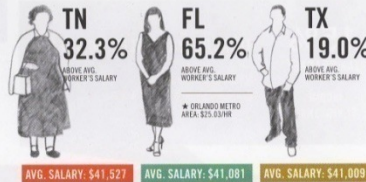
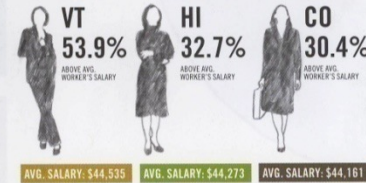
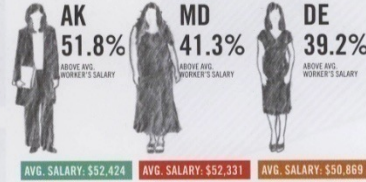
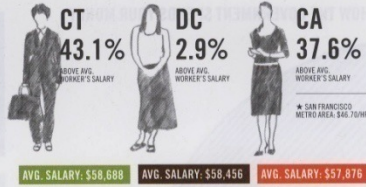


CANADA

AVG. SALARY: \$43,000

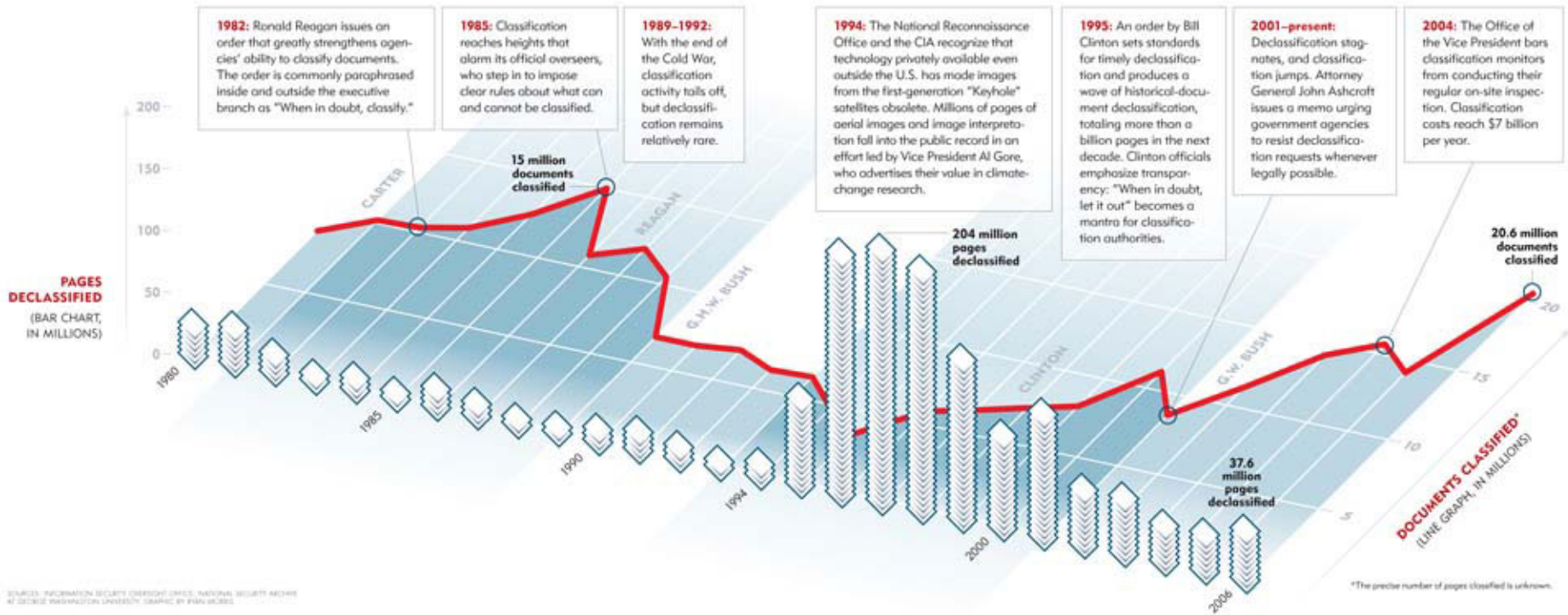
Avg. vacation days: 50

HOURLY: \$30.18
 Hours per week on-site: 55.6

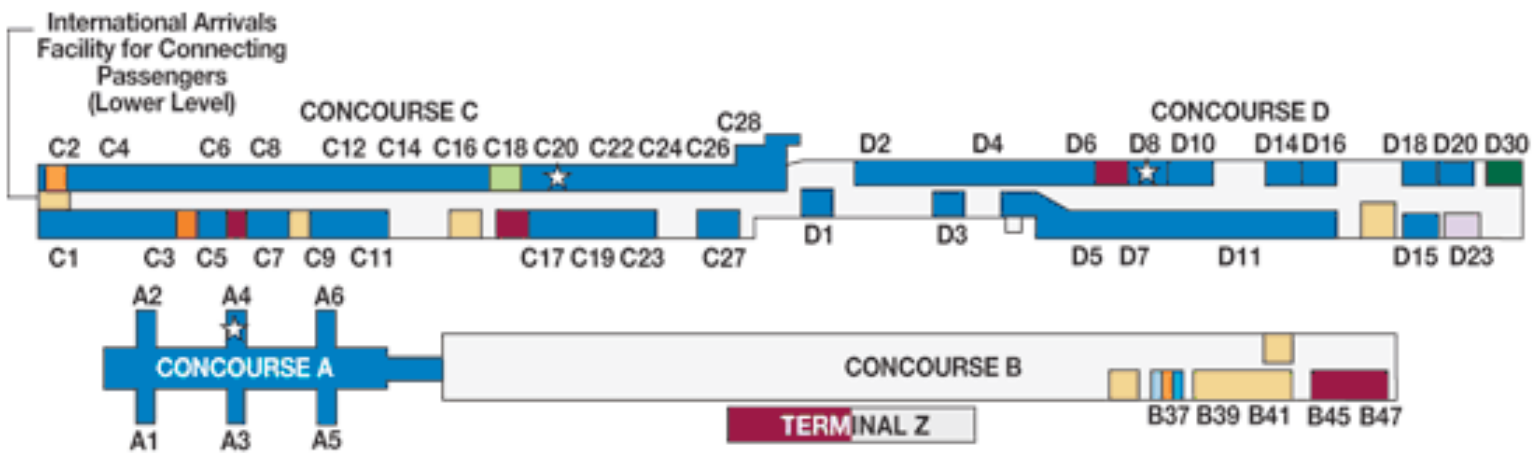


28
GOOD May/June 07
Transparency
SOURCES Manhattan Institute; National Center For Education Statistics; National Education Association; U.S. Bureau of Labor Statistics

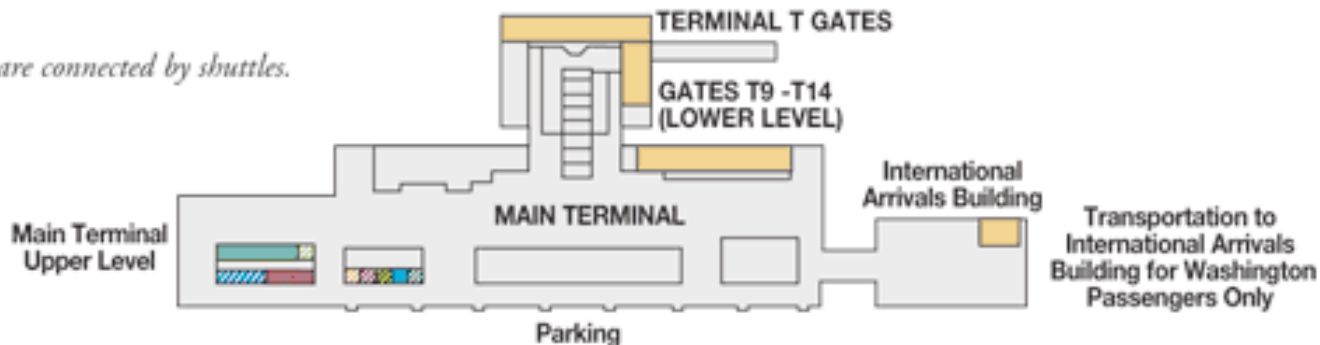
AVERAGE Workers' salaries used for comparison are those of white-collar, nonsales employees.



Source: *The Atlantic* 300 no. 2 (September 2007)
 Number of Classified U.S. Documents



* Terminals are connected by shuttles.



- | | | |
|-----------------------------------|----------------------------|-----------------------------|
| United / TED Gate Area | Lufthansa Check-in | Austrian Airlines Gate Area |
| United Premier Check-in | Air Canada Gate Area | SAS Gate Area |
| United Check-in | Air Canada Check-in | BWIA Gate |
| United International Check-in | Mobile Lounge Dock | South African Airways |
| United Red Carpet Club | ANA Check-in | US Airways Gates |
| United First International Lounge | ANA Fuji Lounge/Gate Area | United EasyCheck-in |
| Lufthansa Gate Area | Austrian Airlines Check-in | US Airways Check-in |

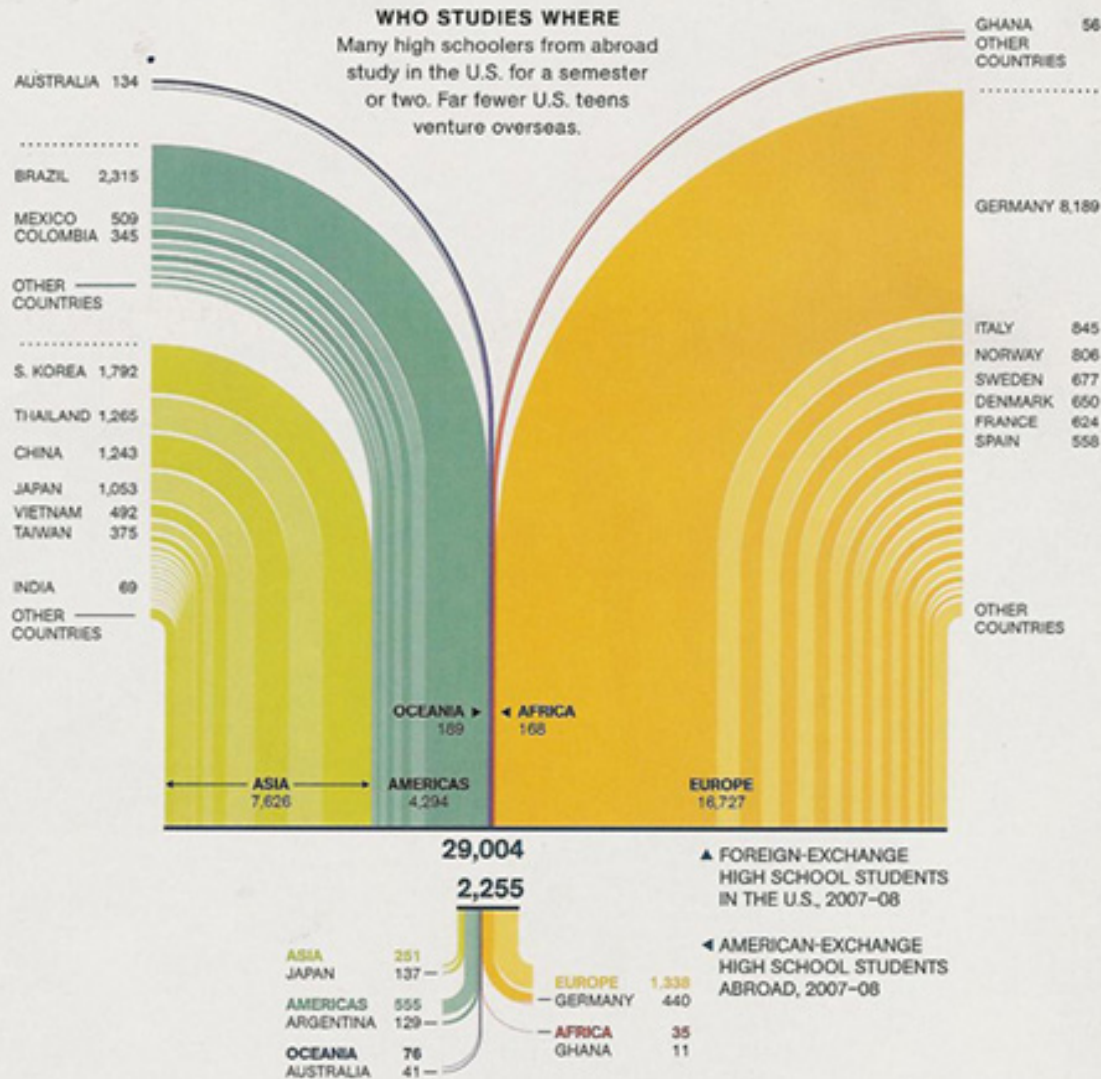
EasyCheck-in is available at this airport.



11 2006

Washington Dulles Airport Map

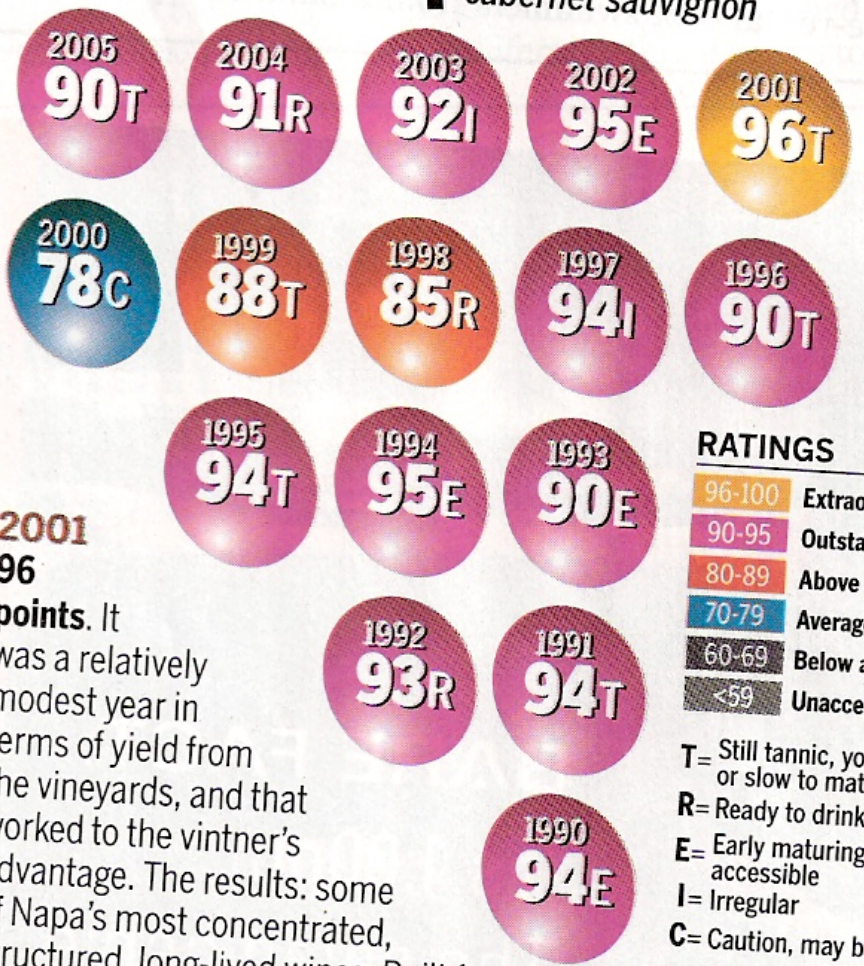
Source: United Airlines Hemispheres



Source: *National Geographic*, September, 2008, p. 22.
Silver, Mark. "High School Give-and-Take."

IT WAS A VERY GOOD YEAR?

Robert Parker's ratings for vintages of Napa Valley cabernet sauvignon



2001
96
points. It was a relatively modest year in terms of yield from the vineyards, and that worked to the vintner's advantage. The results: some of Napa's most concentrated, structured, long-lived wines. Built for aging, they are rich, densely colored, fruity and

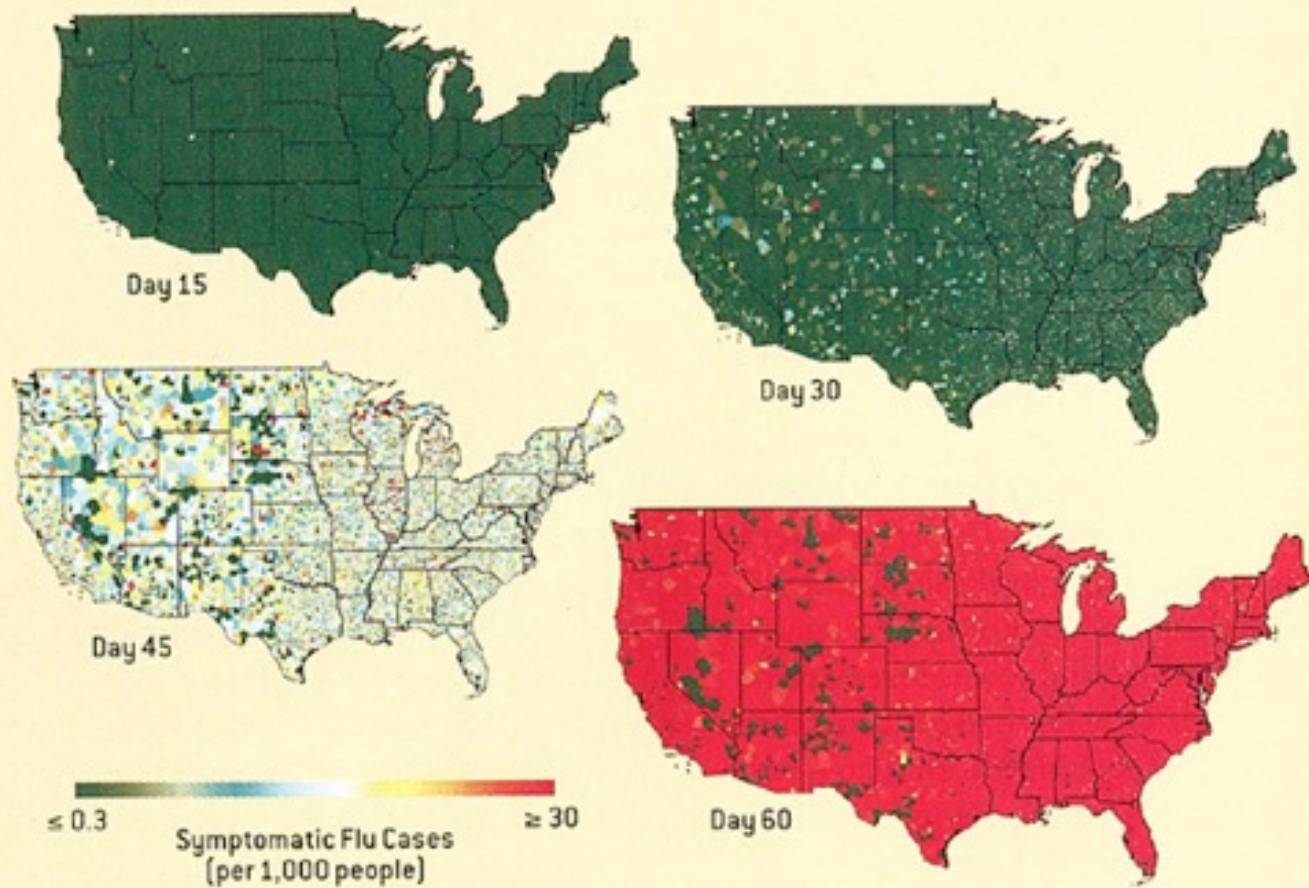
RATINGS

96-100	Extraordinary
90-95	Outstanding
80-89	Above average
70-79	Average
60-69	Below average
<59	Unacceptable

T= Still tannic, youthful, or slow to mature
 R= Ready to drink
 E= Early maturing and accessible
 I= Irregular
 C= Caution, may be too old

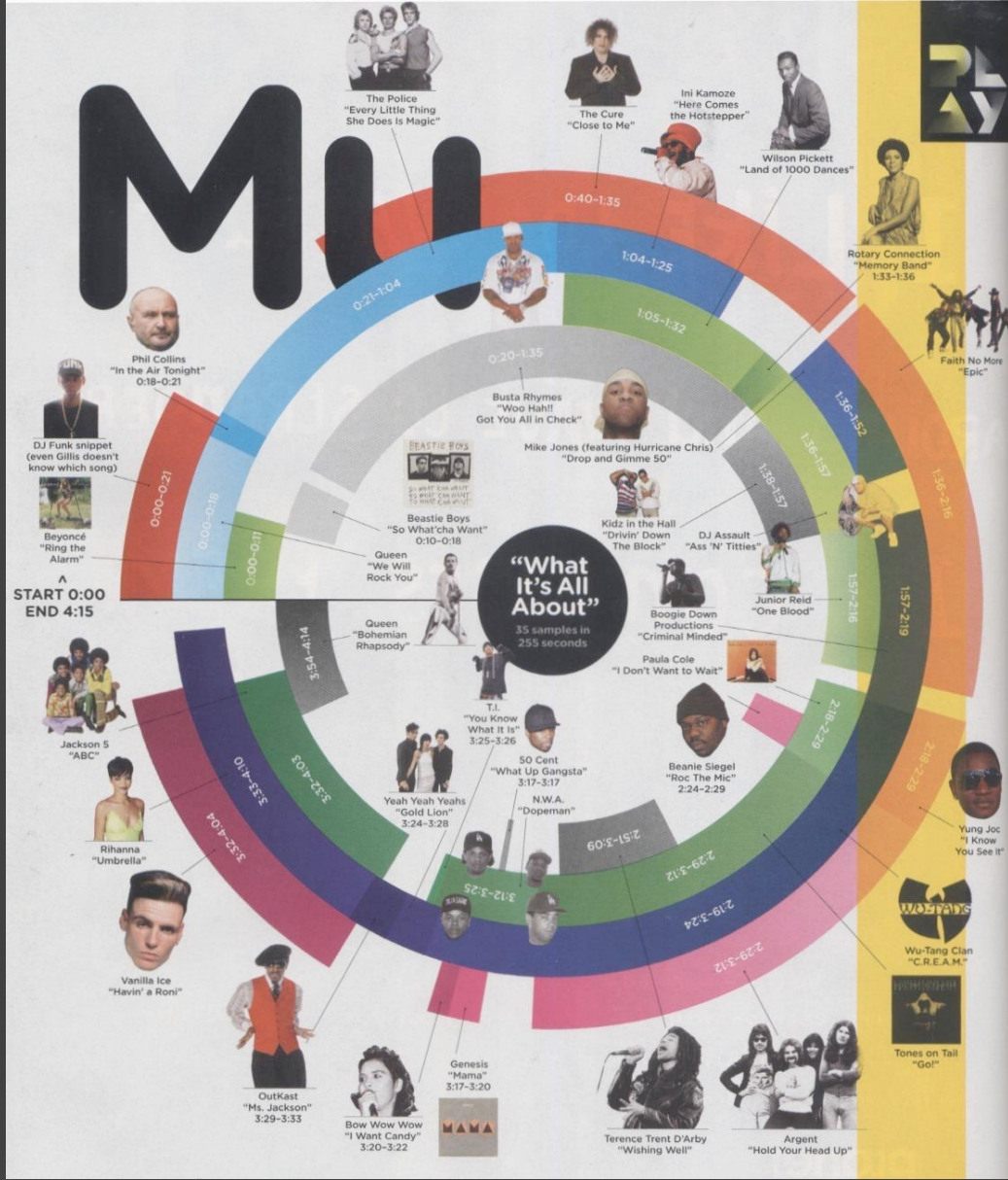
Pandemic Flu Hits the U.S.

A simulation created by researchers from Los Alamos National Laboratory and Emory University shows the first wave of a pandemic spreading rapidly with no vaccine or antiviral drugs employed to slow it down. Colors represent the number of symptomatic flu cases per 1,000 people (see scale). Starting with 40 infected people on the first day, nationwide cases peak around day 60, and the wave subsides after four months with 33 percent of the population having become sick. The scientists are also modeling potential interventions with drugs and vaccines to learn if travel restrictions, quarantines and other disruptive disease-control strategies could be avoided.



Preparing for a Pandemic

Source: *Scientific American*, 293(5). November, 2005, p. 50



Source: *Wired Magazine*, September 2008 Edition
Music: Super Cuts (page 92)

A Design Space of Visual Encodings

Mapping Data to Visual Variables

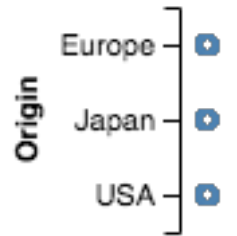
Assign **data fields** (e.g., with N , O , Q types) to **visual channels** (x , y , $color$, $shape$, $size$, ...) for a chosen **graphical mark** type ($point$, bar , $line$, ...).

Additional concerns include choosing appropriate **encoding parameters** ($log\ scale$, $sorting$, ...) and **data transformations** (bin , $group$, $aggregate$, ...).

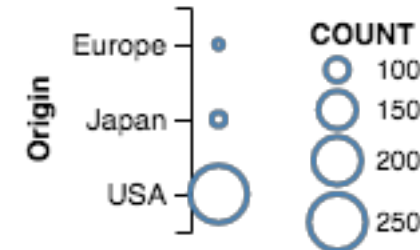
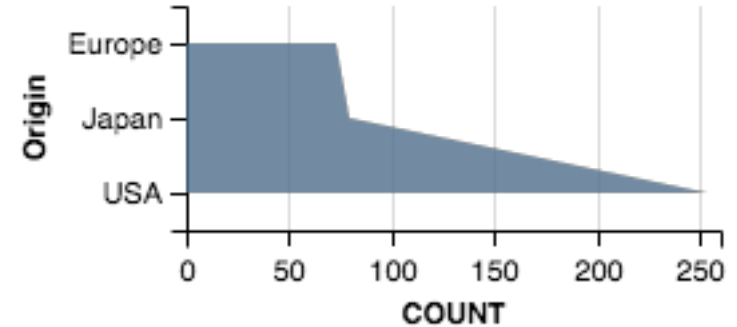
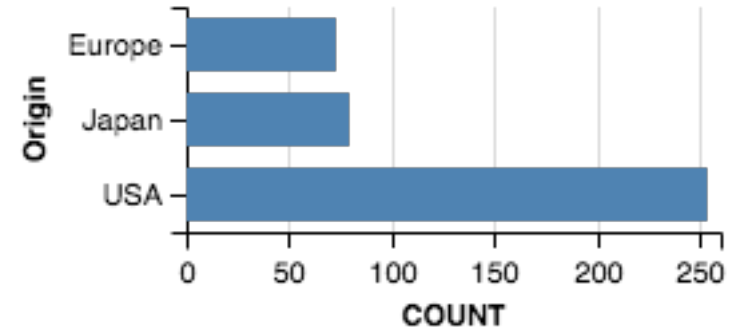
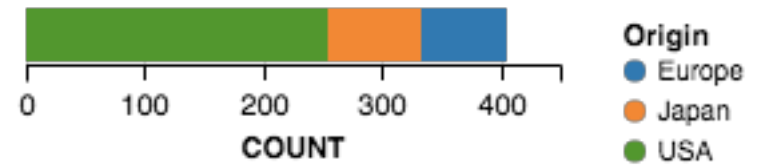
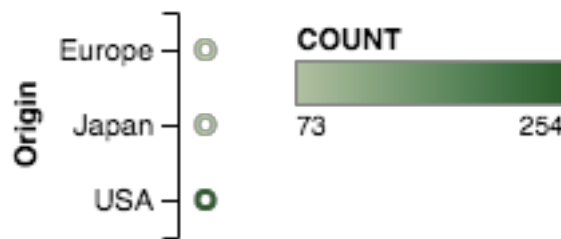
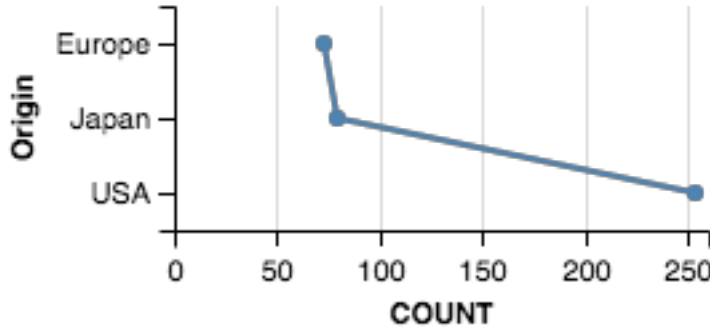
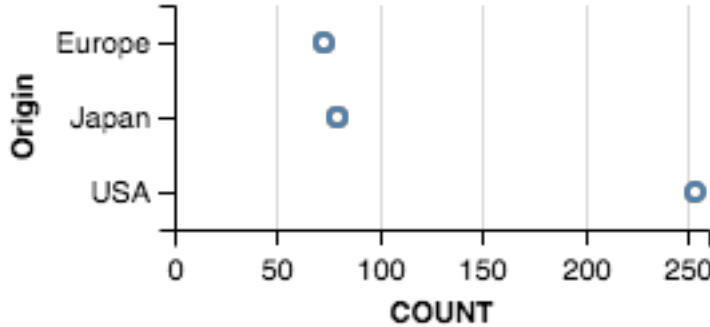
These options define a large combinatorial space, containing both useful and questionable charts!

1D: Nominal

Raw



Aggregate (Count)

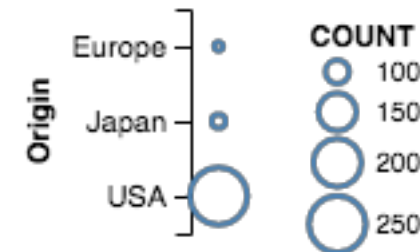
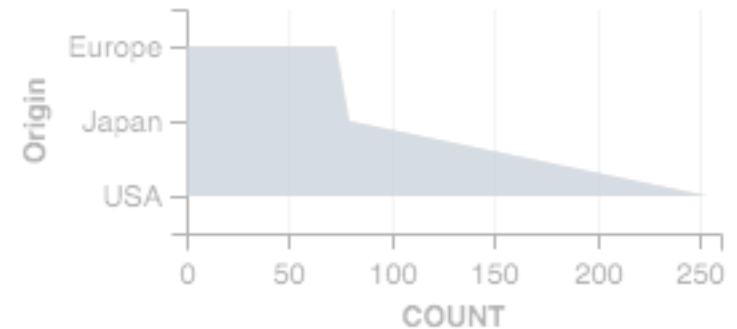
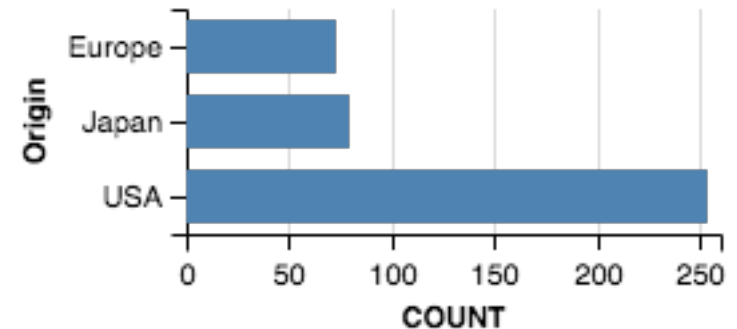
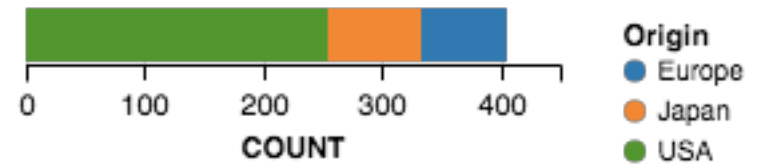
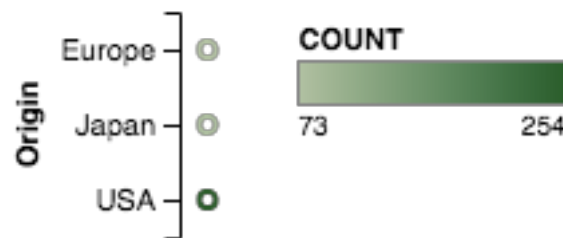
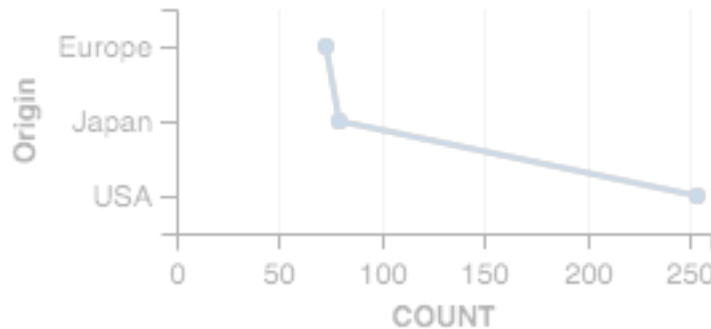
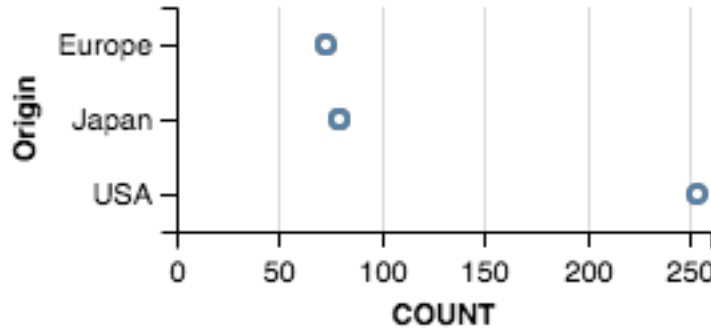


Expressive?

Raw

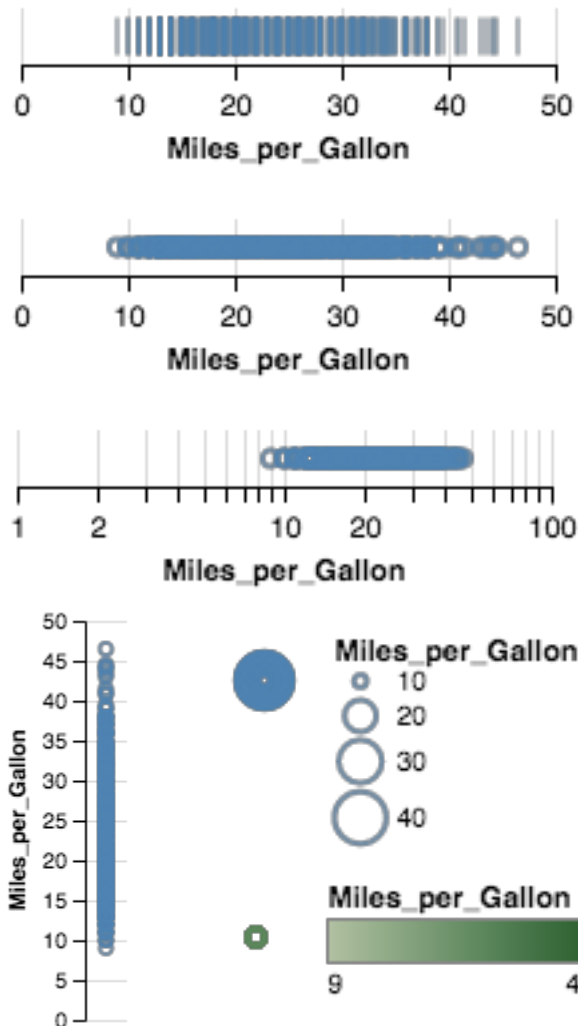


Aggregate (Count)

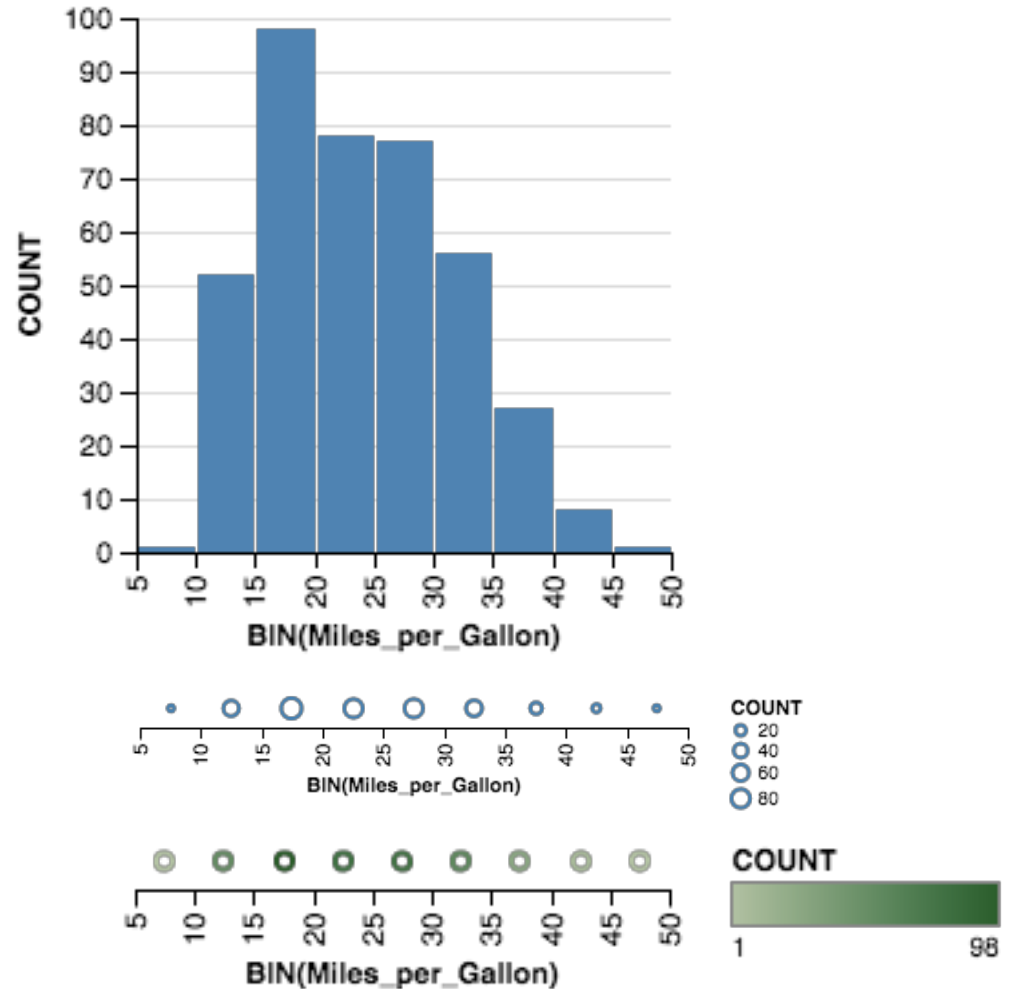


1D: Quantitative

Raw

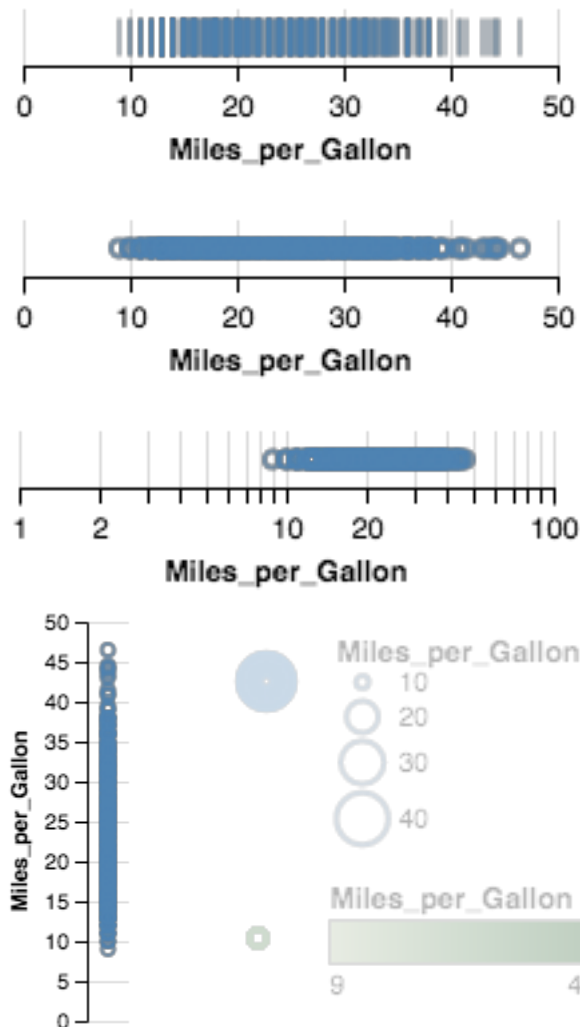


Aggregate (Count)

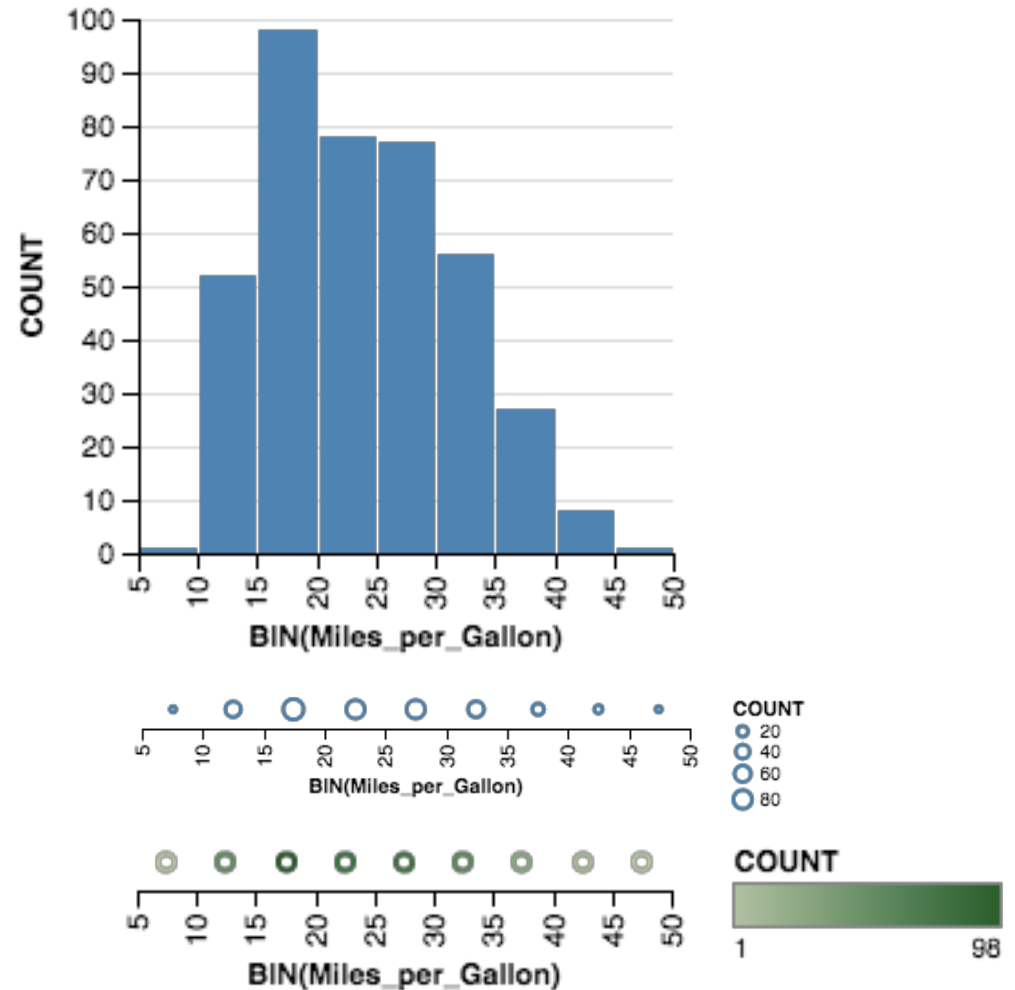


Expressive?

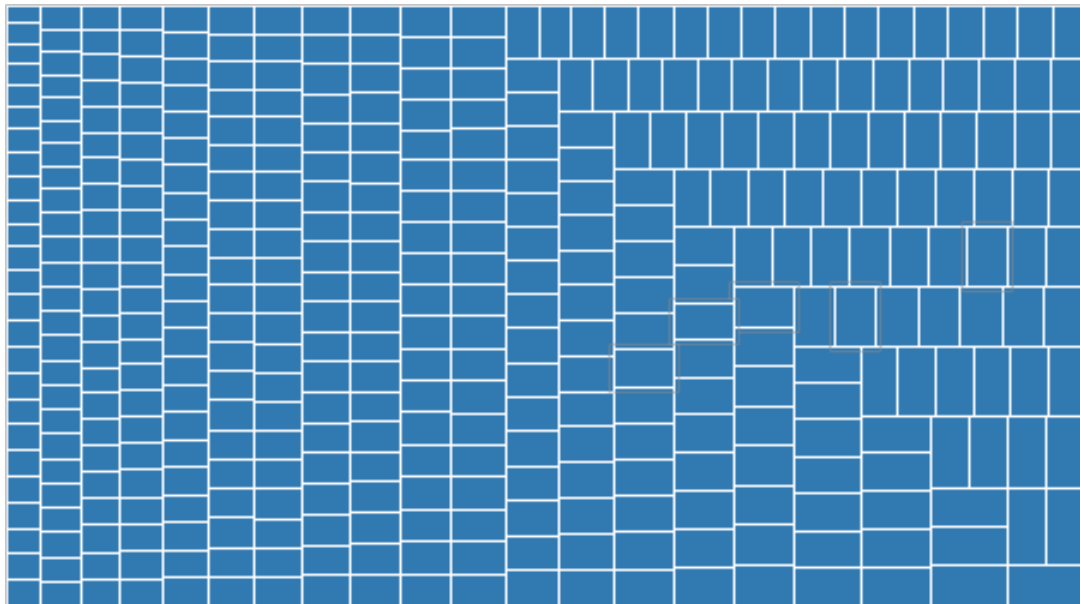
Raw



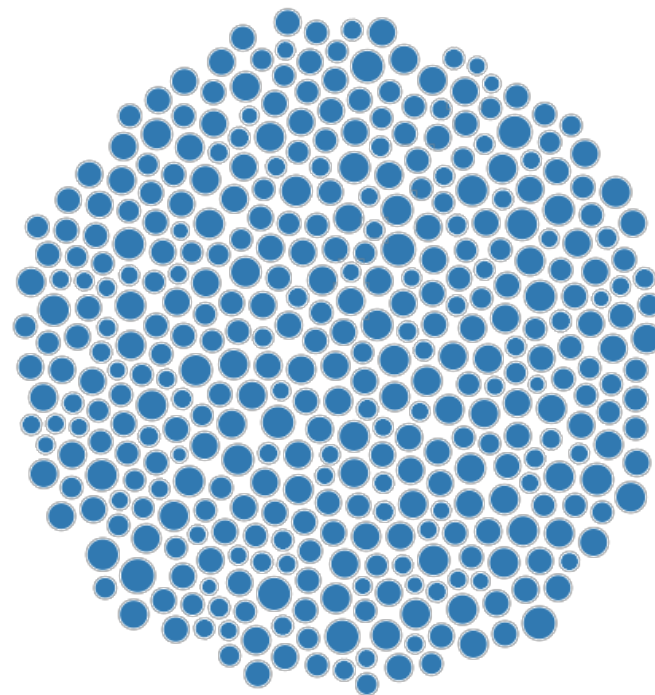
Aggregate (Count)



Raw (with Layout Algorithm)

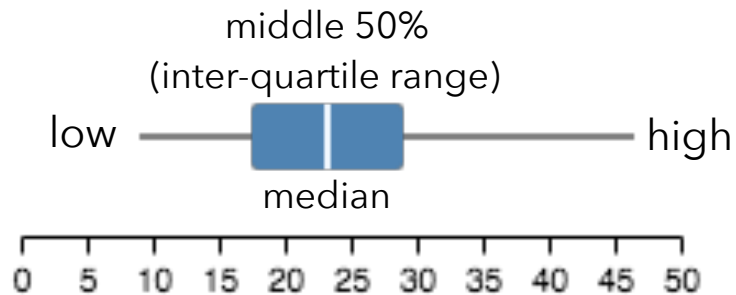


Treemap

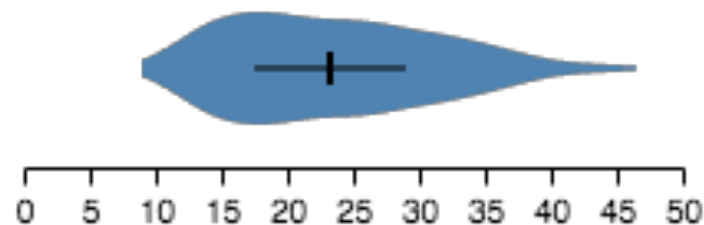


Bubble Chart

Aggregate (Distributions)



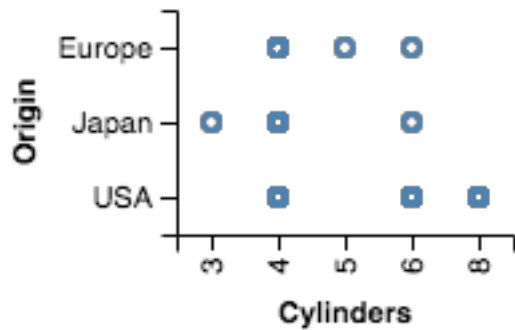
Box Plot



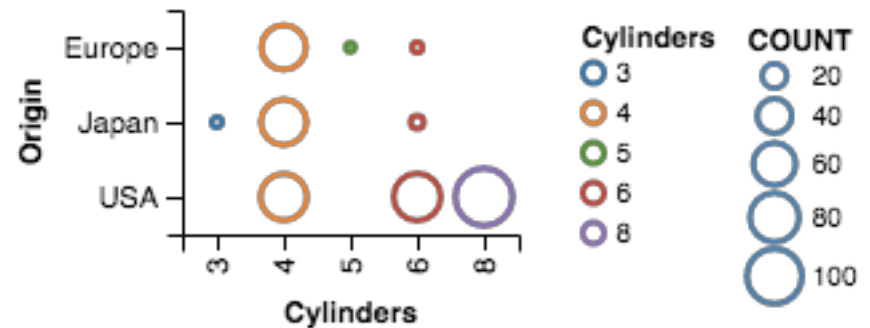
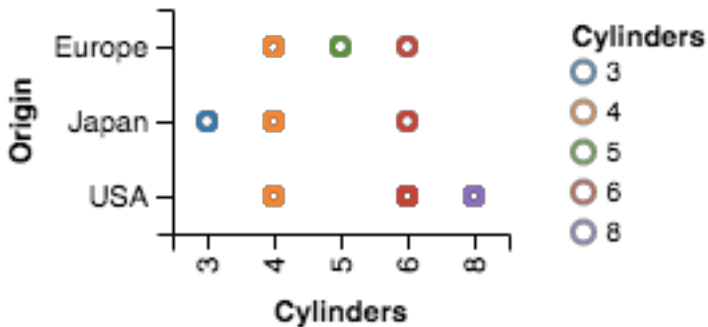
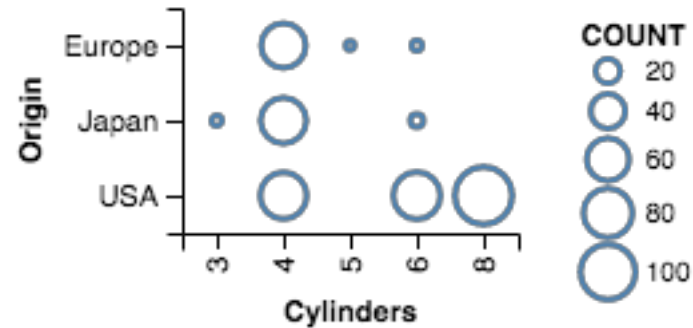
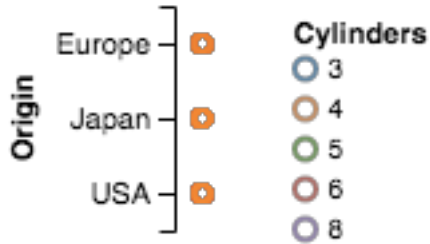
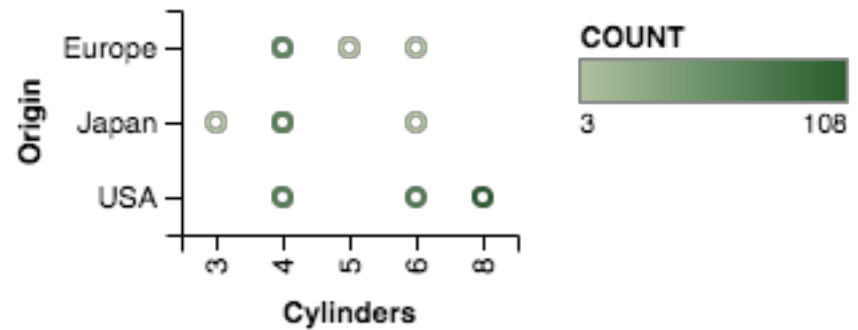
Violin Plot

2D: Nominal x Nominal

Raw

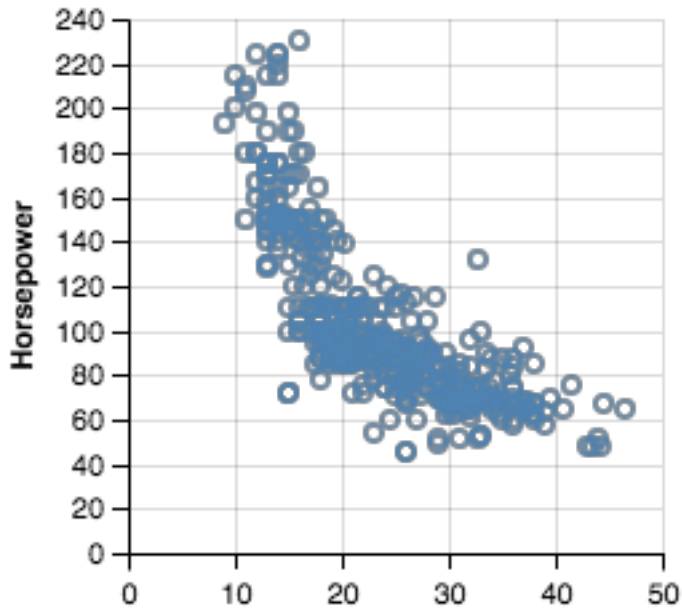


Aggregate (Count)

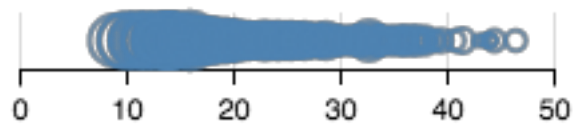


2D: Quantitative x Quantitative

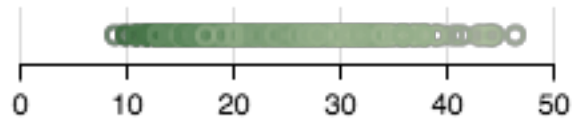
Raw



Miles_per_Gallon



Miles_per_Gallon



Miles_per_Gallon

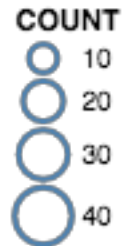
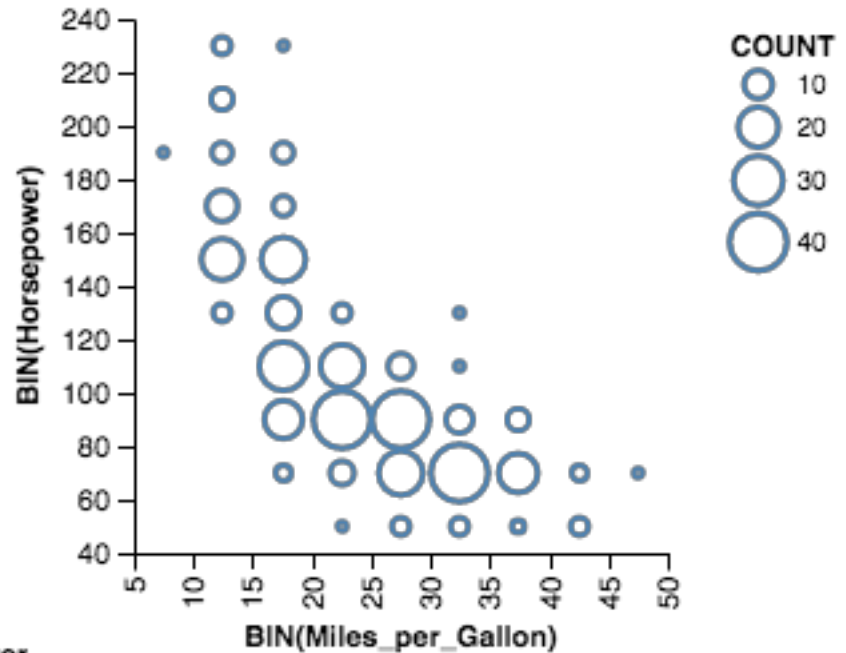
Horsepower



Horsepower

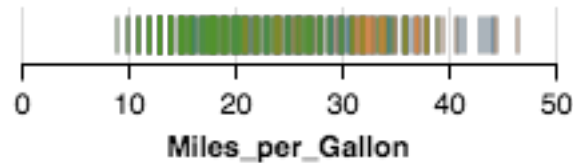
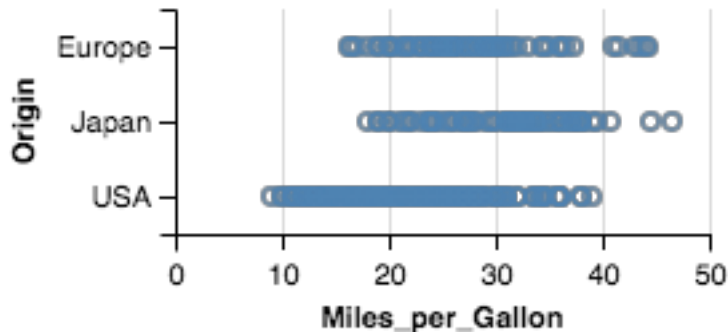


Aggregate (Count)

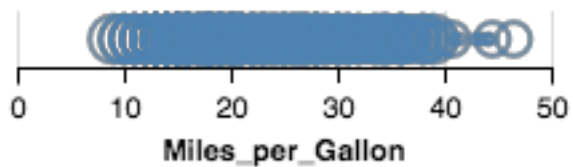


2D: Nominal x Quantitative

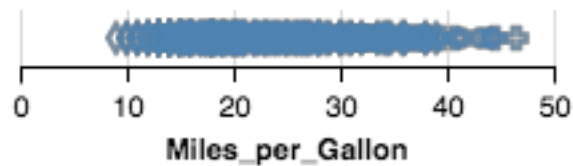
Raw



Origin
● Europe
● Japan
● USA

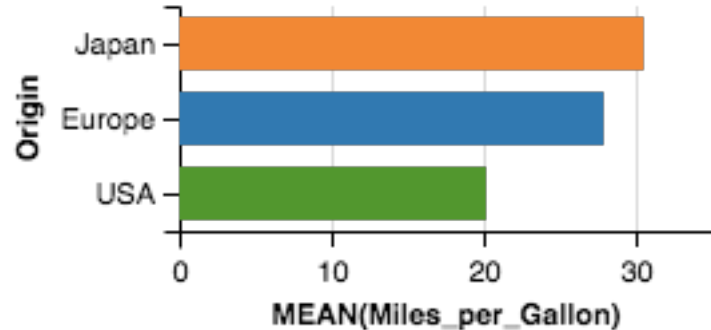
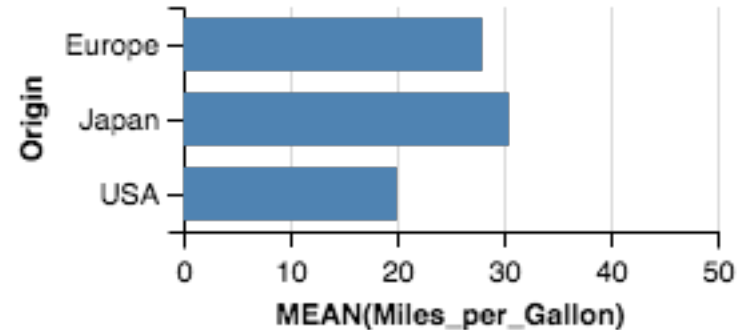


Origin
○ Europe
○ Japan
○ USA



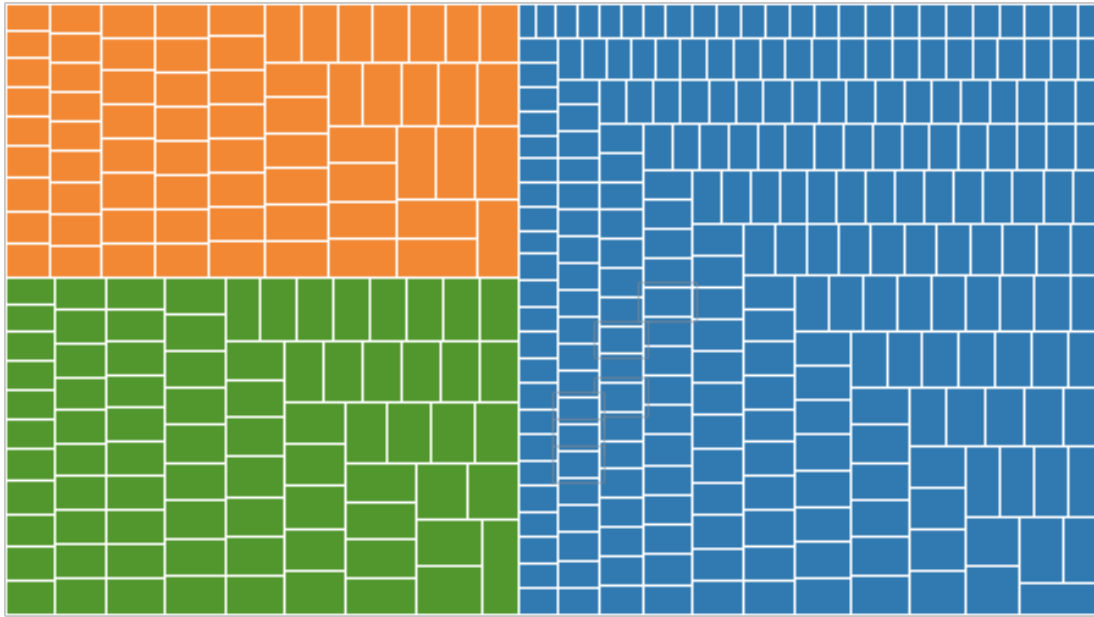
Origin
○ Europe
+ Japan
◇ USA

Aggregate (Mean)

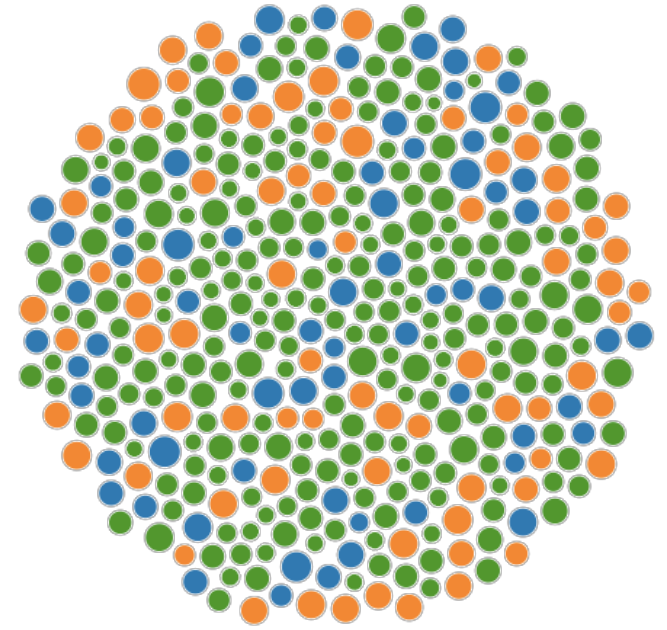


Origin
● Europe
● Japan
● USA

Raw (with Layout Algorithm)

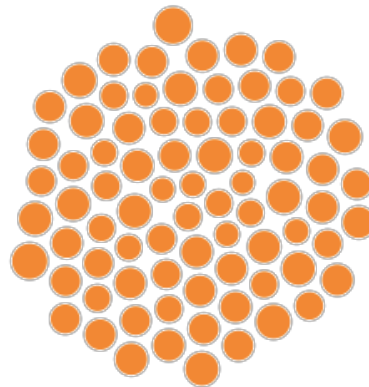
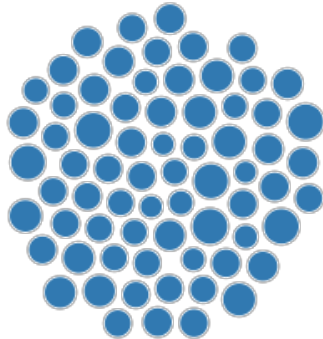


Treemap

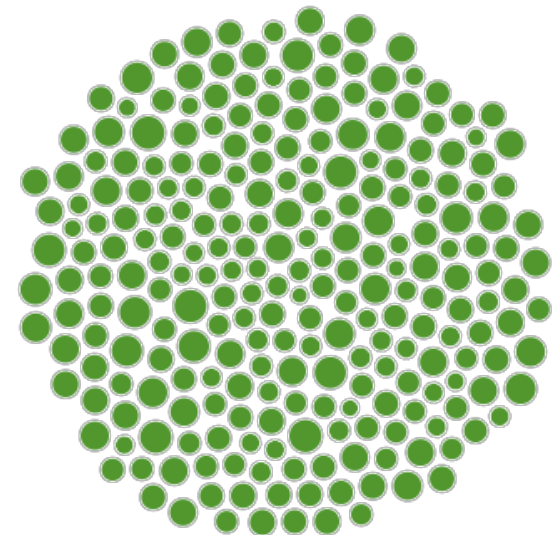


Bubble Chart

Origin
● Europe
● Japan
● USA



Beeswarm Plot



3D and Higher

Two variables $[x,y]$

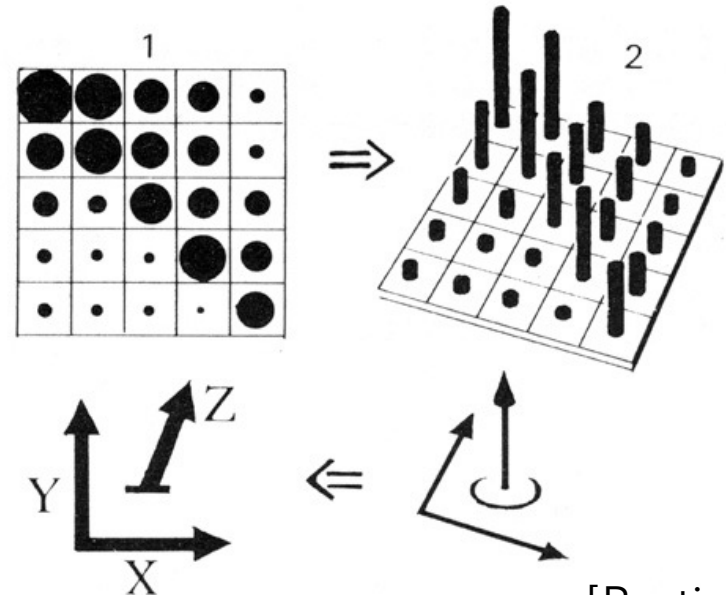
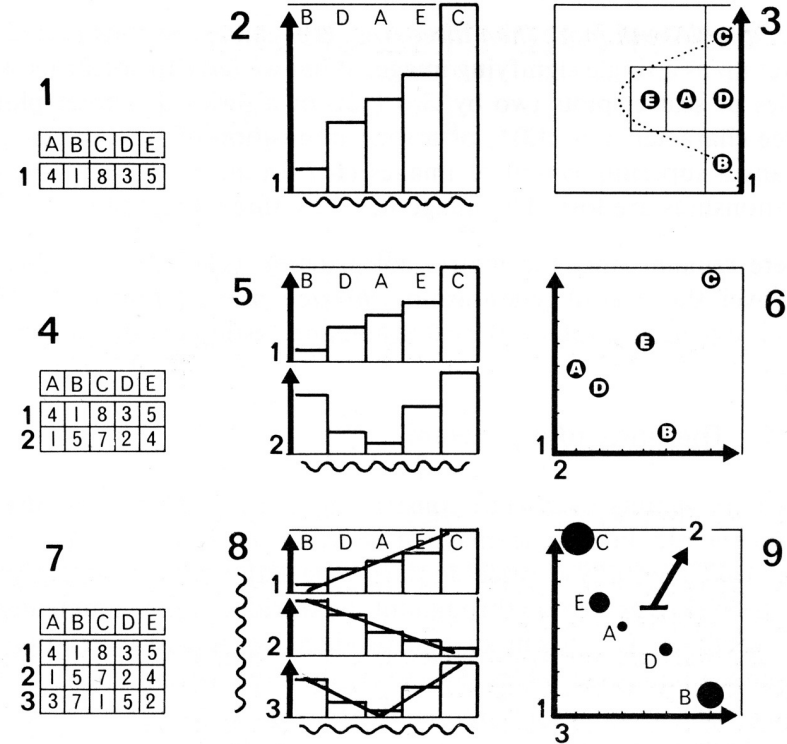
Can map to 2D points.

Scatterplots, maps, ...

Third variable $[z]$

Often use one of size, color, opacity, shape, etc. Or, one can further partition space.

What about 3D rendering?



Other Visual Encoding Channels?

wind map

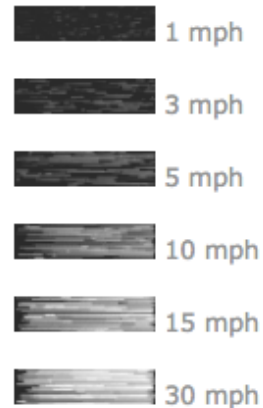
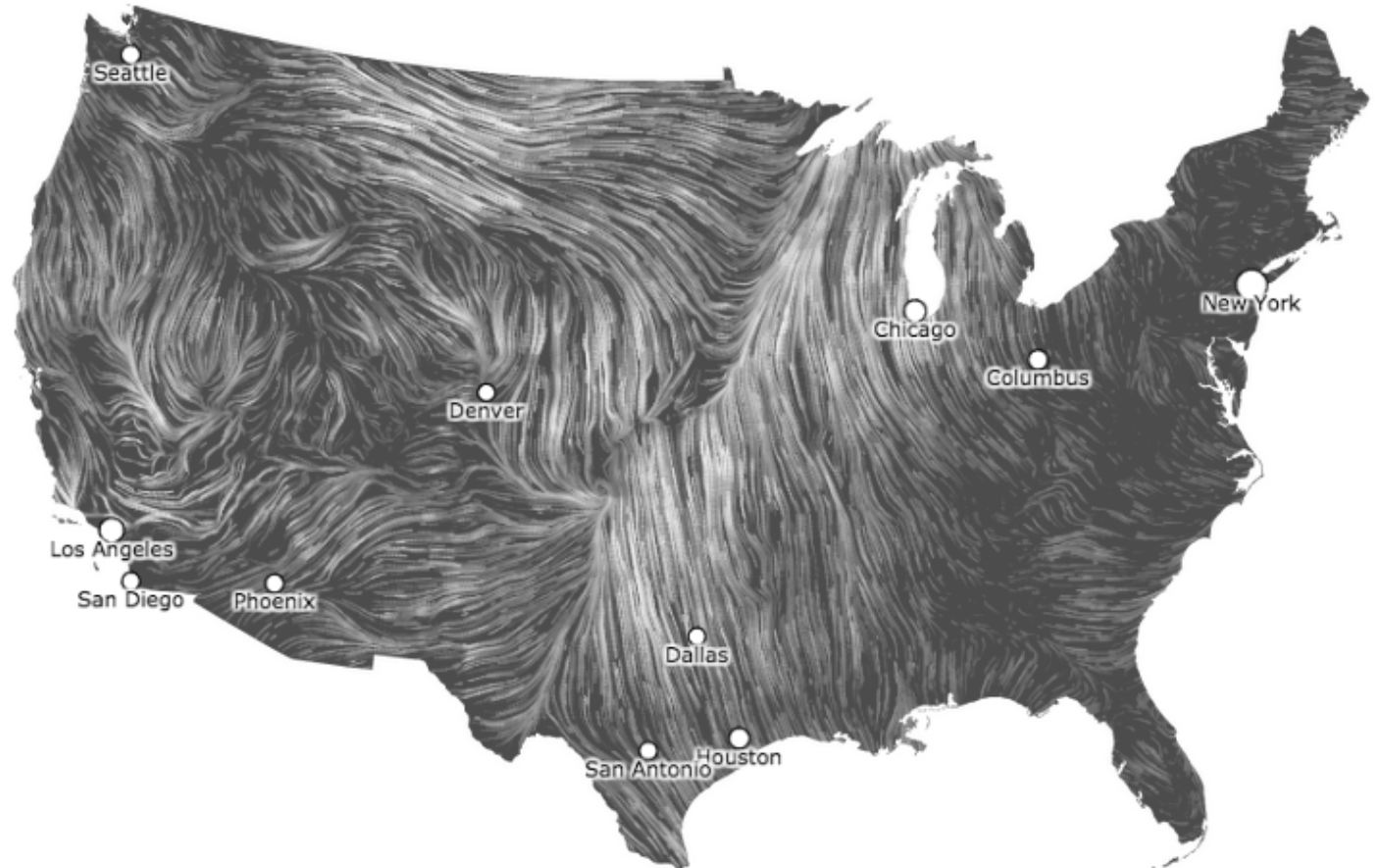
April 1, 2015

11:35 pm EST

(time of forecast download)

top speed: **30.5 mph**

average: **10.2 mph**



Encoding Effectiveness

Effectiveness Rankings [Mackinlay 86]

QUANTITATIVE

Position
Length
Angle
Slope
Area (Size)
Volume
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Shape

ORDINAL

Position
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Length
Angle
Slope
Area (Size)
Volume
Shape

NOMINAL

Position
Color Hue
Texture
Connection
Containment
Density (Value)
Color Sat
Shape
Length
Angle
Slope
Area
Volume

Effectiveness Rankings [Mackinlay 86]

QUANTITATIVE

Position

Length
Angle
Slope
Area (Size)
Volume
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Shape

ORDINAL

Position

Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Length
Angle
Slope
Area (Size)
Volume
Shape

NOMINAL

Position

Color Hue
Texture
Connection
Containment
Density (Value)
Color Sat
Shape
Length
Angle
Slope
Area
Volume

Effectiveness Rankings [Mackinlay 86]

QUANTITATIVE

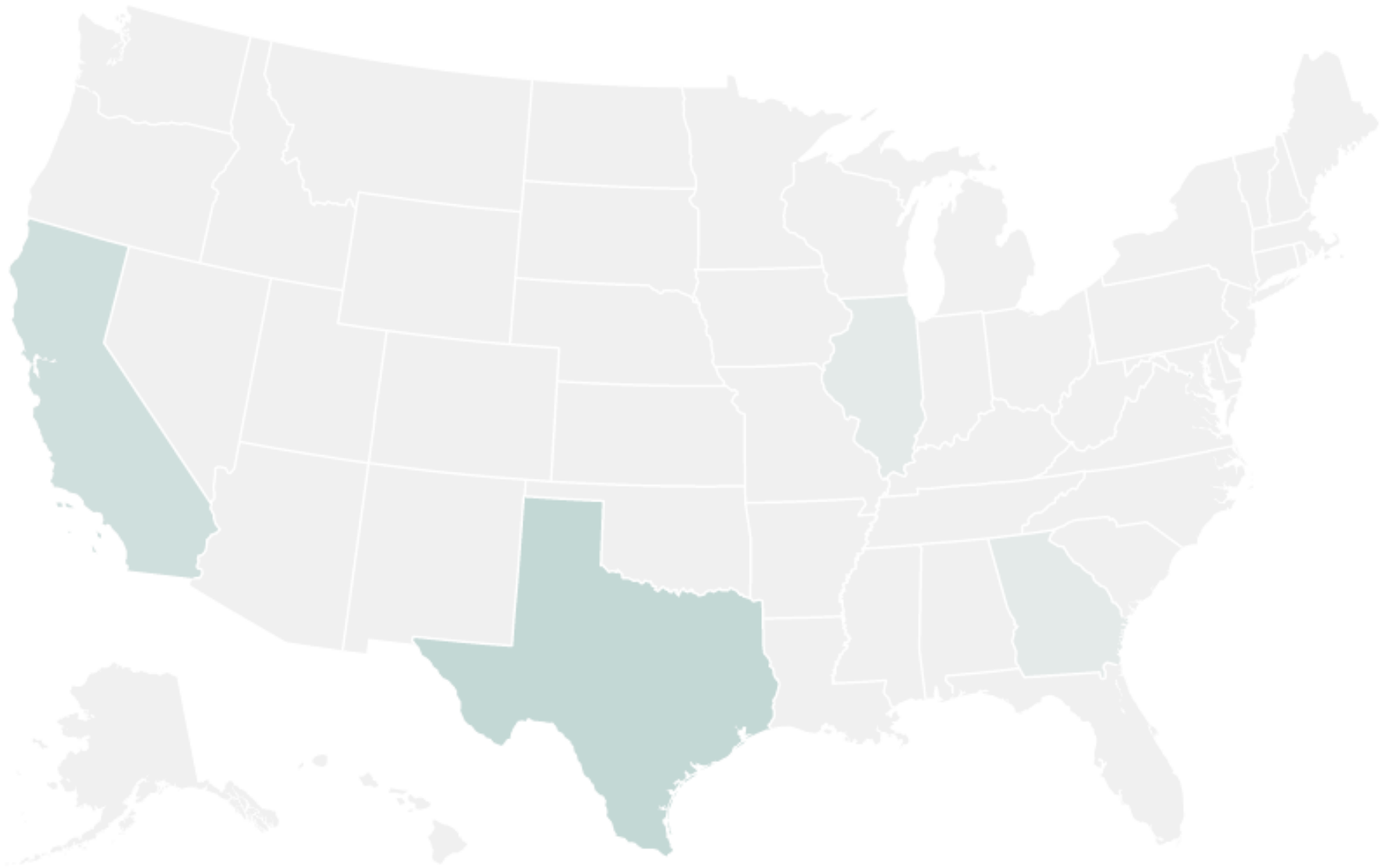
Position
Length
Angle
Slope
Area (Size)
Volume
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Shape

ORDINAL

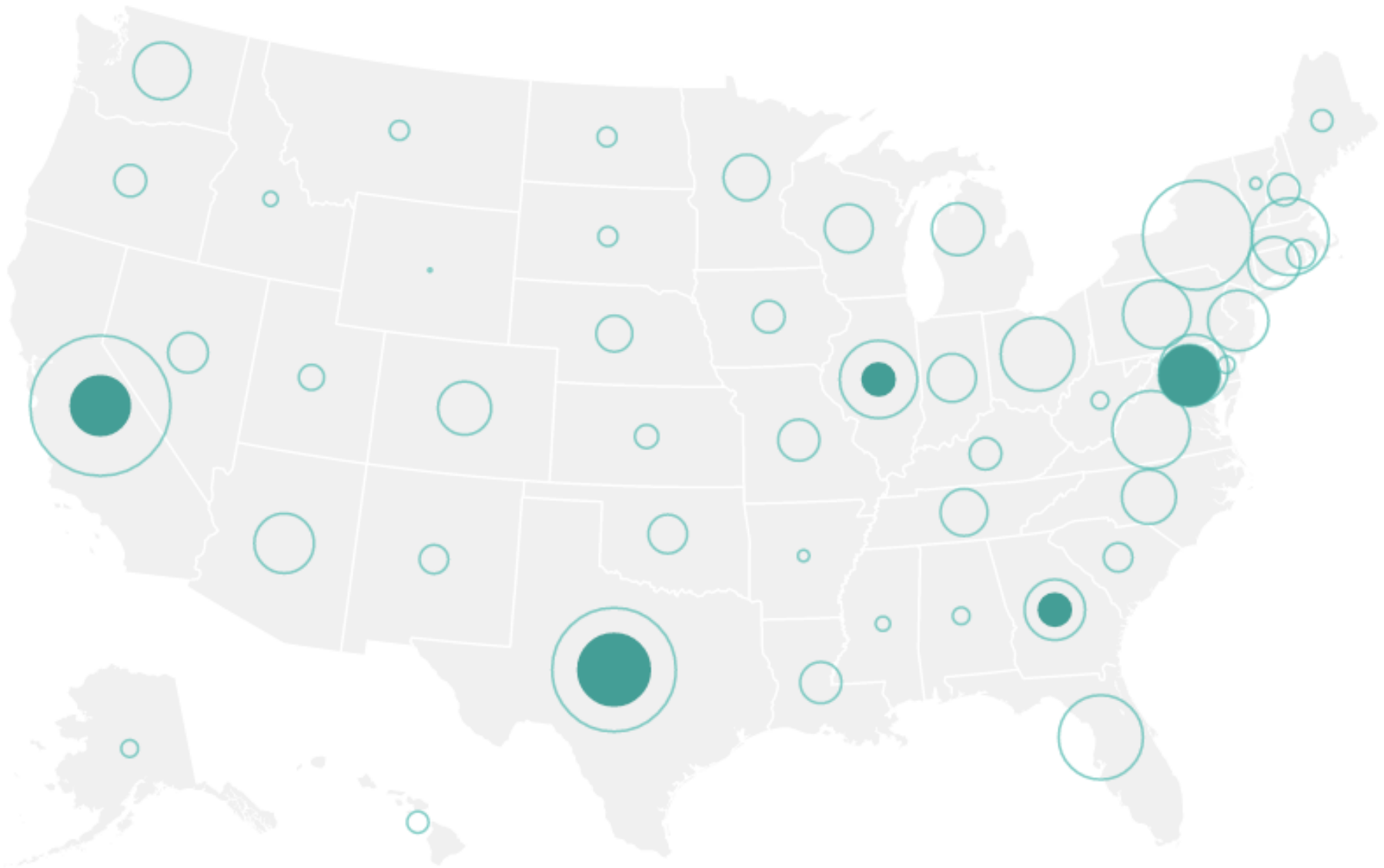
Position
Density (Value)
Color Sat
Color Hue
Texture
Connection
Containment
Length
Angle
Slope
Area (Size)
Volume
Shape

NOMINAL

Position
Color Hue
Texture
Connection
Containment
Density (Value)
Color Sat
Shape
Length
Angle
Slope
Area
Volume



Color Encoding



Area Encoding

Effectiveness Rankings

QUANTITATIVE

Position

Length

Angle

Slope

Area (Size)

Volume

Density (Value)

Color Sat

Color Hue

Texture

Connection

Containment

Shape

ORDINAL

Position

Density (Value)

Color Sat

Color Hue

Texture

Connection

Containment

Length

Angle

Slope

Area (Size)

Volume

Shape

NOMINAL

Position

Color Hue

Texture

Connection

Containment

Density (Value)

Color Sat

Shape

Length

Angle

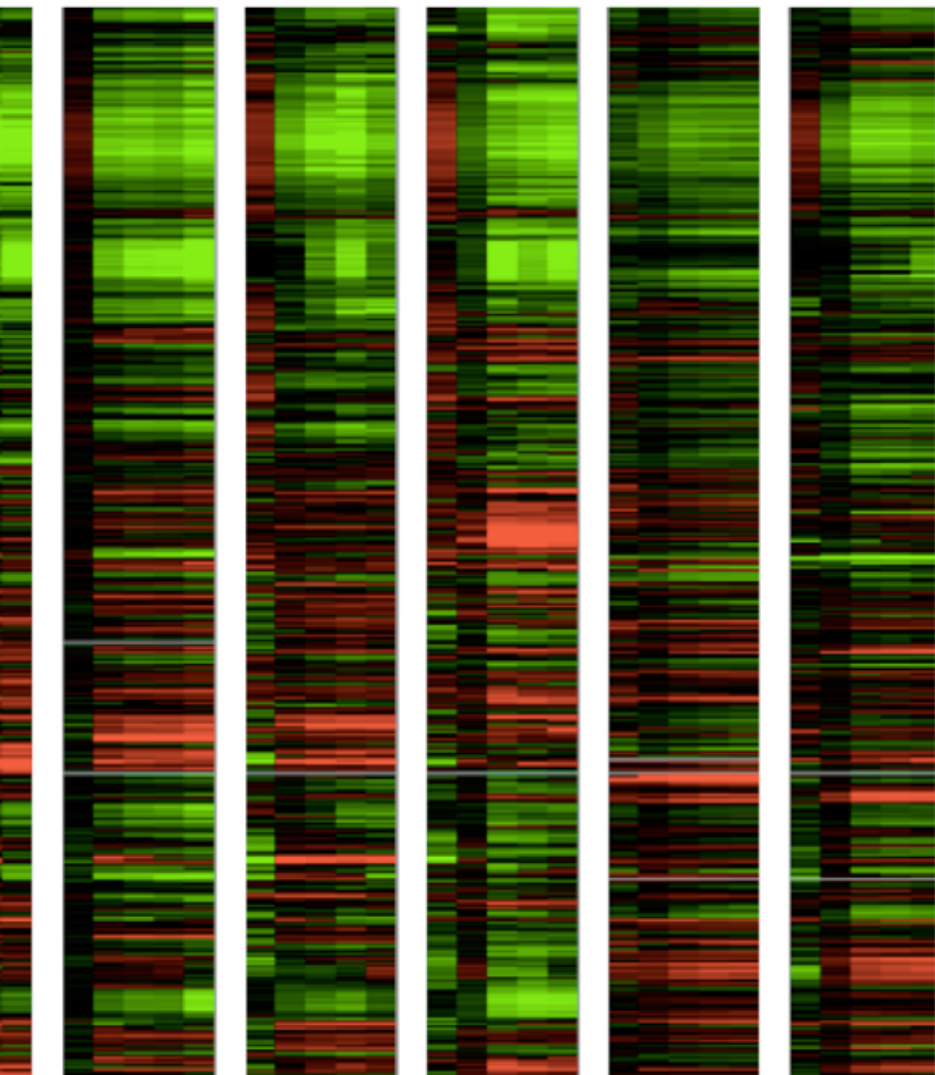
Slope

Area

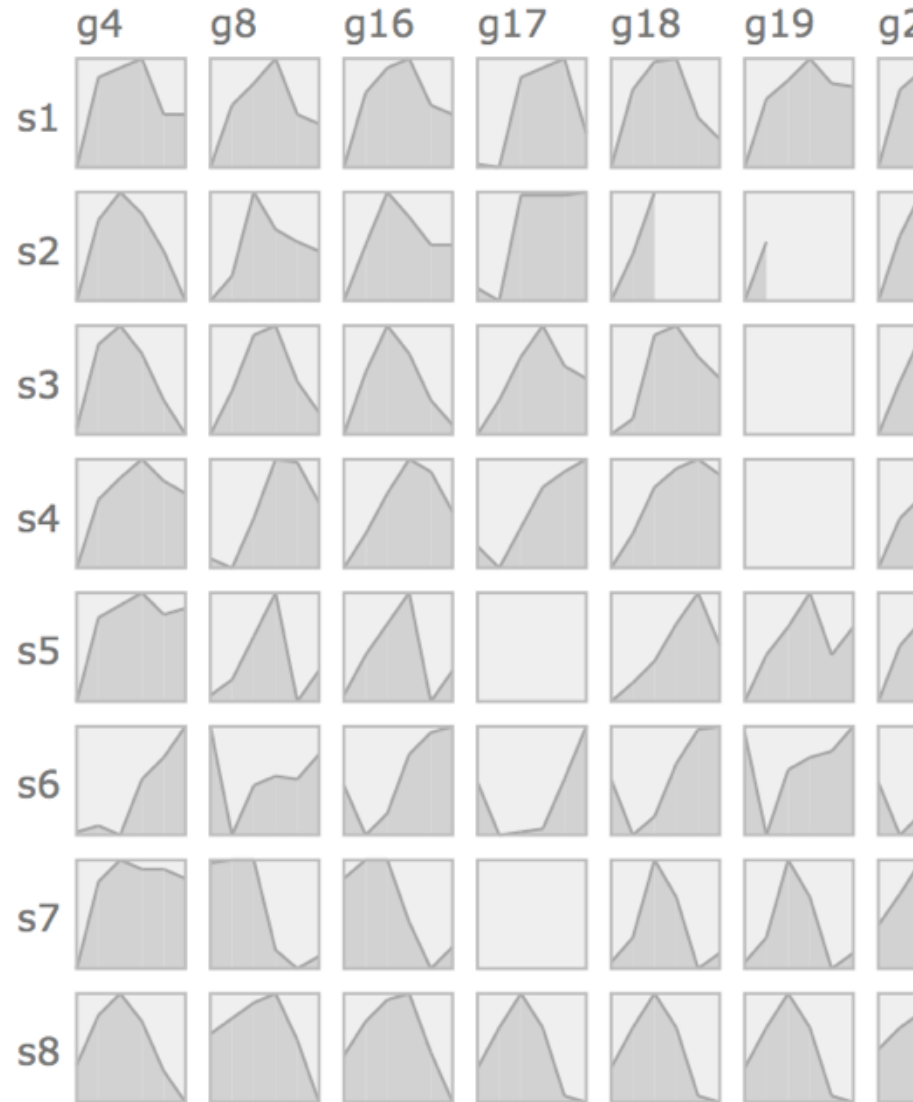
Volume

Gene Expression Time-Series [Meyer et al '11]

Color Encoding



Position Encoding



Effectiveness Rankings

QUANTITATIVE

Position

Length

Angle

Slope

Area (Size)

Volume

~~Density (Value)~~

Color Sat

~~Color Hue~~

Texture

Connection

Containment

Shape

ORDINAL

Position

Density (Value)

Color Sat

Color Hue

Texture

Connection

Containment

Length

Angle

Slope

Area (Size)

Volume

Shape

NOMINAL

Position

Color Hue

Texture

Connection

Containment

Density (Value)

Color Sat

Shape

Length

Angle

Slope

Area

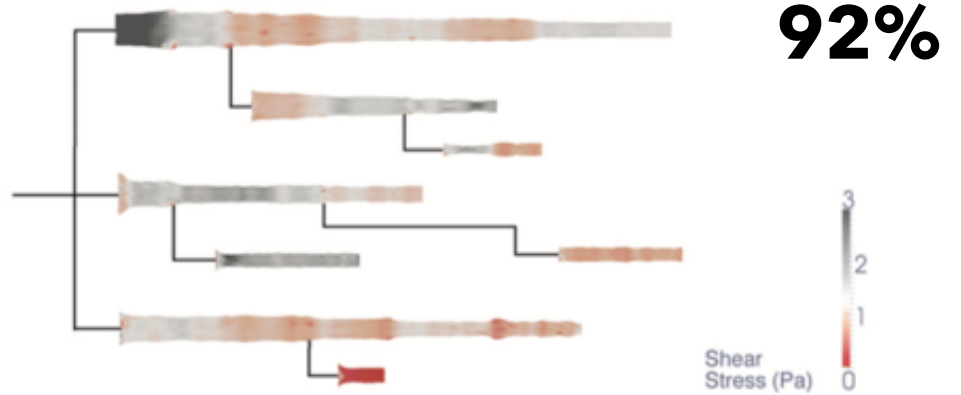
Volume

Artery Visualization [Borkin et al '11]

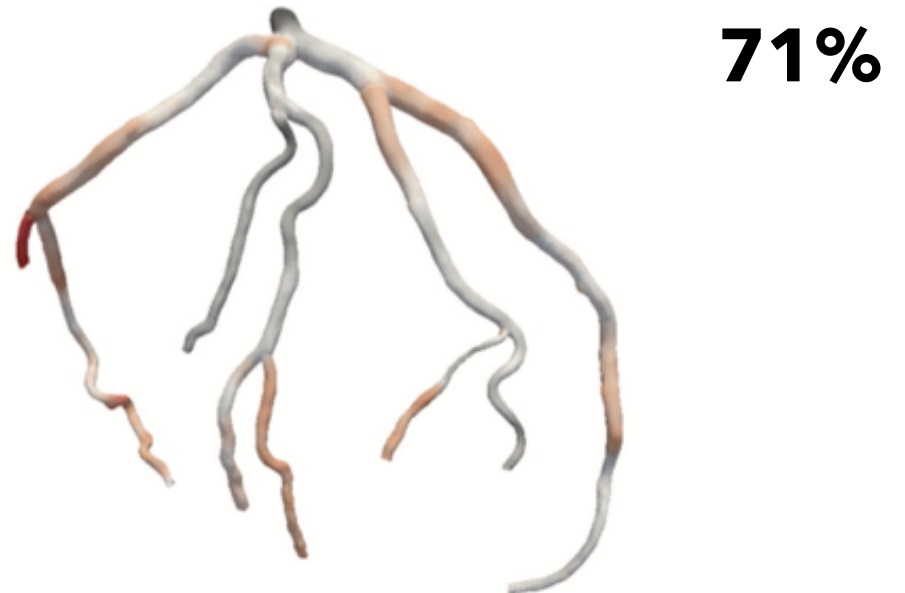
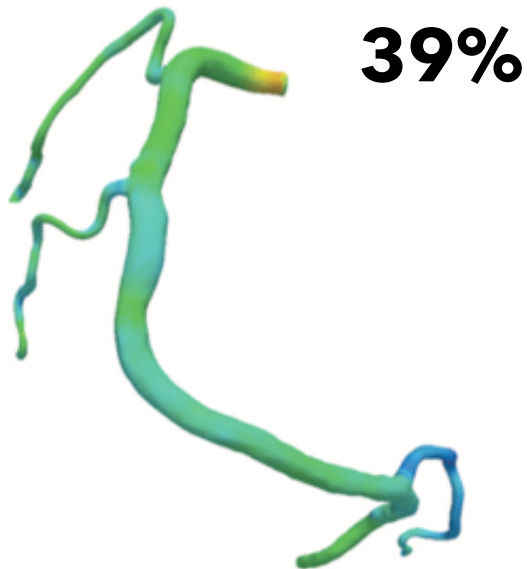
Rainbow Palette

Diverging Palette

2D



3D



Effectiveness Rankings

QUANTITATIVE

Position ↻

Length

Angle

Slope

Area (Size)

Volume

Density (Value)

Color Sat

Color Hue

Texture

Connection

Containment

Shape

ORDINAL

Position

Density (Value)

Color Sat

Color Hue

Texture

Connection

Containment

Length

Angle

Slope

Area (Size)

Volume

Shape

NOMINAL

Position

Color Hue

Texture

Connection

Containment

Density (Value)

Color Sat

Shape

Length

Angle

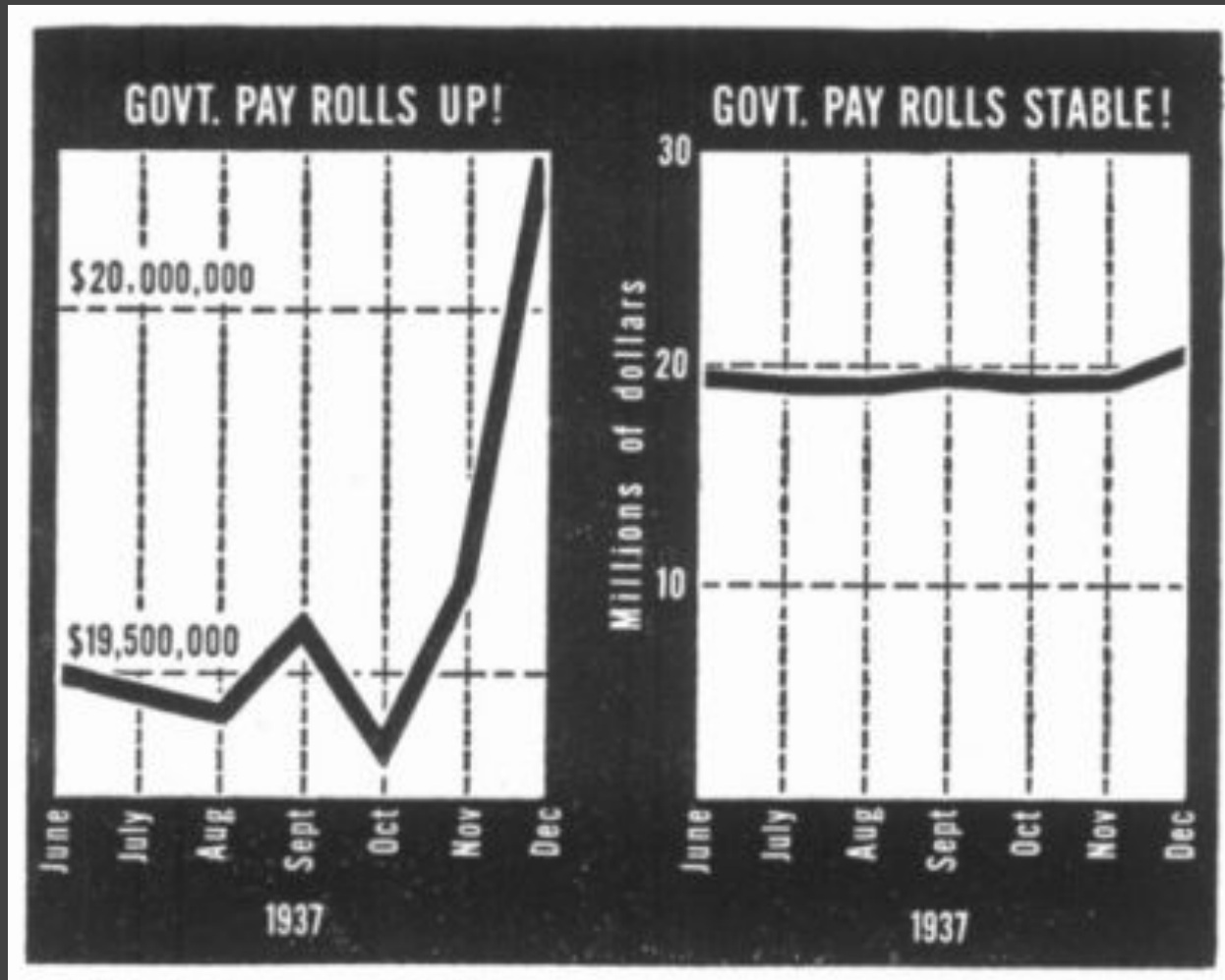
Slope

Area

Volume

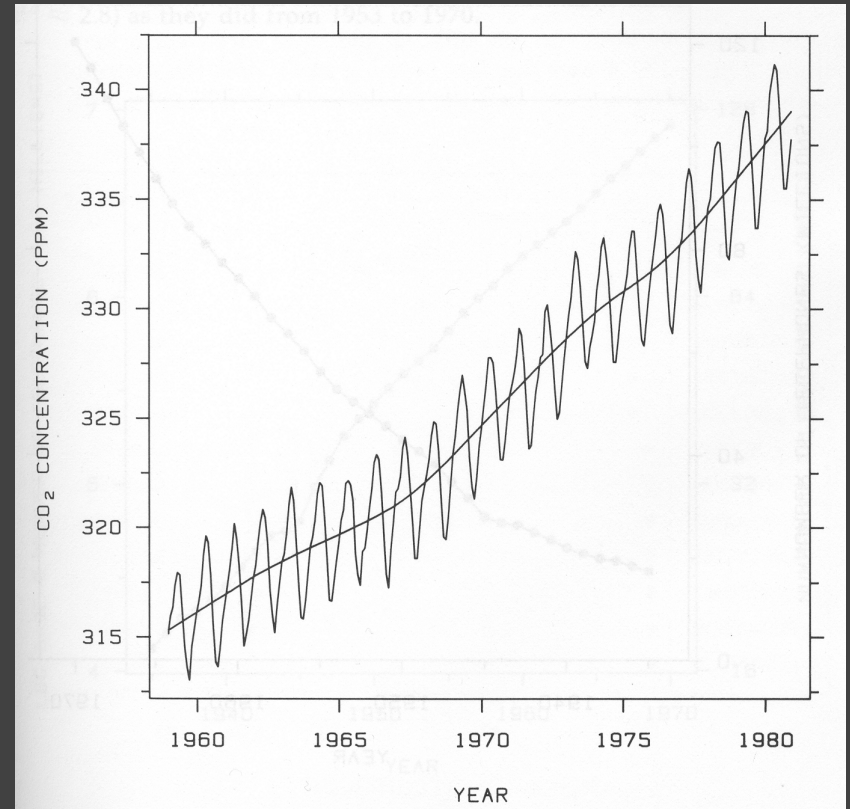
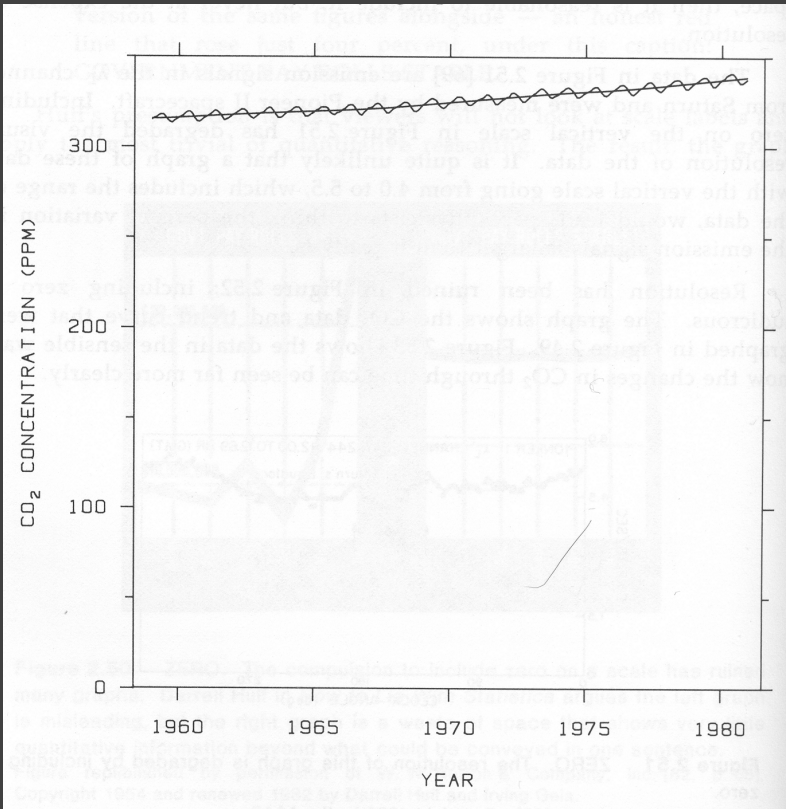
Scales & Axes

Include Zero in Axis Scale?



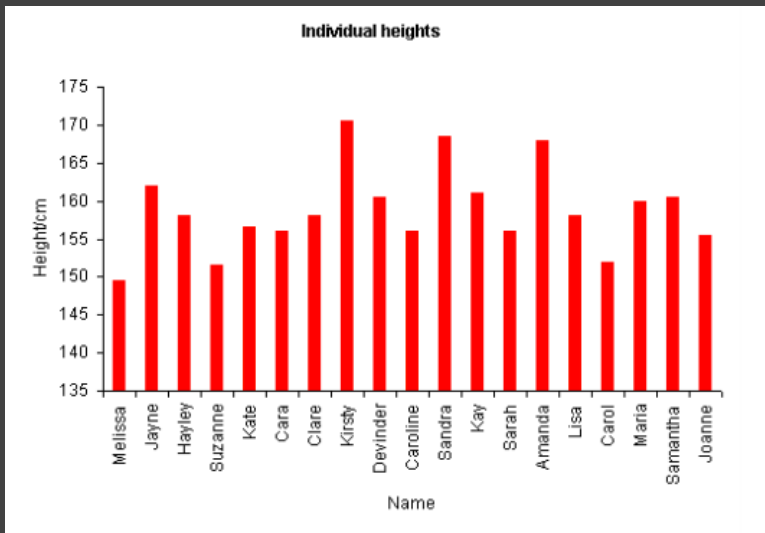
Government payrolls in 1937 [How To Lie With Statistics. Huff]

Include Zero in Axis Scale?



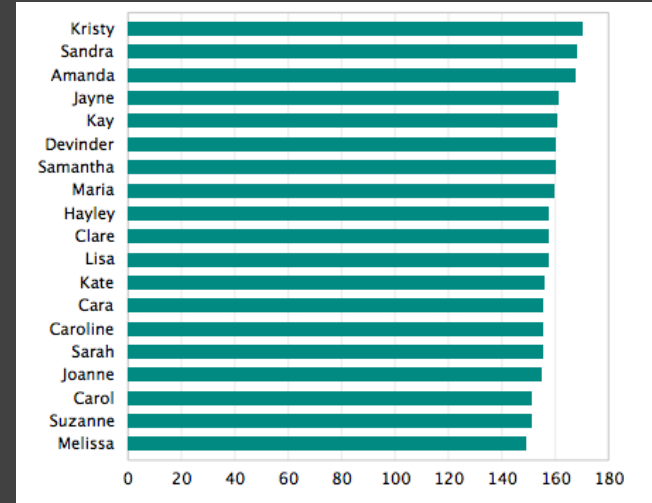
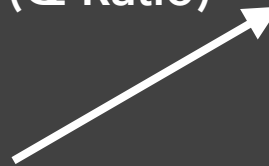
Yearly CO₂ concentrations [Cleveland 85]

Include Zero in Axis Scale?

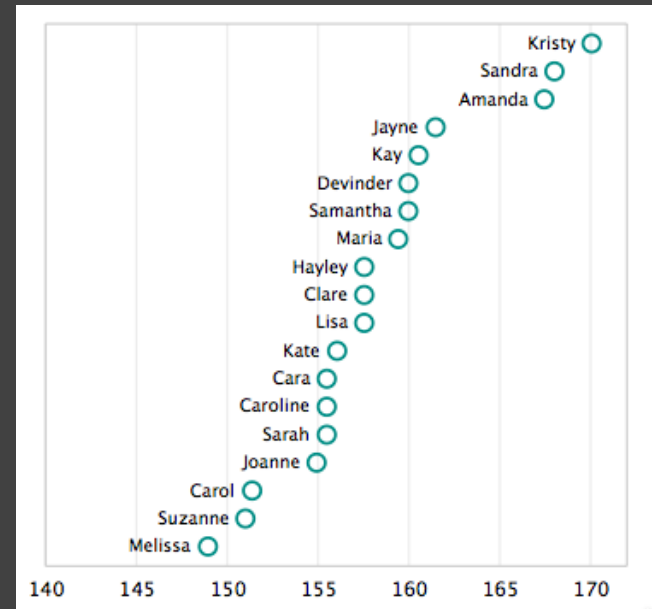


Violates Expressiveness Principle!

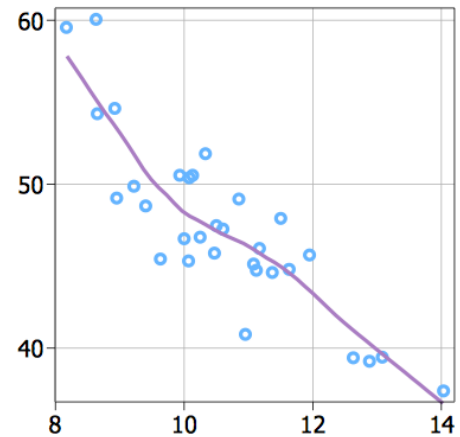
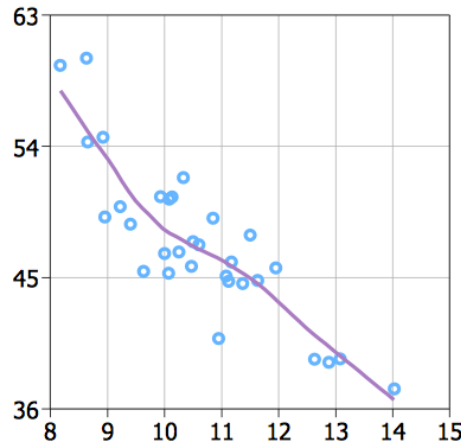
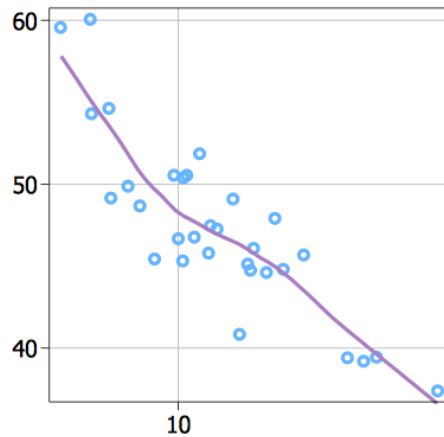
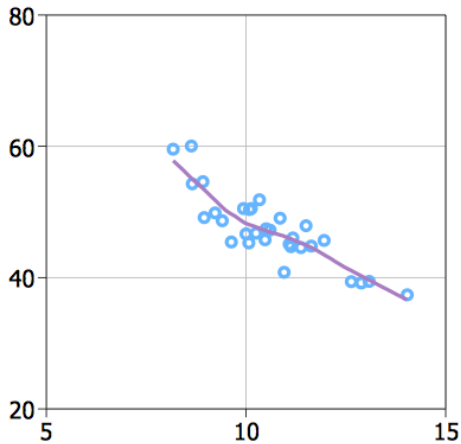
Compare Proportions (Q-Ratio)



Compare Relative Position (Q-Interval)

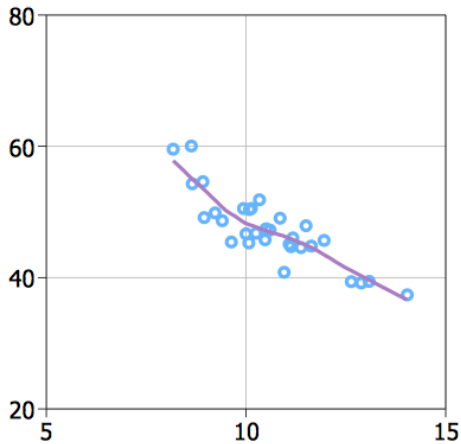


Axis Tick Mark Selection

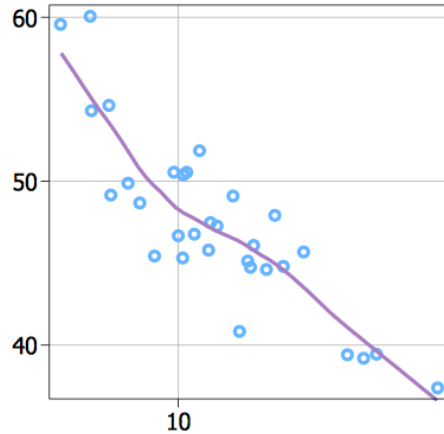


What are some properties of "good" tick marks?

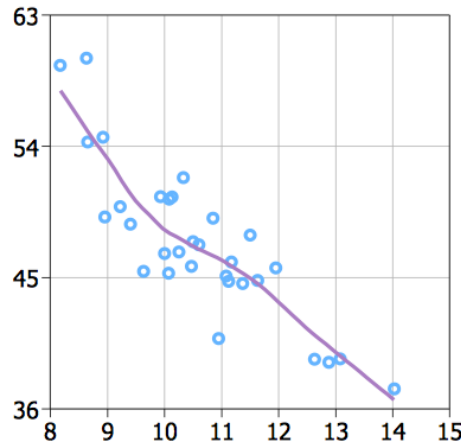
Axis Tick Mark Selection



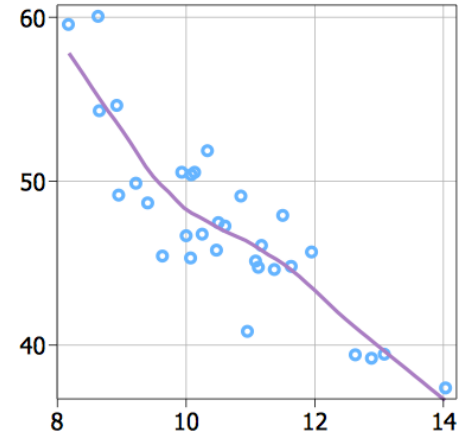
(a) Heckbert



(b) R's pretty



(c) Wilkinson



(d) Extended

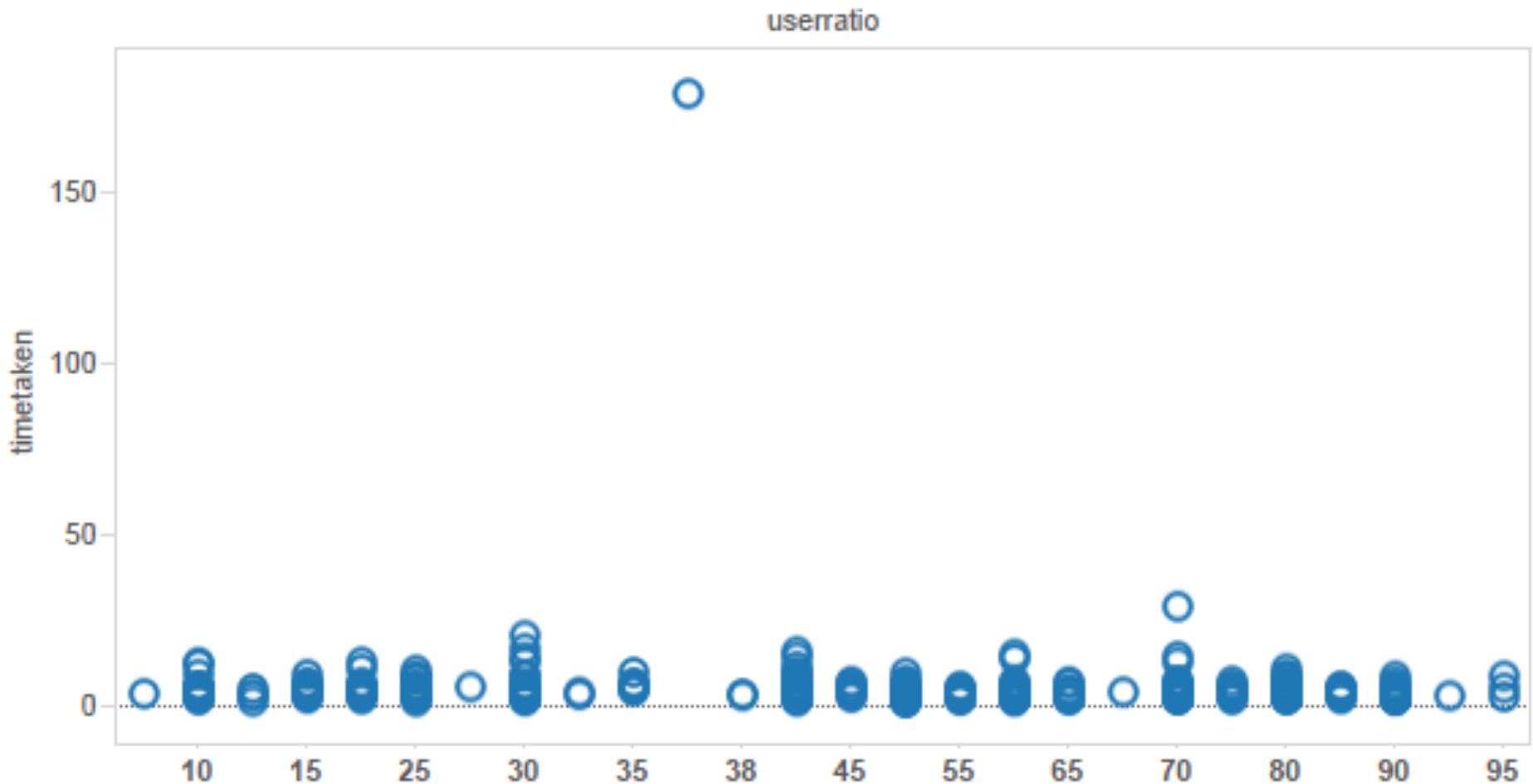
Simplicity - numbers are multiples of 10, 5, 2

Coverage - ticks near the ends of the data

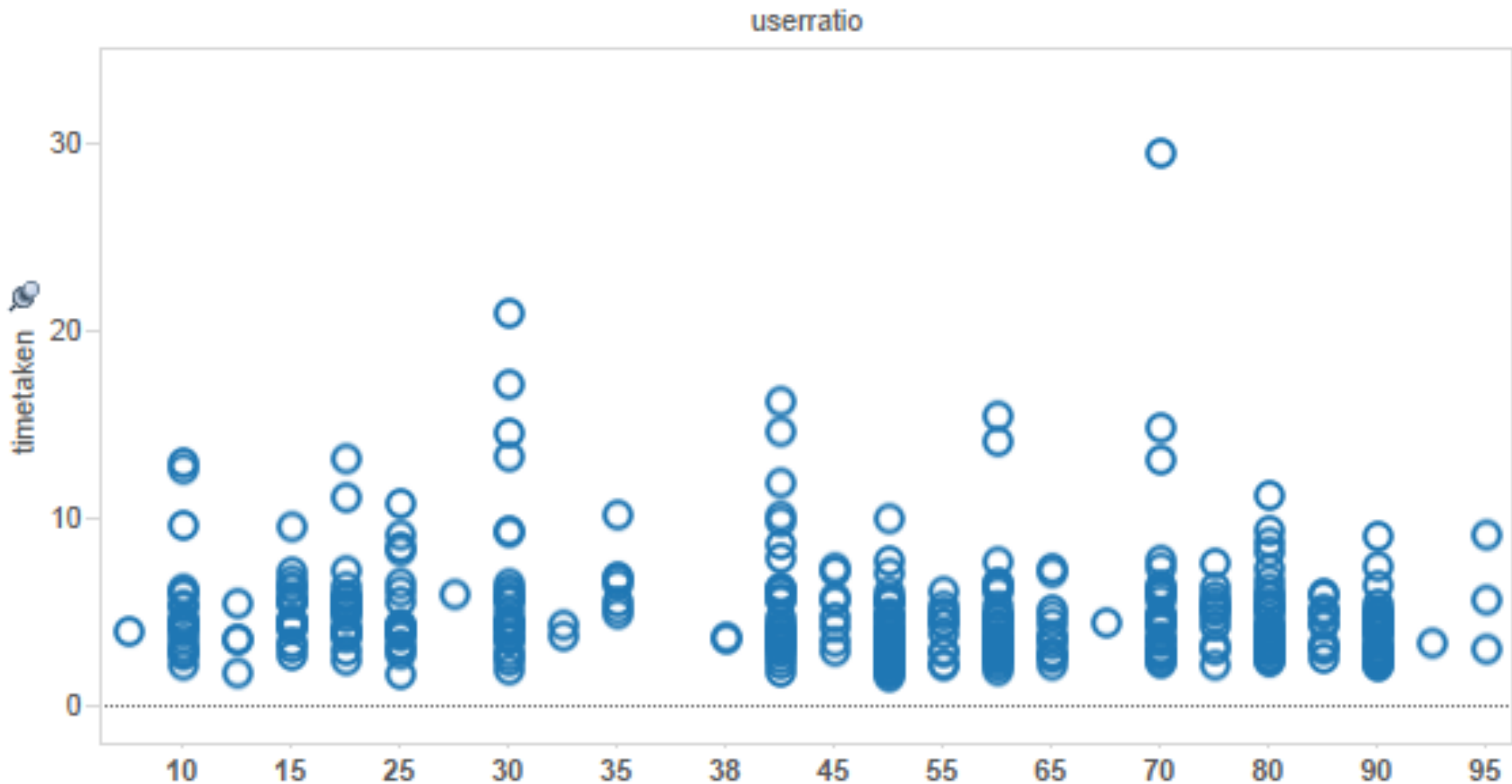
Density - not too many, nor too few

Legibility - whitespace, horizontal text, size

How to Scale the Axis?

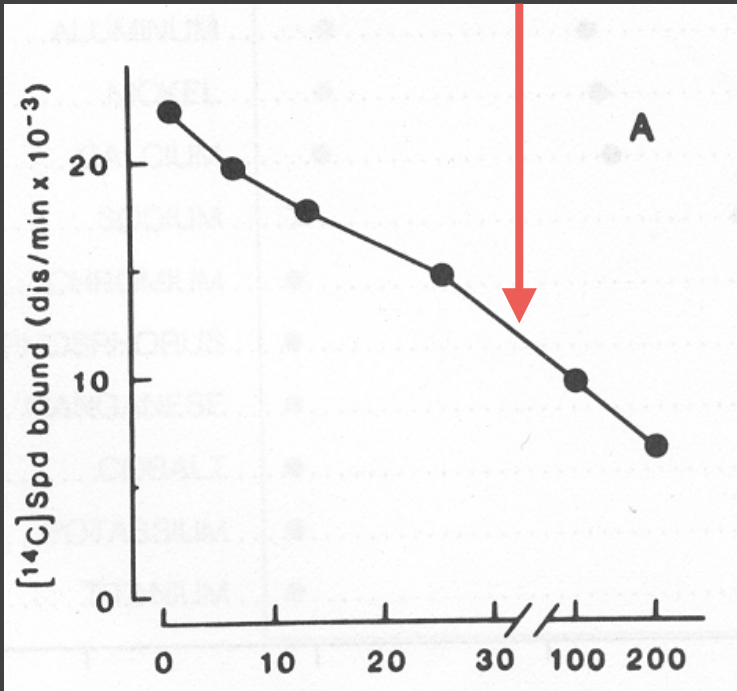


One Option: Clip Outliers

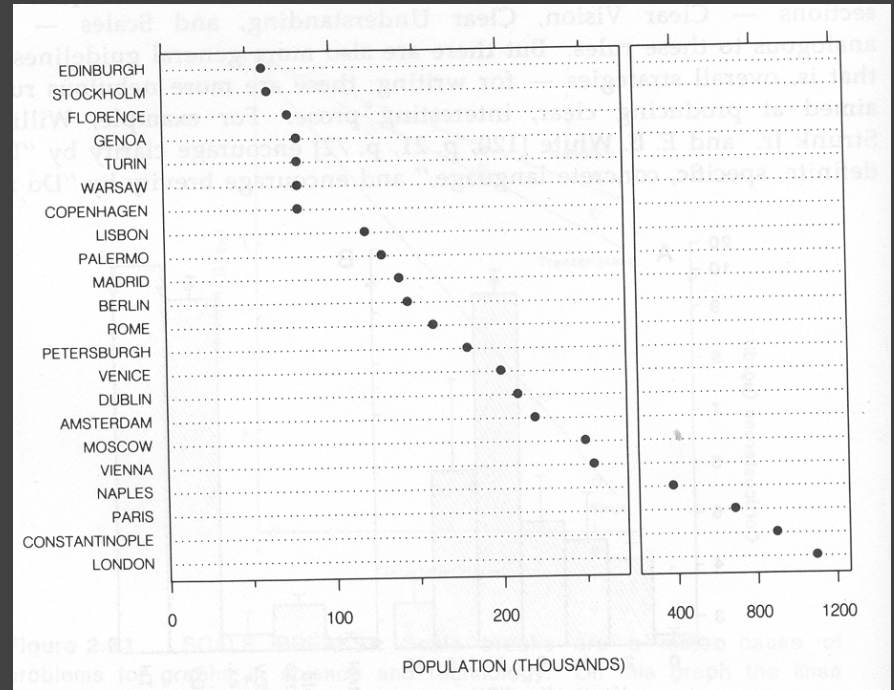


Clearly Mark Scale Breaks

Violates Expressiveness Principle!

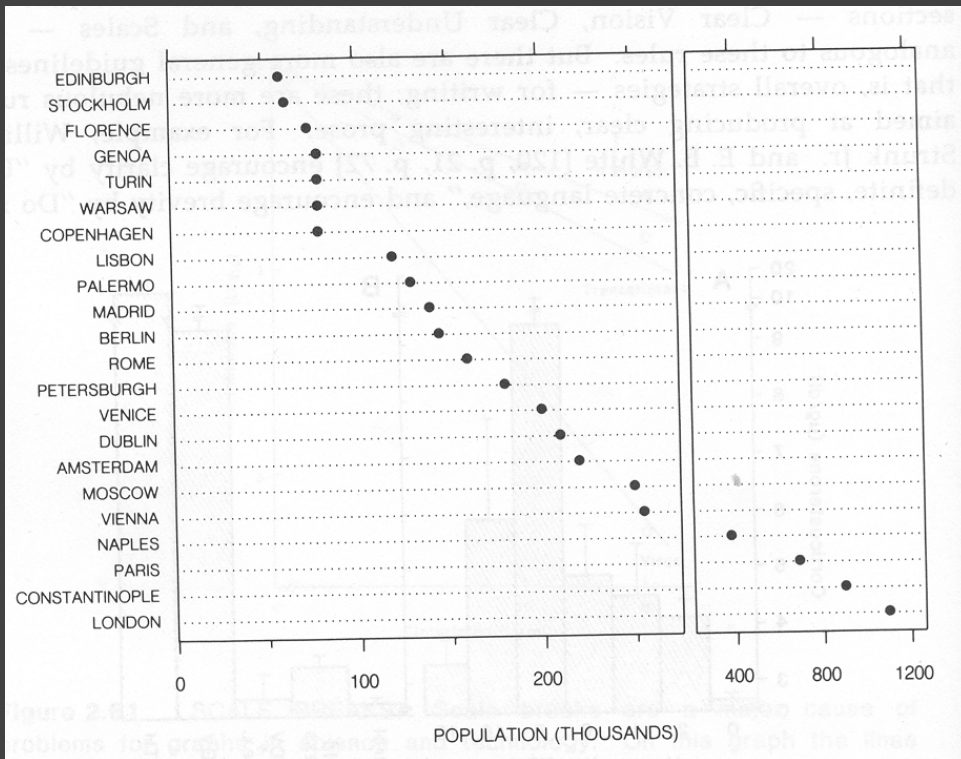


Poor scale break [Cleveland 85]

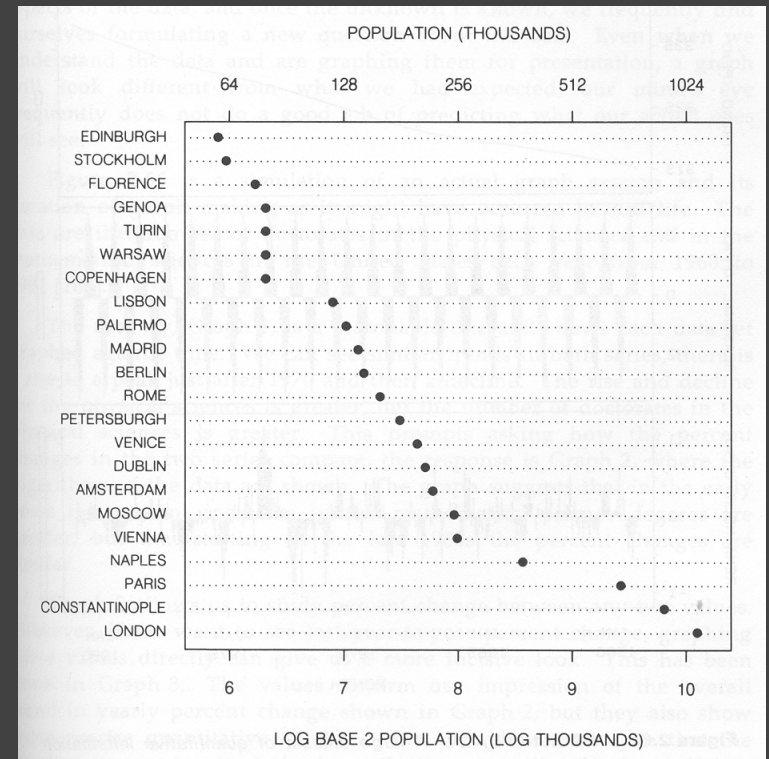


Well-marked scale break [Cleveland 85]

Scale Break vs. Log Scale

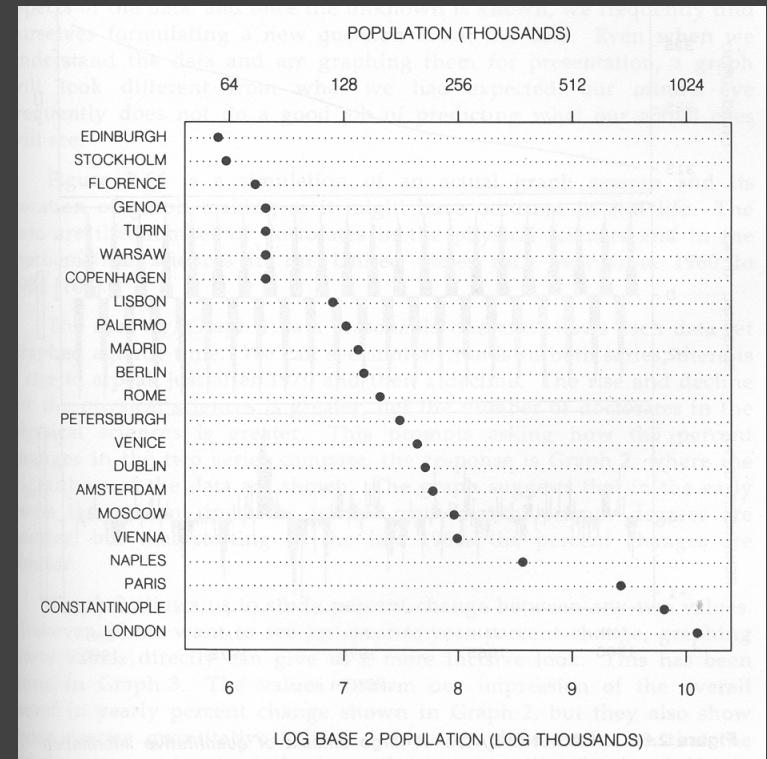
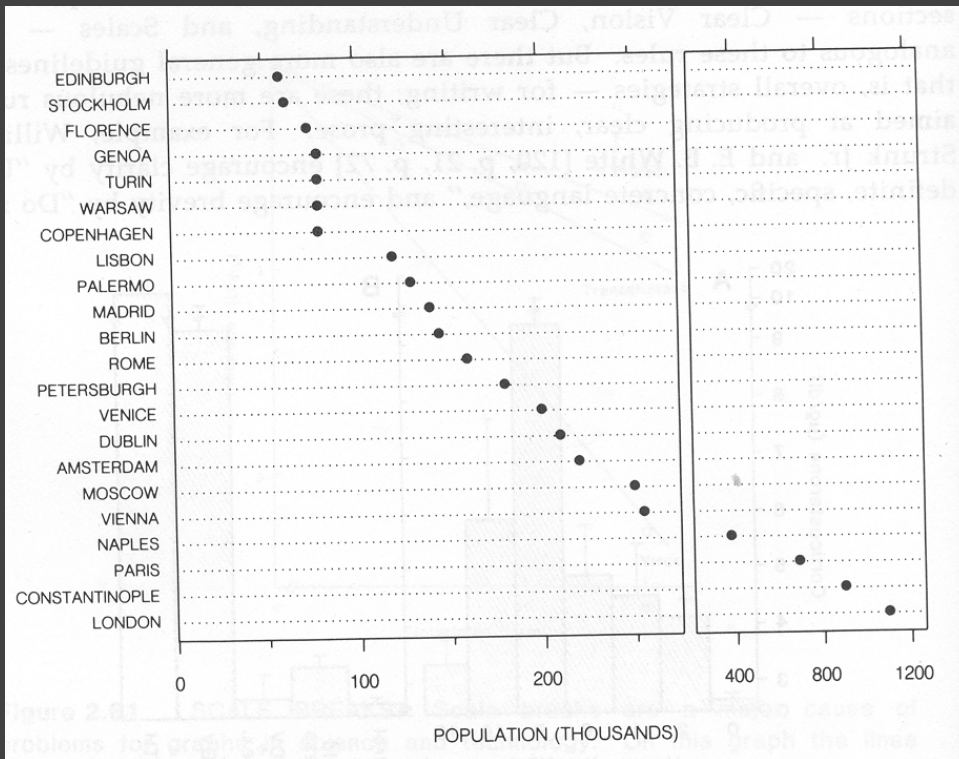


Scale Break



Log Scale

Scale Break vs. Log Scale



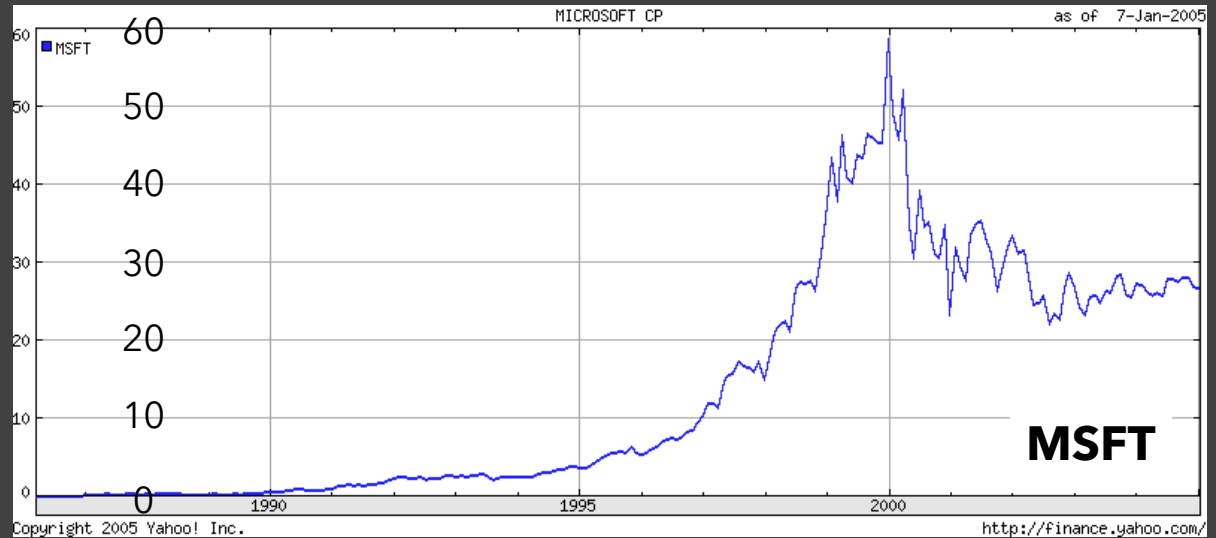
Both increase visual resolution

Scale break: difficult to compare (*cognitive* – not *perceptual* – work)

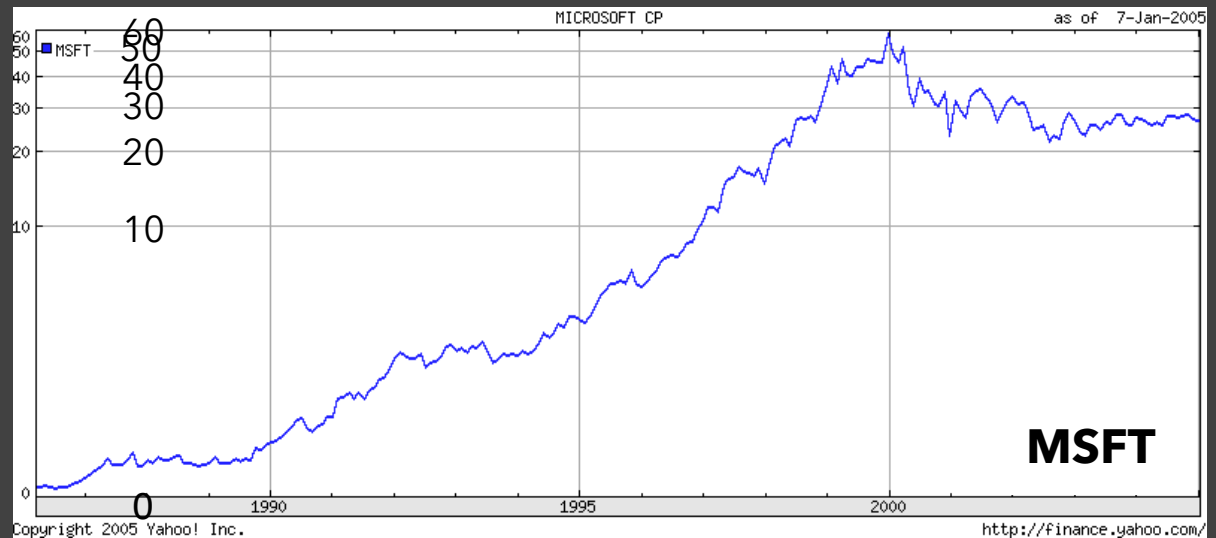
Log scale: direct comparison of all data

Linear Scale vs. Log Scale

Linear Scale



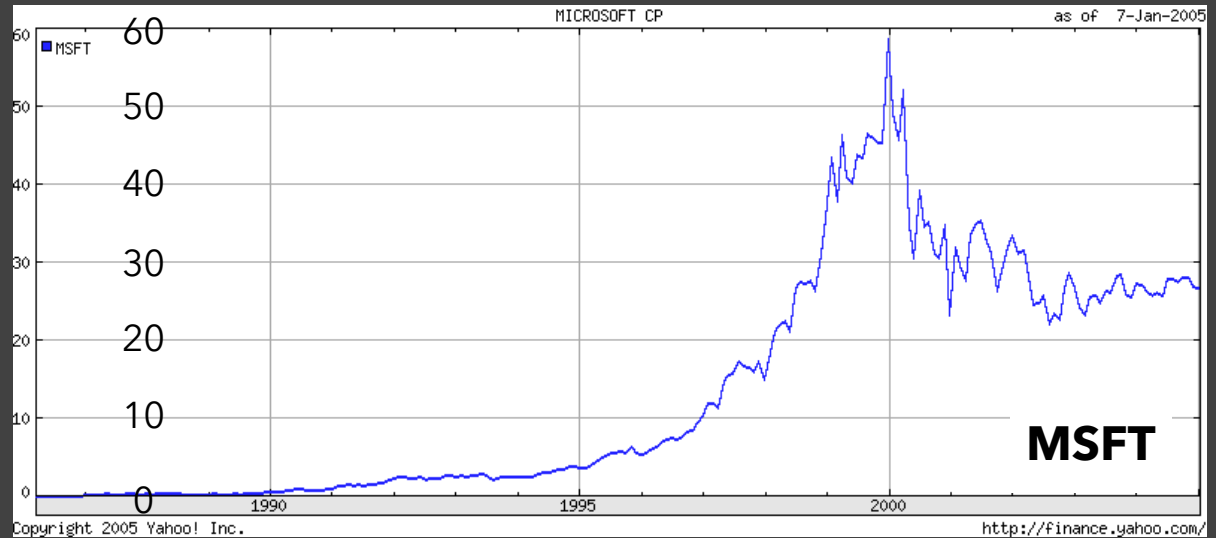
Log Scale



Linear Scale vs. Log Scale

Linear Scale

Absolute change

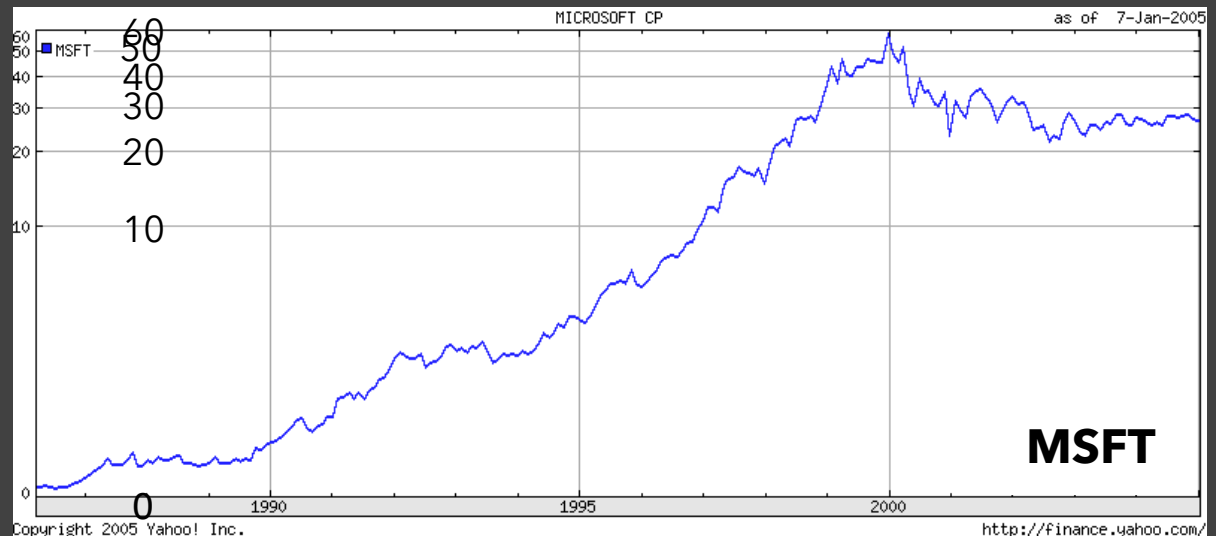


Log Scale

Small fluctuations

Percent change

$$d(10,20) = d(30,60)$$



When To Apply a Log Scale?

Address data skew (e.g., long tails, outliers)

Enables comparison within and across multiple orders of magnitude.

Focus on multiplicative factors (not additive)

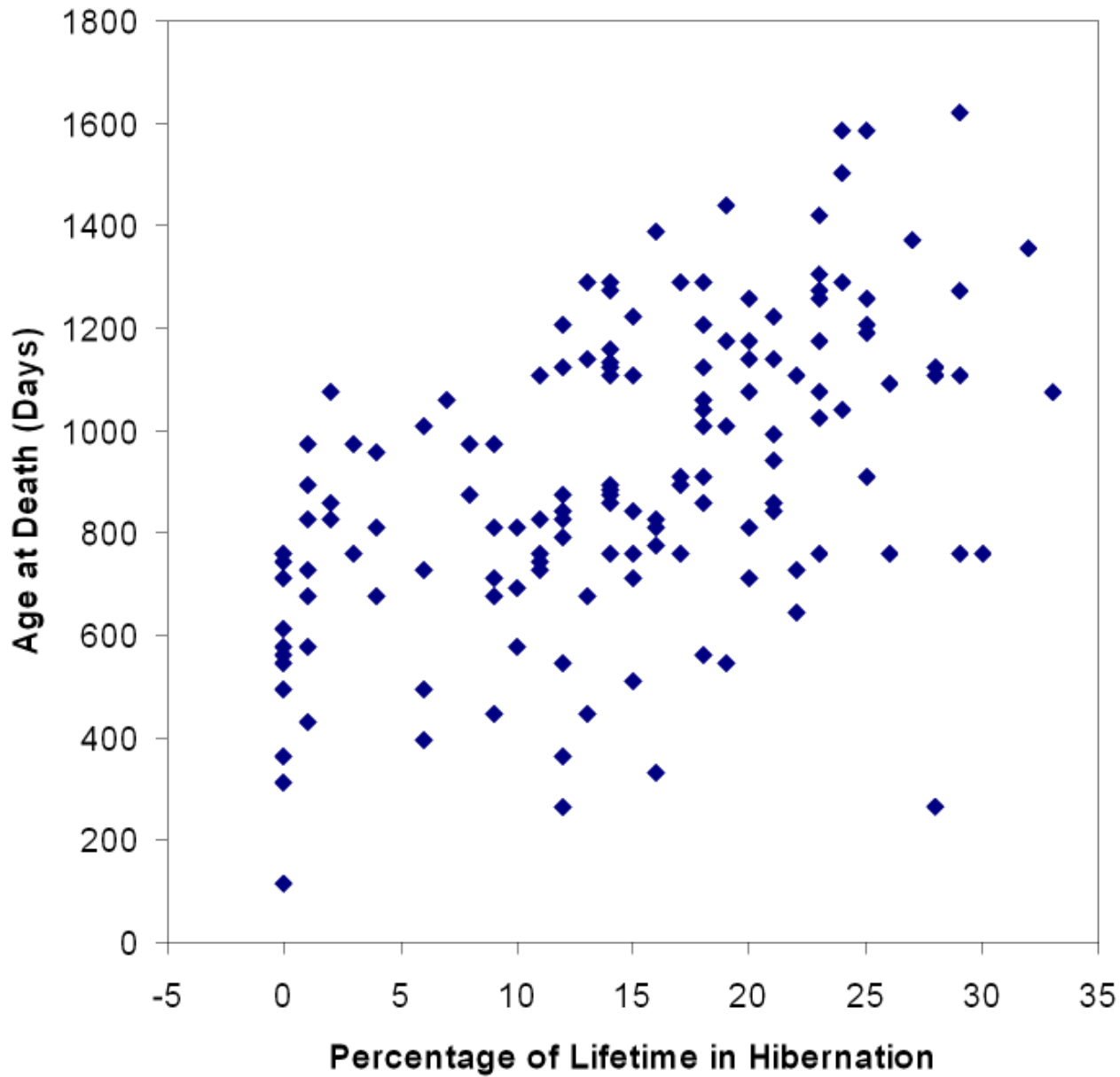
Recall that the logarithm transforms \times to $+$!

Percentage change, not absolute value.

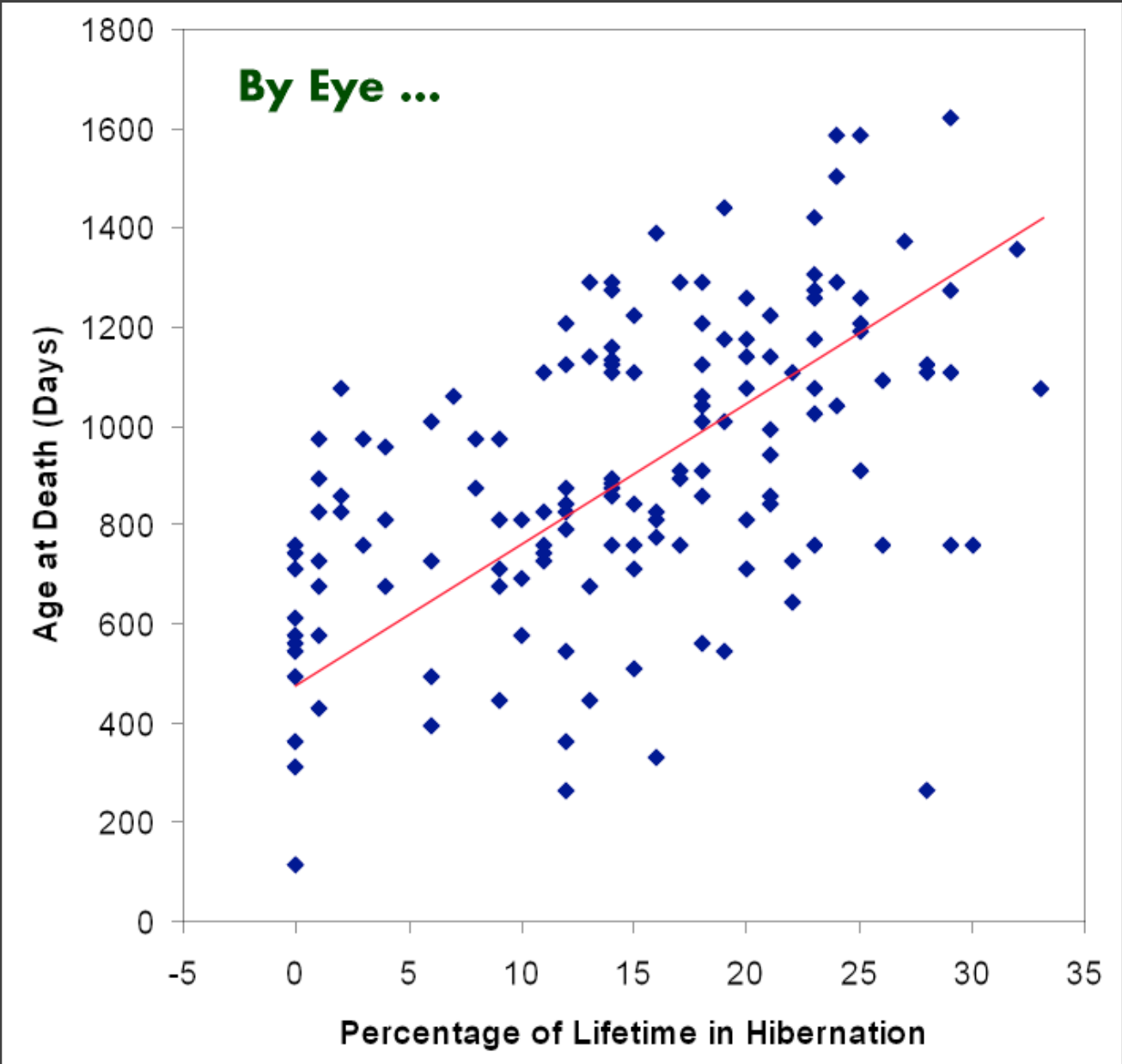
Constraint: **positive, non-zero values**

Constraint: **audience familiarity?**

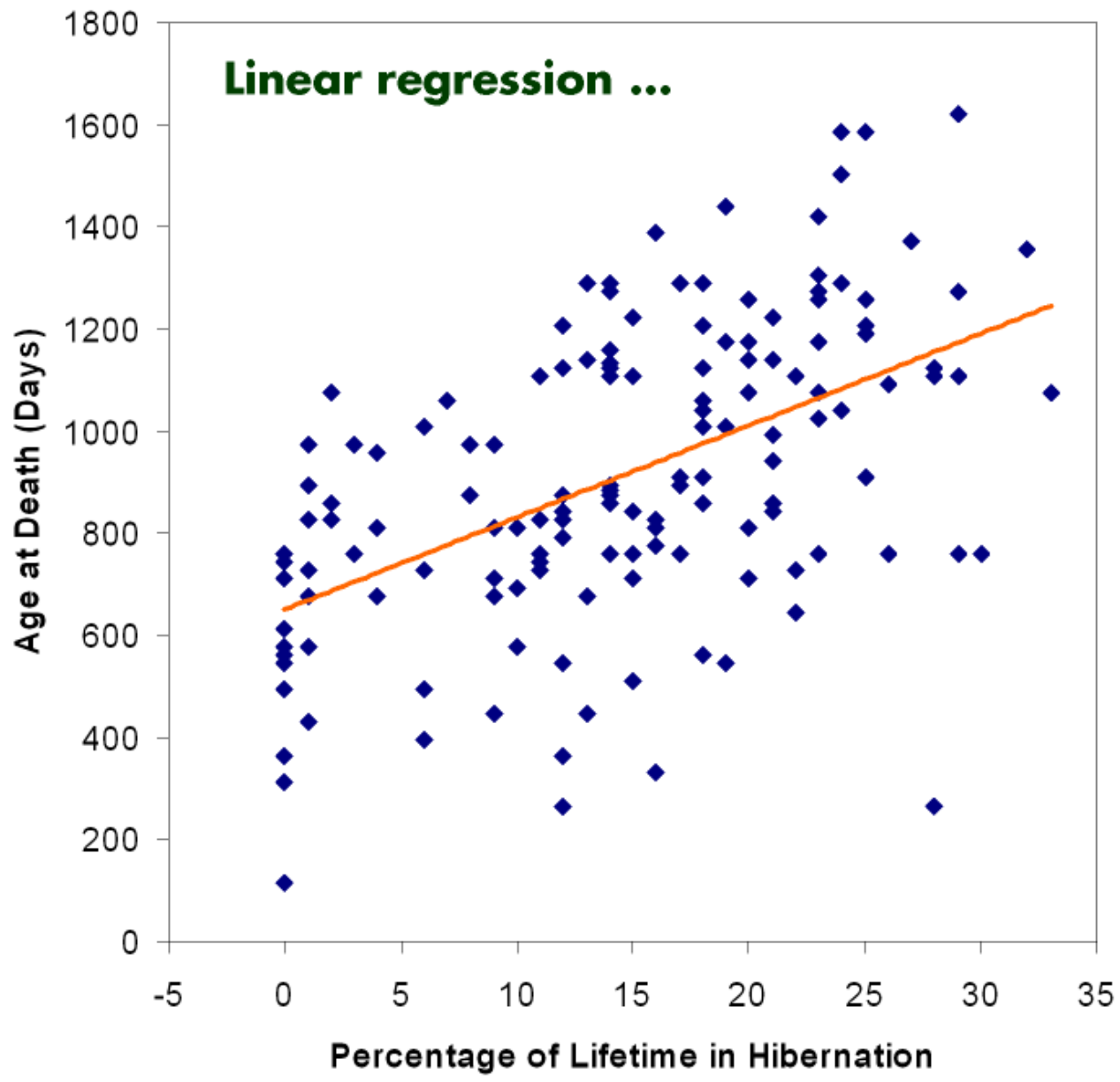
Regression Lines



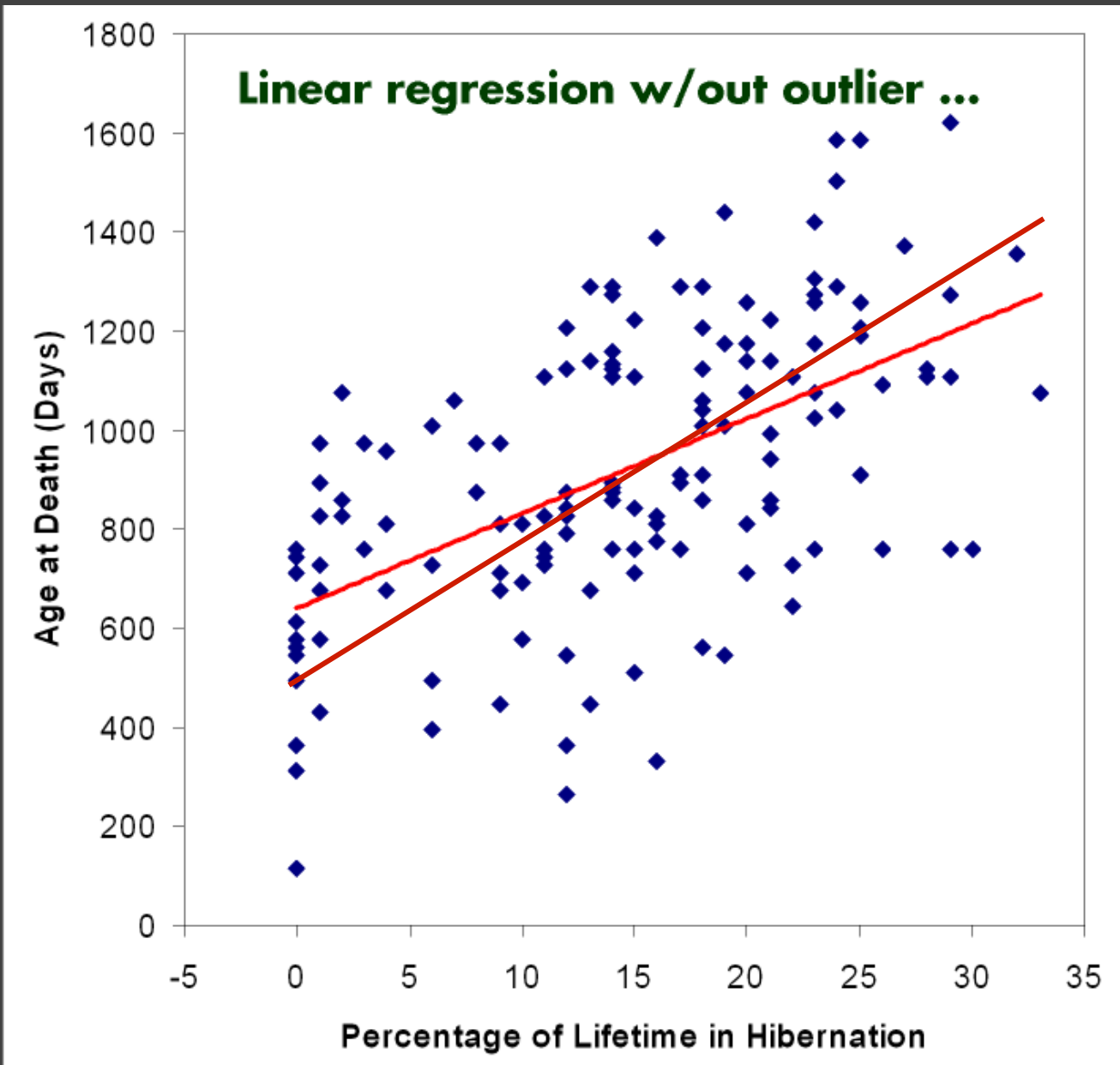
[The Elements of Graphing Data. Cleveland 94]



[The Elements of Graphing Data. Cleveland 94]



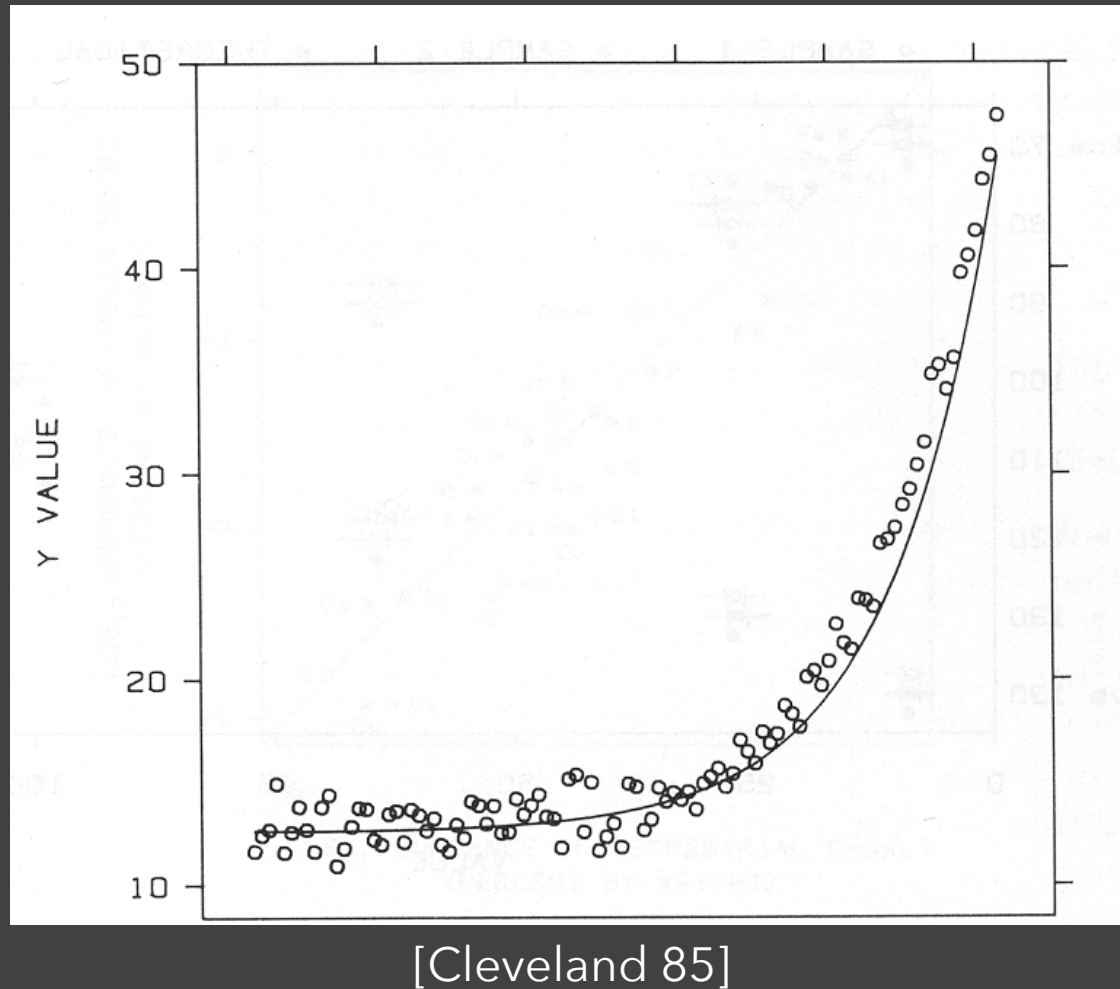
[The Elements of Graphing Data. Cleveland 94]



[The Elements of Graphing Data. Cleveland 94]

Transforming Data

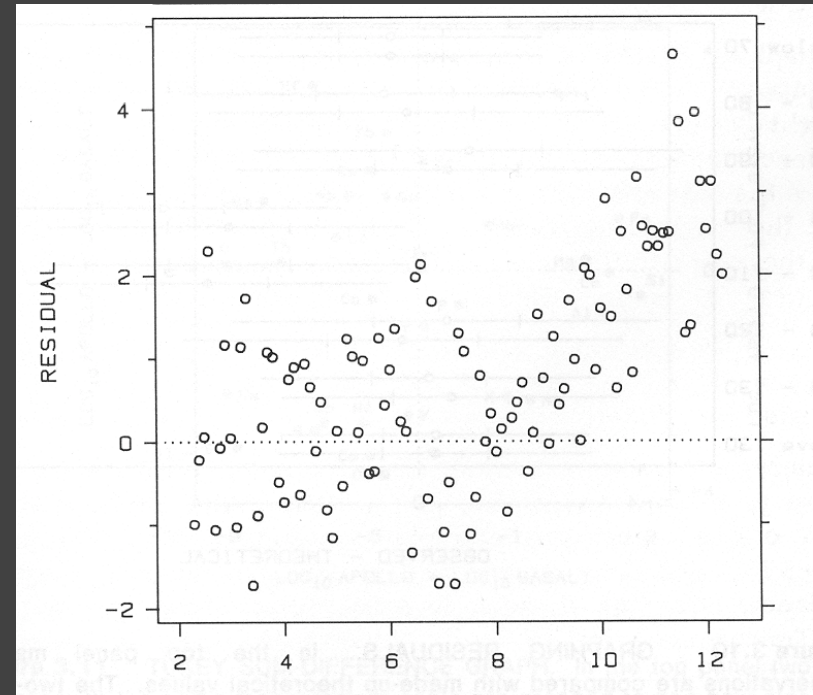
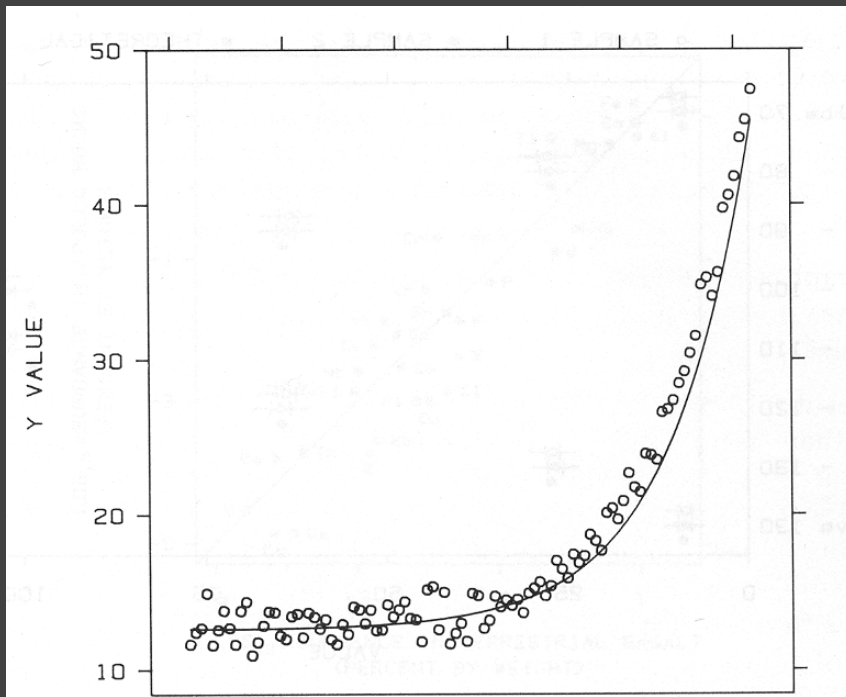
How well does the curve fit the data?



Plot the Residuals

Plot vertical distance from best fit curve

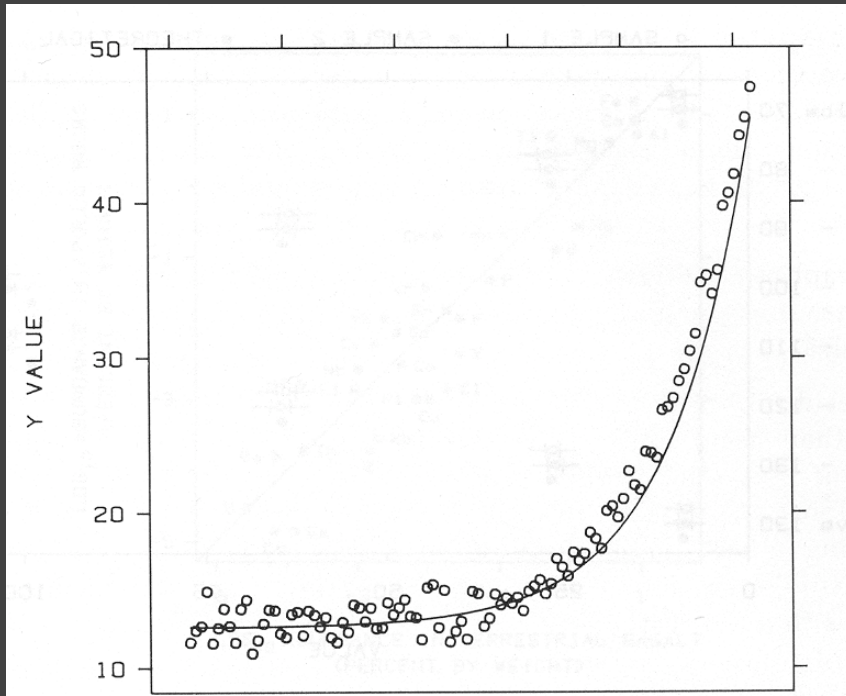
Residual graph shows accuracy of fit



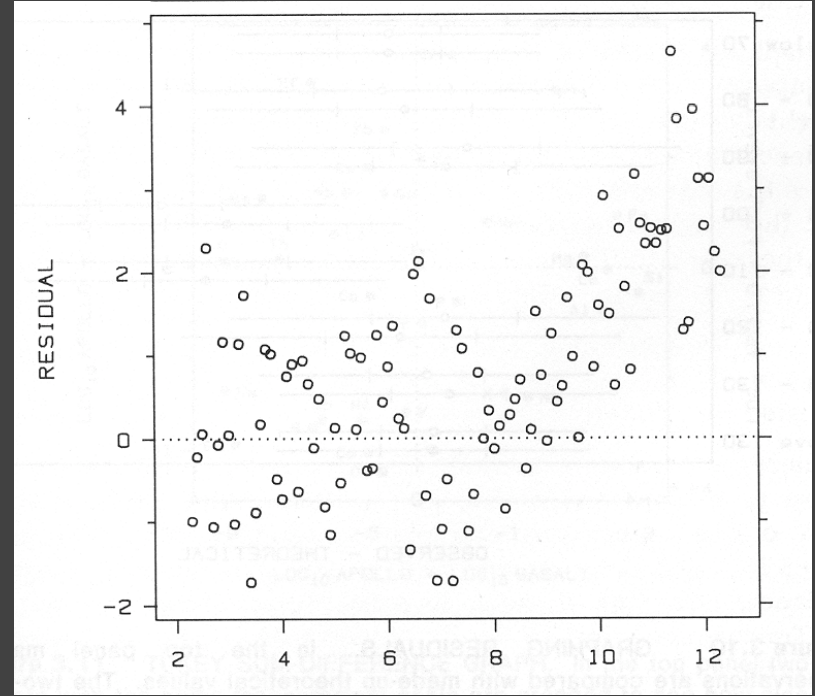
[Cleveland 85]

Multiple Plotting Options

Plot model in data space



Plot data in model space



[Cleveland 85]

Administrivia

A2: Exploratory Data Analysis

Use visualization software to form & answer questions

First steps:

Step 1: Pick domain & data

Step 2: Pose questions

Step 3: Profile the data

Iterate as needed

Create visualizations

Interact with data

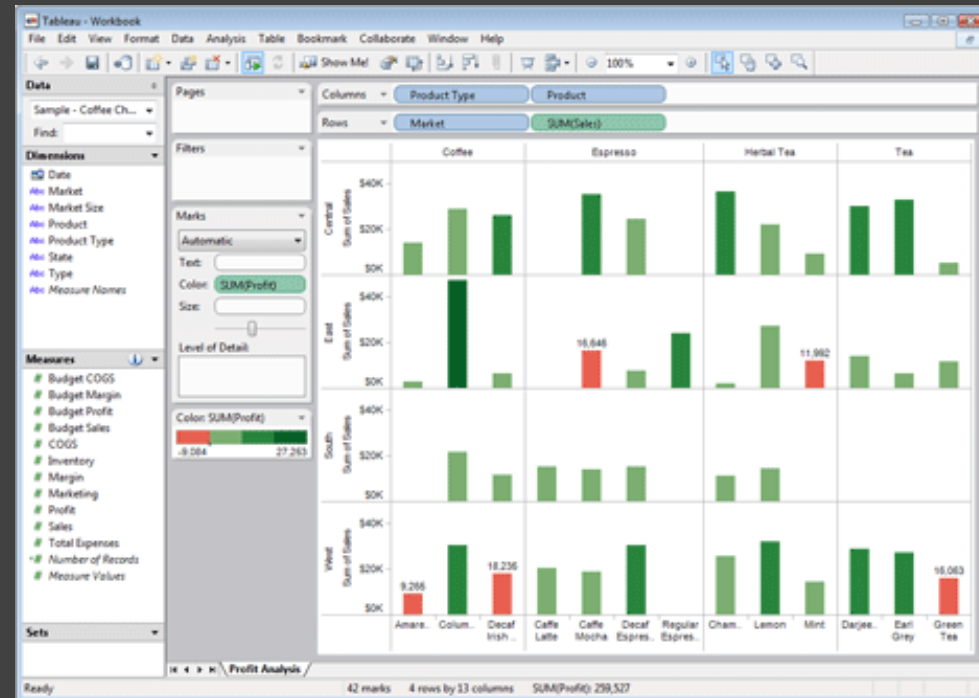
Refine your questions

Author a report

Images of annotated visualizations

(8+ images; min 4 views of the data)

Include titles and labels for each view



Due by 11:59pm
Tuesday, Jan 28

Multidimensional Data

Visual Encoding Variables

Position (X)

Position (Y)

Size

Value

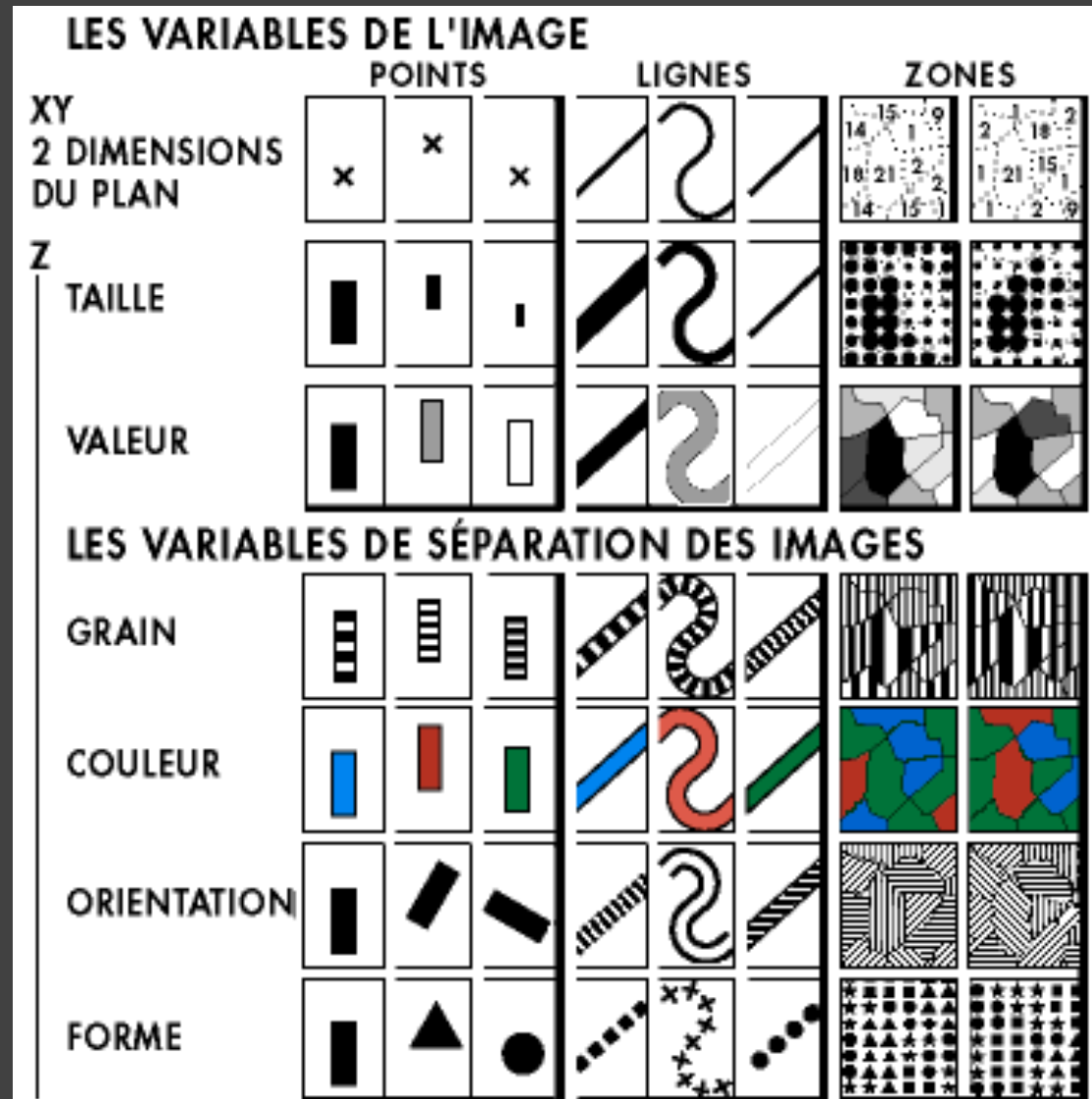
Texture

Color

Orientation

Shape

~8 dimensions?



Example: Coffee Sales

Sales figures for a fictional coffee chain

Sales	Q-Ratio
Profit	Q-Ratio
Marketing	Q-Ratio
Product Type	N {Coffee, Espresso, Herbal Tea, Tea}
Market	N {Central, East, South, West}

Filters

YEAR(Date): 2010

Marks

x+ Automatic

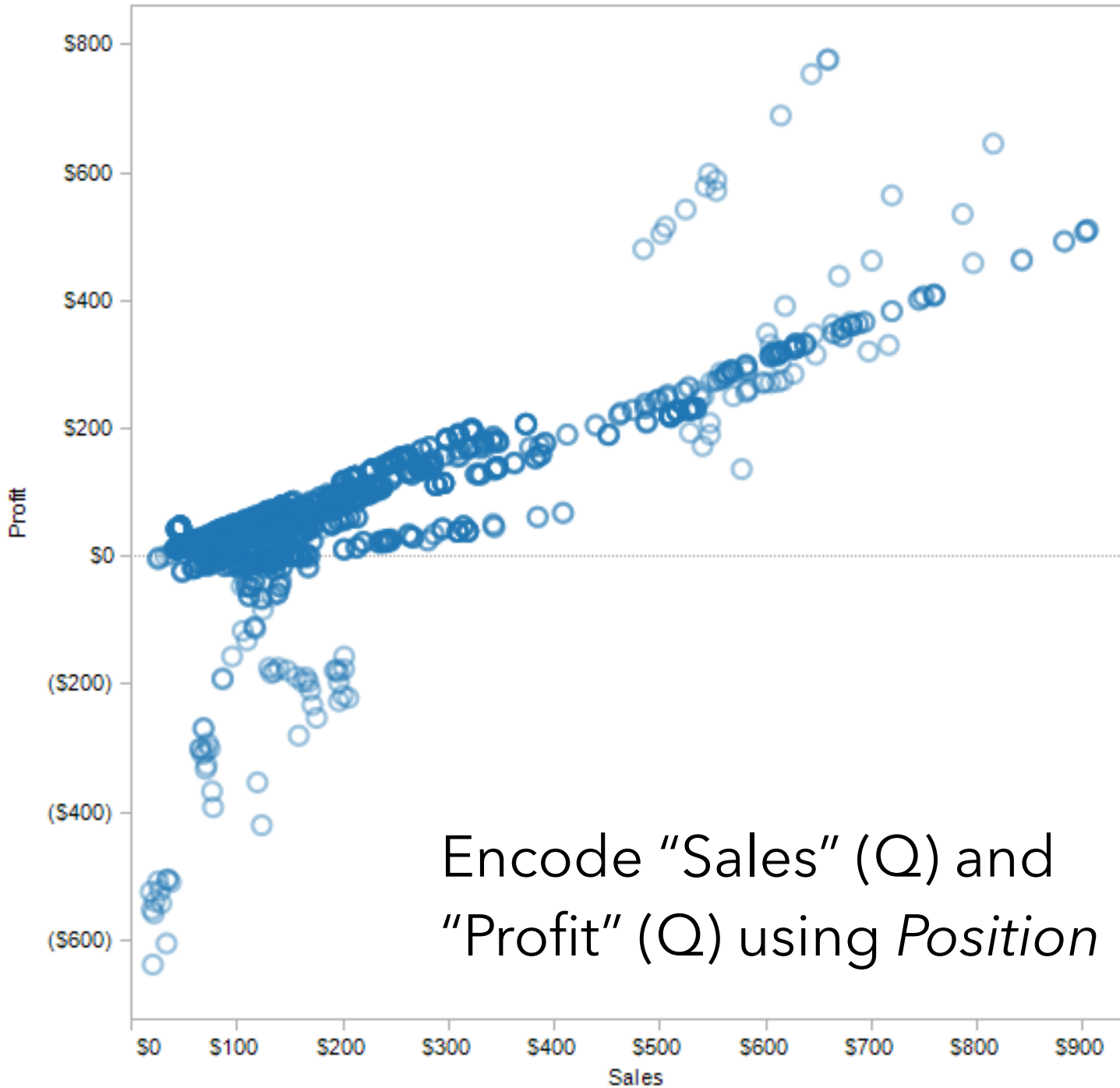
Shape Circle

Label

Color

Size

Level of Detail



Filters

YEAR(Date): 2010

Marks

x+ Automatic

Shape

Label

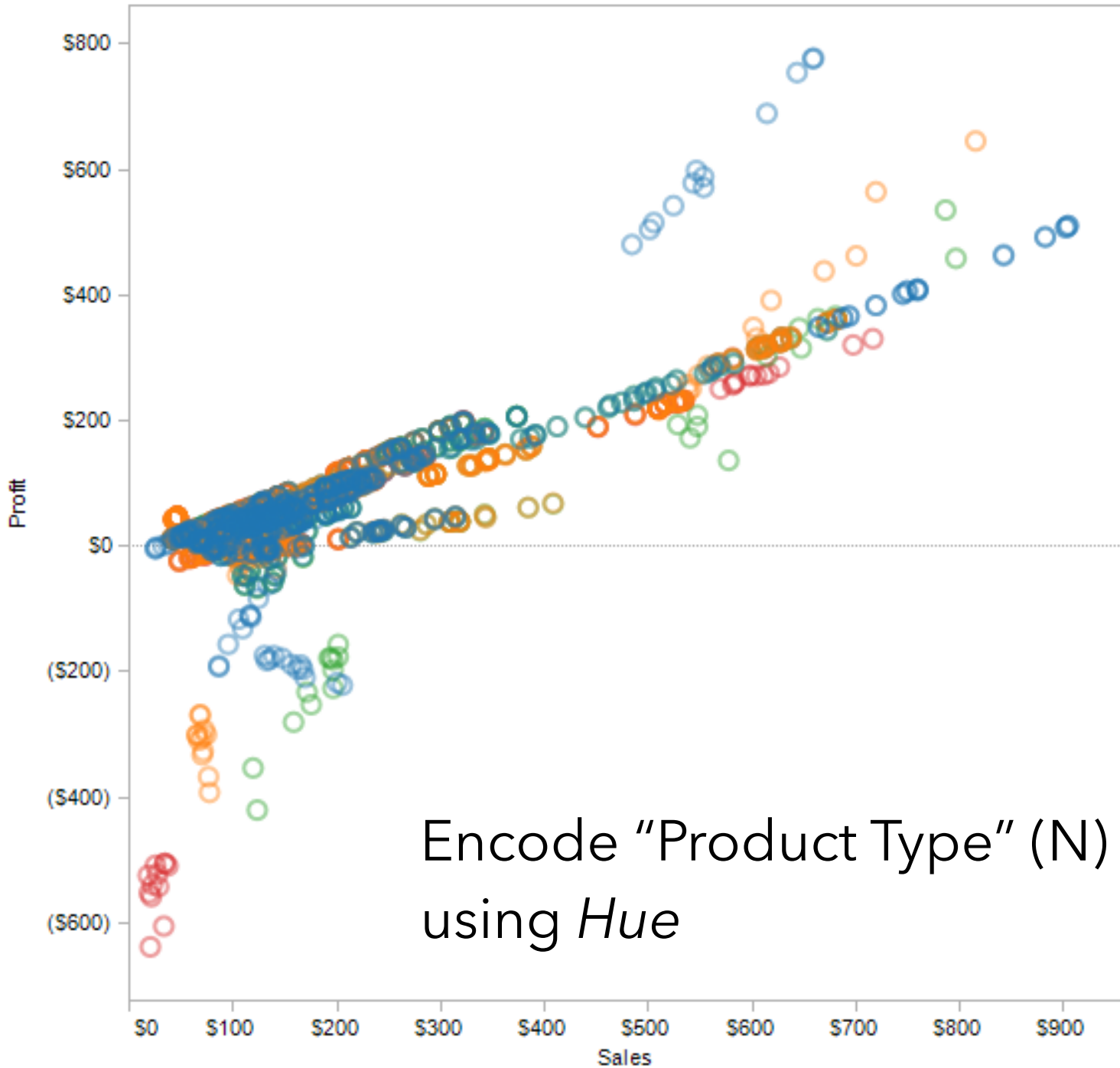
Color

Size

Level of Detail

Product Type

- Coffee
- Espresso
- Herbal Tea
- Tea



Filters

YEAR(Date): 2010

Marks

Automatic

Shape Market

Label Market

Color Product Type

Size

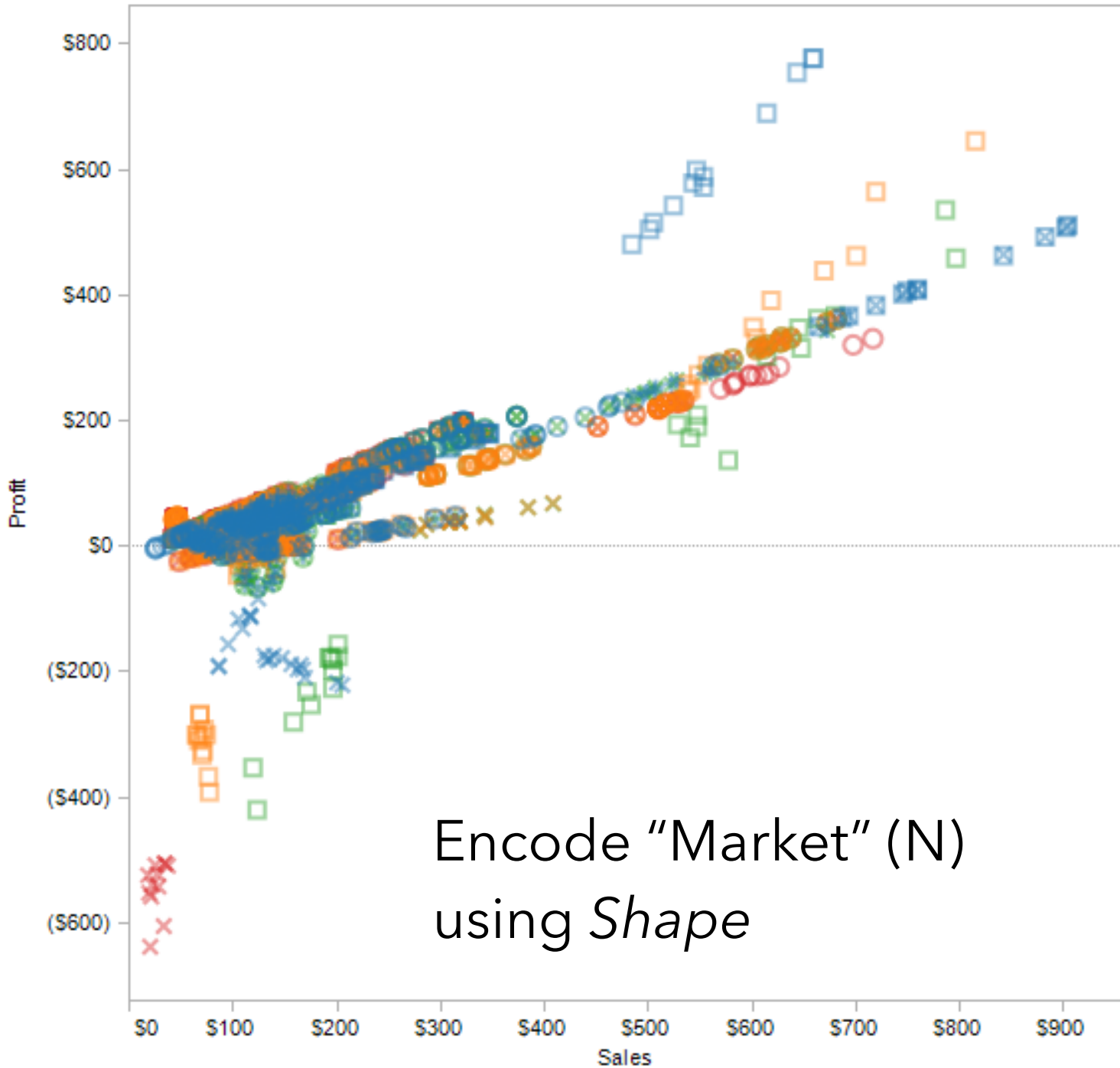
Level of Detail

Product Type

- Coffee
- Espresso
- Herbal Tea
- Tea

Market

- Central
- East
- South
- West



Filters

YEAR(Date): 2010

Marks

Automatic

Shape Market

Label

Color Product Type

Size Marketing

Marketing

Level of Detail

Product Type

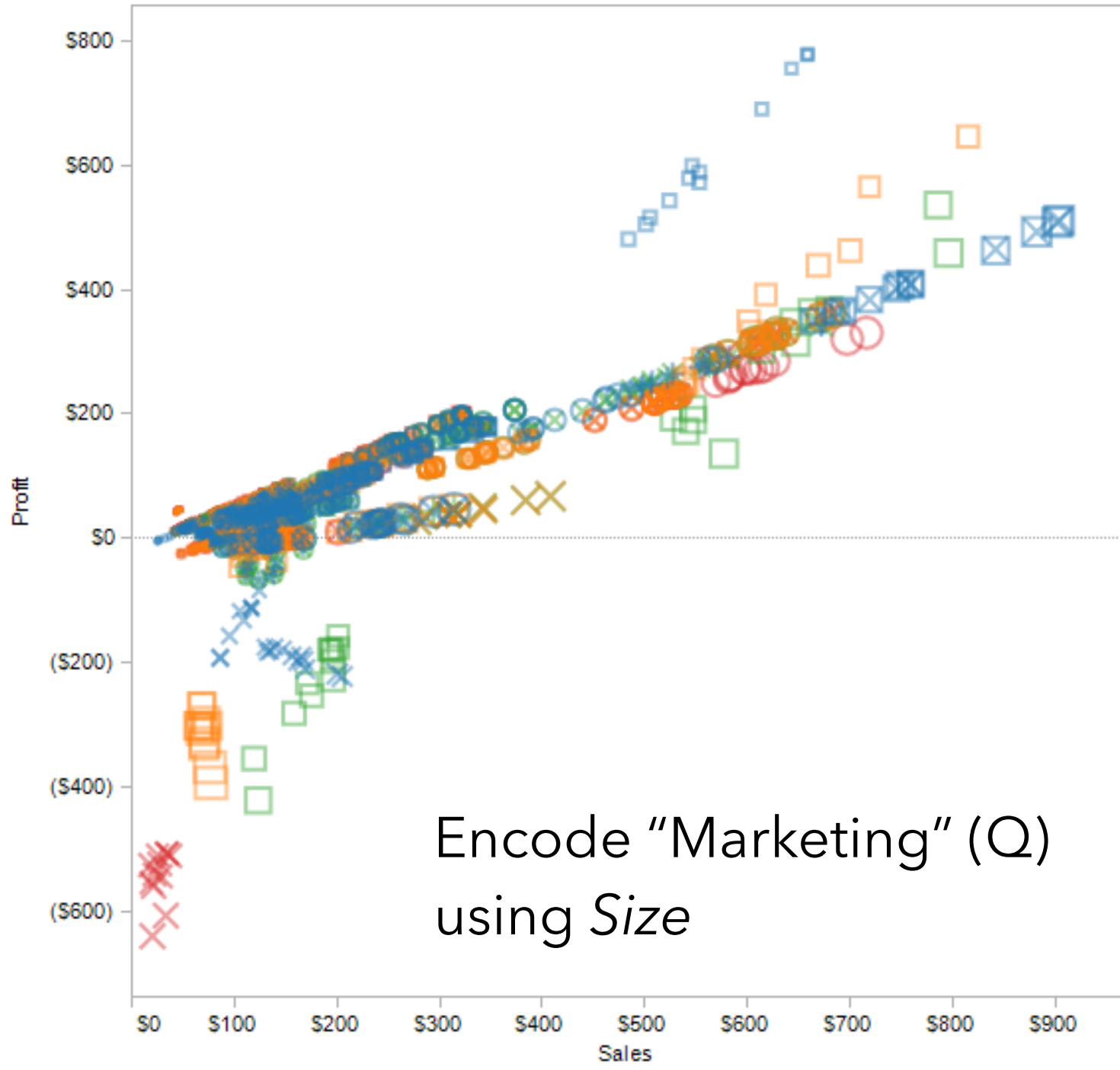
- Coffee
- Espresso
- Herbal Tea

Market

- Central
- East
- South

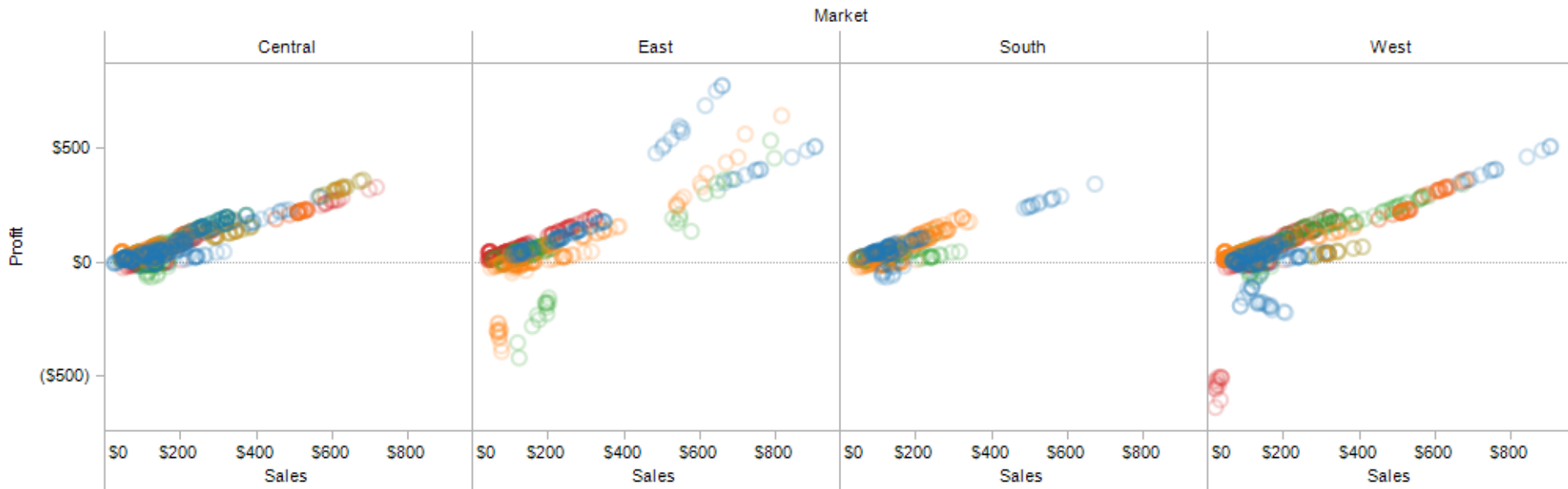
Marketing

- \$0
- \$50
- \$100



Encode "Marketing" (Q) using *Size*

Trellis Plots



A *trellis plot* subdivides space to enable comparison across multiple plots.

Typically nominal or ordinal variables are used as dimensions for subdivision.

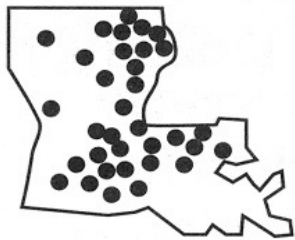
Small Multiples



[MacEachren '95, Figure 2.11, p. 38]

Small Multiples

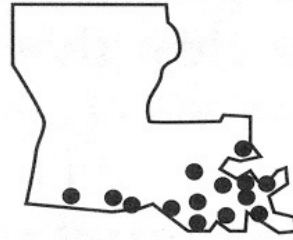
alfisol



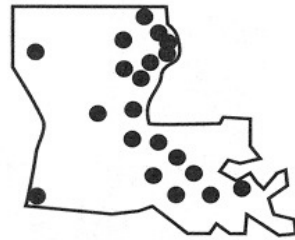
entisol



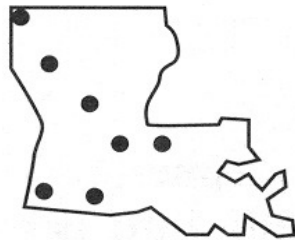
histosol



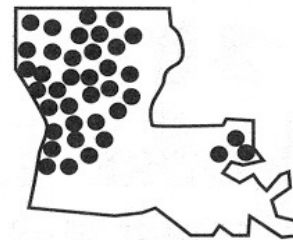
inceptisol



mollisol

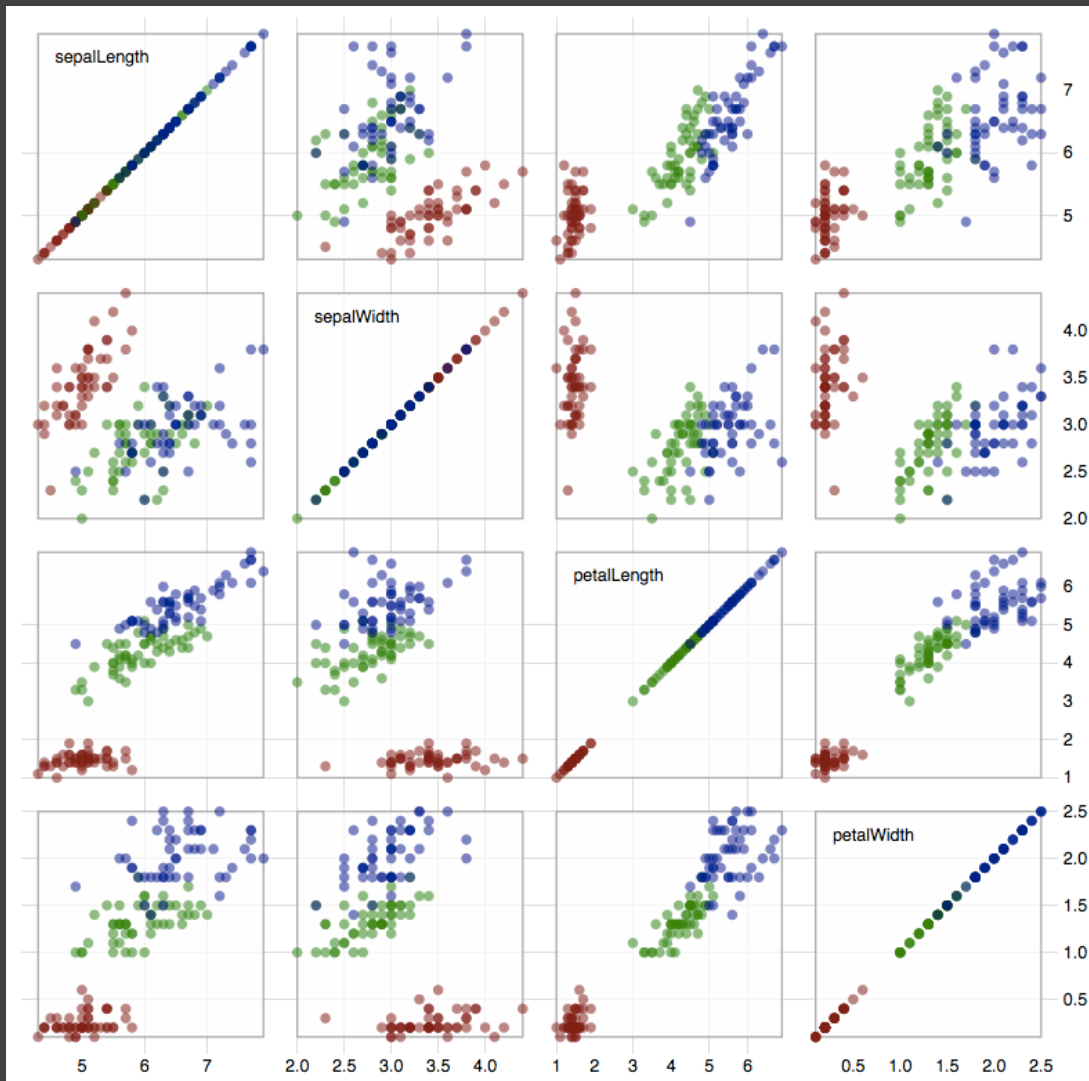


ultisol



[MacEachren '95, Figure 2.11, p. 38]

Scatterplot Matrix (SPLOM)



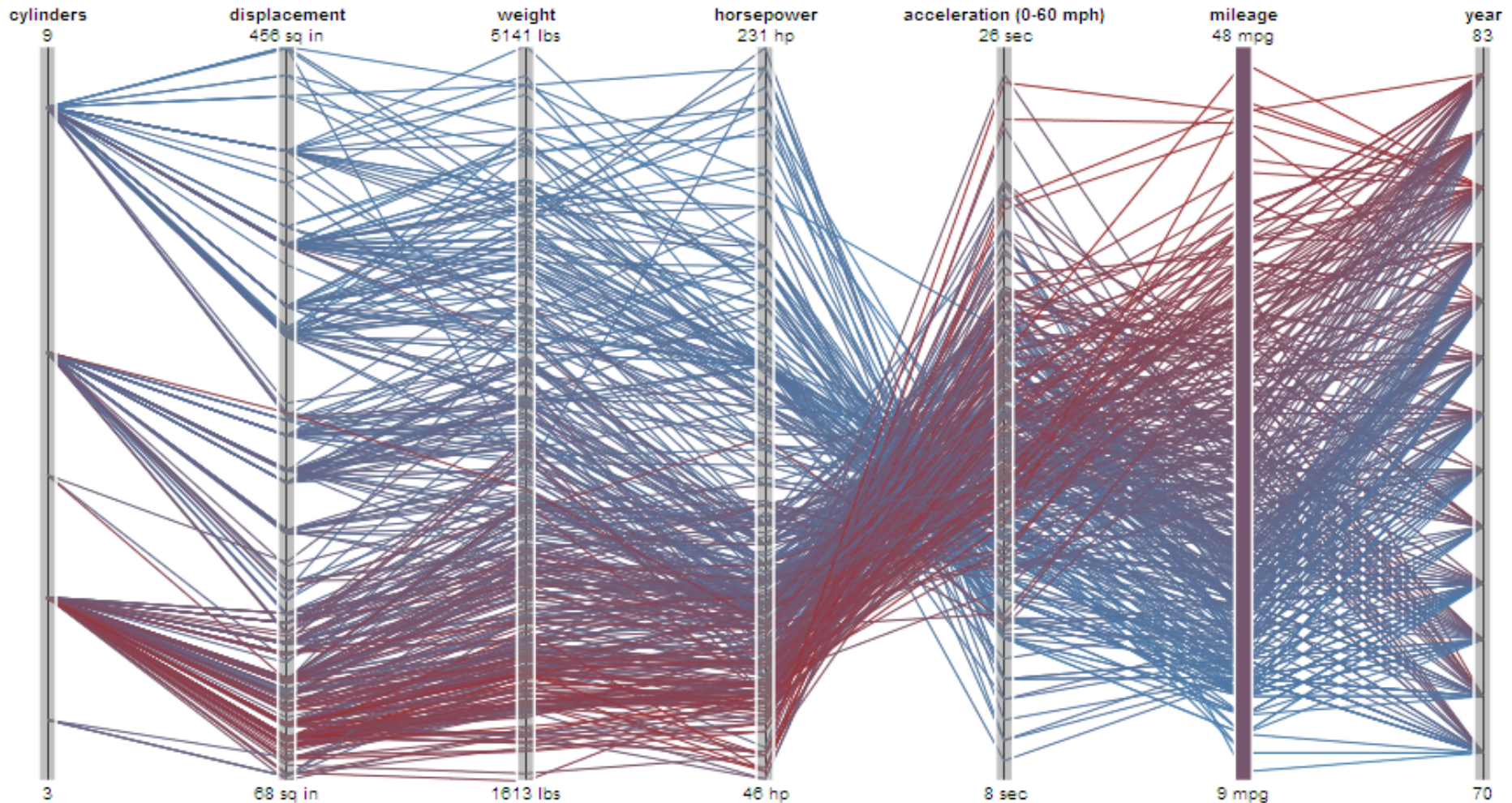
Scatter plots for pairwise comparison of each data dimension.

Multiple Coordinated Views



Parallel Coordinates

Parallel Coordinates [Inselberg]



Parallel Coordinates [Inselberg]

Visualize up to ~two dozen dimensions at once

1. Draw parallel axes for each variable
2. For each tuple, connect points on each axis

Between adjacent axes: line crossings imply neg. correlation, shared slopes imply pos. correlation.

Full plot can be cluttered. **Interactive selection** can be used to assess multivariate relationships.

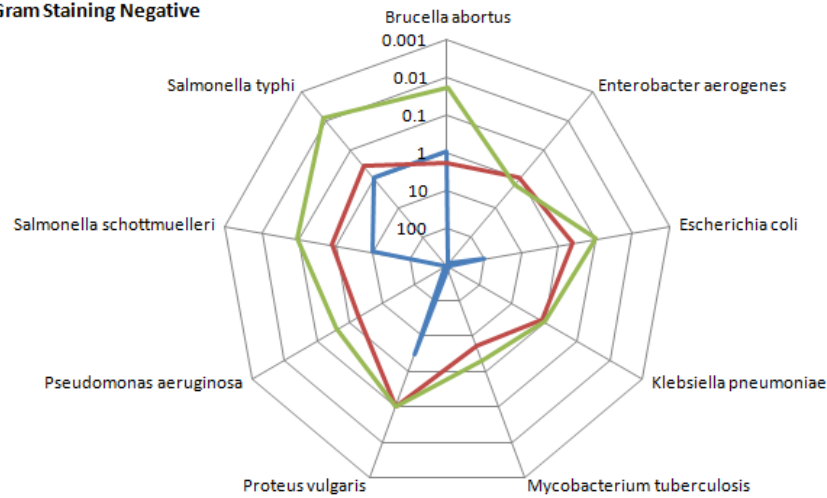
Highly sensitive to axis **scale** and **ordering**.

Expertise required to use effectively!

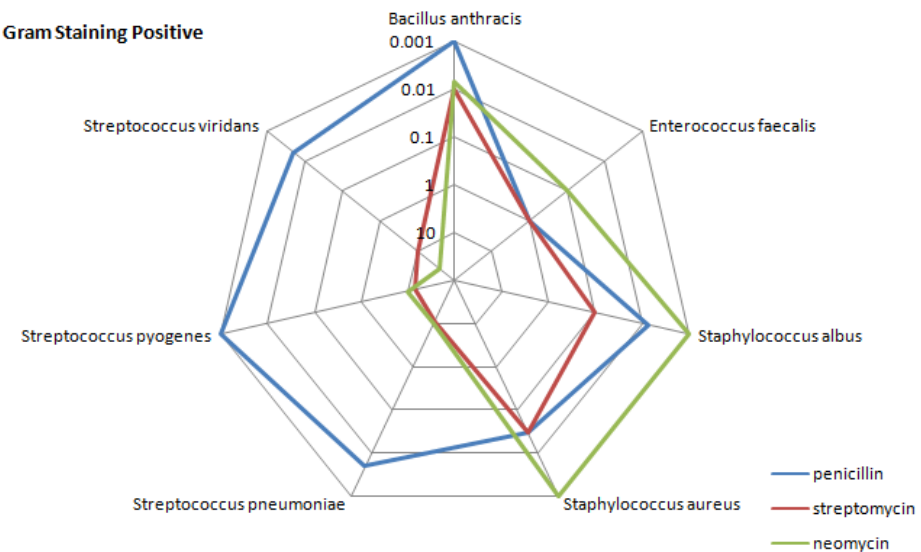
Radar Plot / Star Graph

Antibiotics MIC Concentrations

Gram Staining Negative



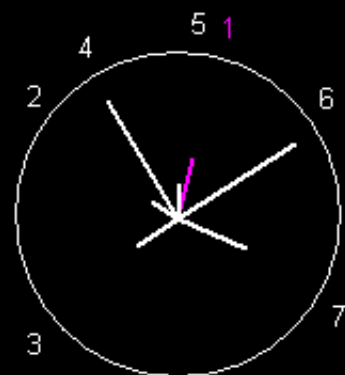
Gram Staining Positive



“Parallel” dimensions in polar coordinate space
Best if same units apply to each axis

Dimensionality Reduction

Dimensionality Reduction



<http://www.ggobi.org/>

1:0.099,0.367(243.00)

2:-0.157,0.106(47.74)

3:-0.251,-0.178(9.00)

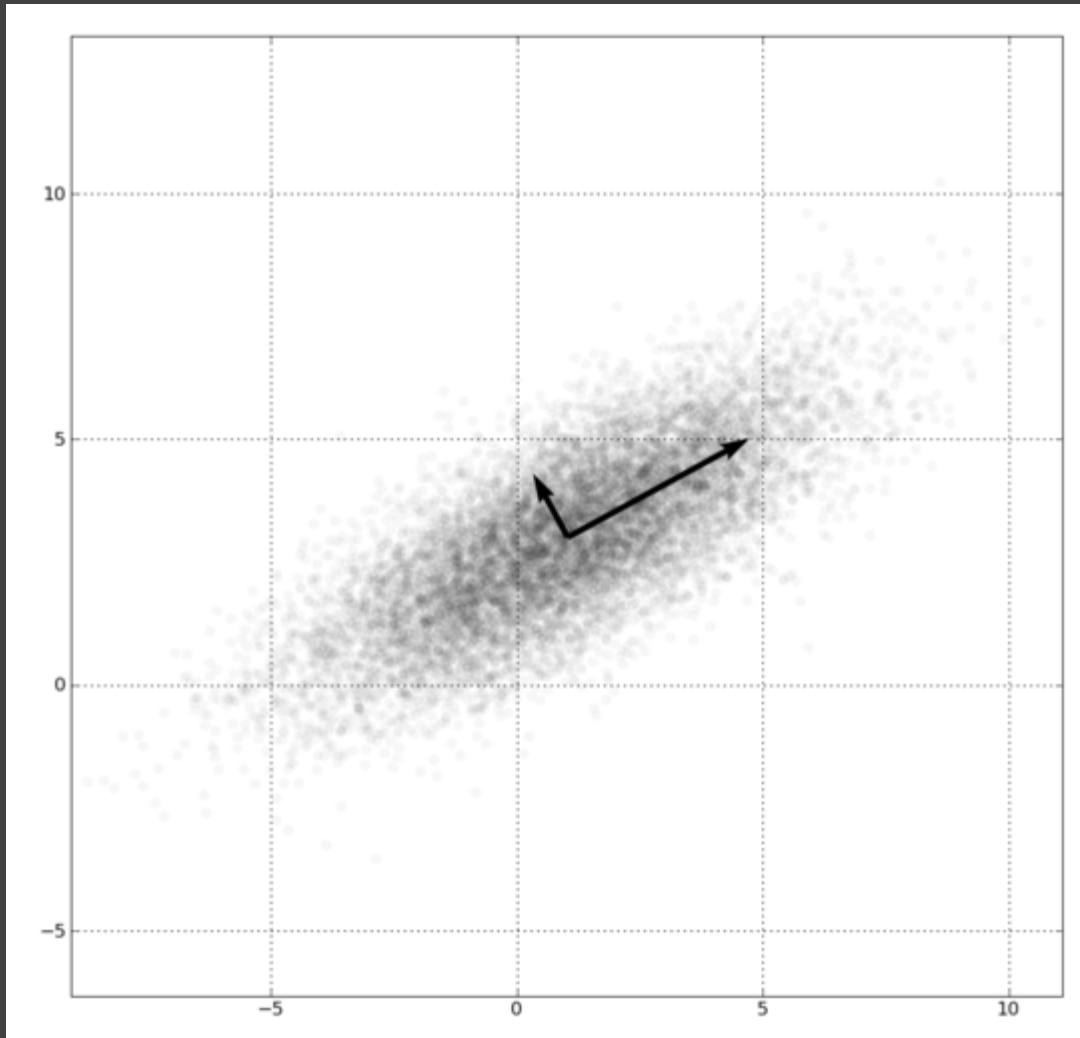
4:-0.442,0.723(1.00)

5:0.016,0.222(1.00)

6:0.726,0.461(3.00)

7:0.424,-0.195(1.00)

Principal Components Analysis



1. Mean-center the data.
2. Find \perp basis vectors that maximize the data variance.
3. Plot the data using the top vectors.

PCA of Genomes [Demiralp et al. '13]



Many Reduction Techniques!

General Strategies:

Matrix Factorization

Nearest Neighbor (Topological) Methods

Popular Techniques:

Principal Components Analysis (PCA)

t-Dist. Stochastic Neighbor Embedding (t-SNE)

Uniform Manifold Approx. & Projection (UMAP)

The Beginner's Guide to Dimensionality Reduction

Explore the methods that data scientists use to visualize high-dimensional data.

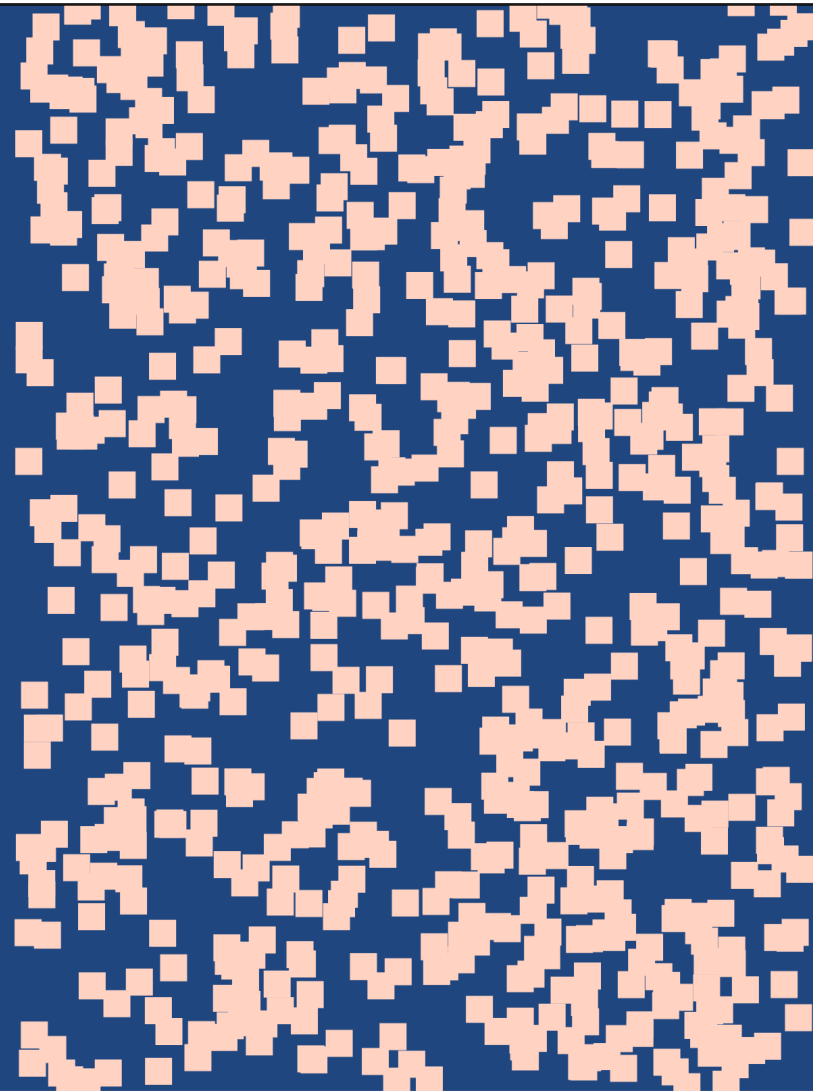
By: [Matthew Conlen](#) and [Fred Hohman](#)

July 16, 2018

Dimensionality reduction is a powerful technique used by data scientists to look for hidden structure in data. The method is useful in a number of domains, for example document categorization, protein disorder prediction, and machine learning model debugging^[2].

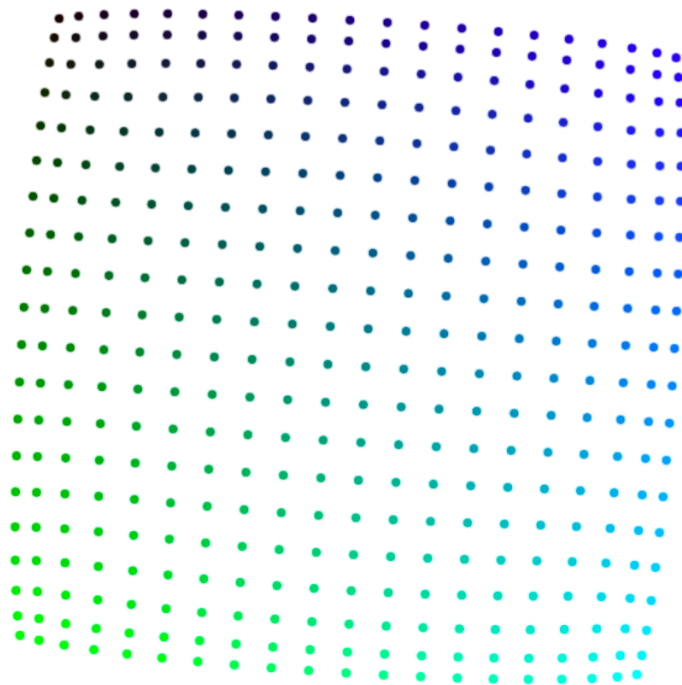
The results of a dimensionality reduction algorithm can be visualized to reveal patterns and clusters of similar or dissimilar data. Even though the data is displayed in only two or three dimensions, structures present in higher dimensions are maintained, at least roughly^[7].

The technique is available in many applications, for



How to Use t-SNE Effectively

Although extremely useful for visualizing high-dimensional data, t-SNE plots can sometimes be mysterious or misleading. By exploring how it behaves in simple cases, we can learn to use it more effectively.



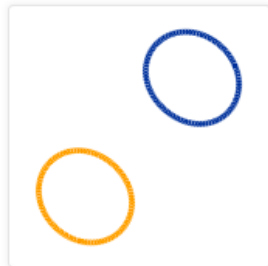
⏸ ↻ Step 1,910
 Points Per Side 20
 Perplexity 10
 Epsilon 5

A square grid with equal spacing between points. Try convergence at different sizes.

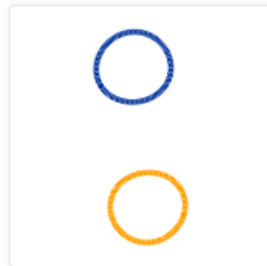
Visualizing t-SNE [Wattenberg et al. '16]



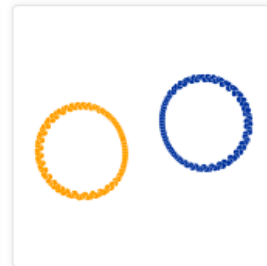
Original



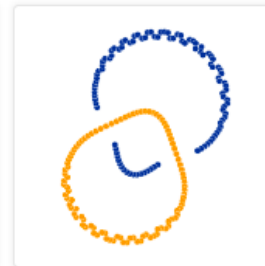
Perplexity: 2
Step: 5,000



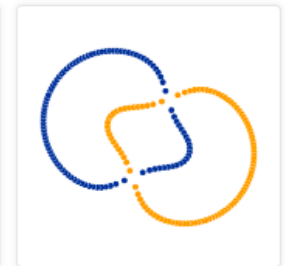
Perplexity: 5
Step: 5,000



Perplexity: 30
Step: 5,000



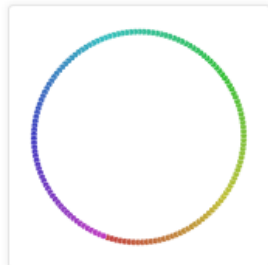
Perplexity: 50
Step: 5,000



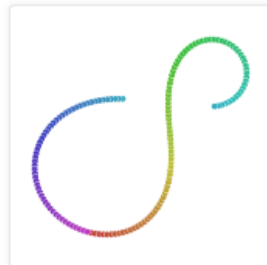
Perplexity: 100
Step: 5,000



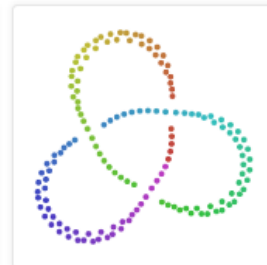
Original



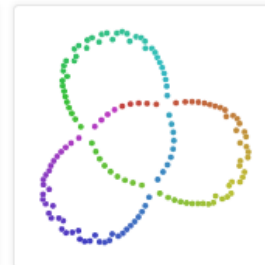
Perplexity: 2
Step: 5,000



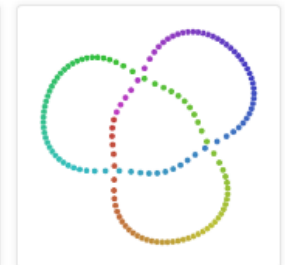
Perplexity: 5
Step: 5,000



Perplexity: 30
Step: 5,000



Perplexity: 50
Step: 5,000



Perplexity: 100
Step: 5,000

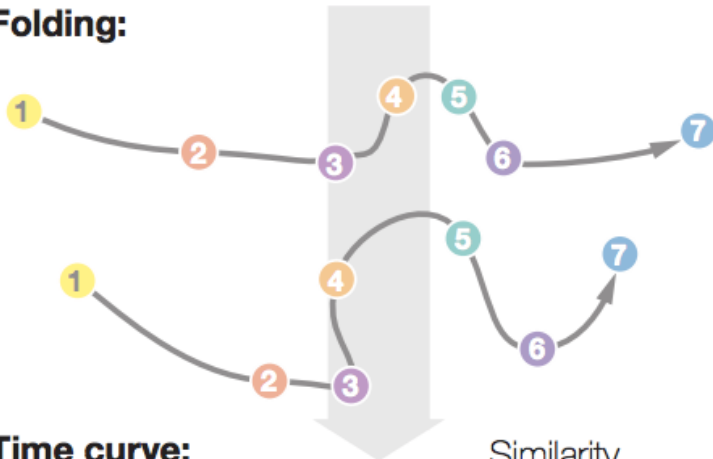
Time Curves [Bach et al. '16]

Timeline:



Circles are data cases with a time stamp.
Similar colors indicate similar data cases.

Folding:

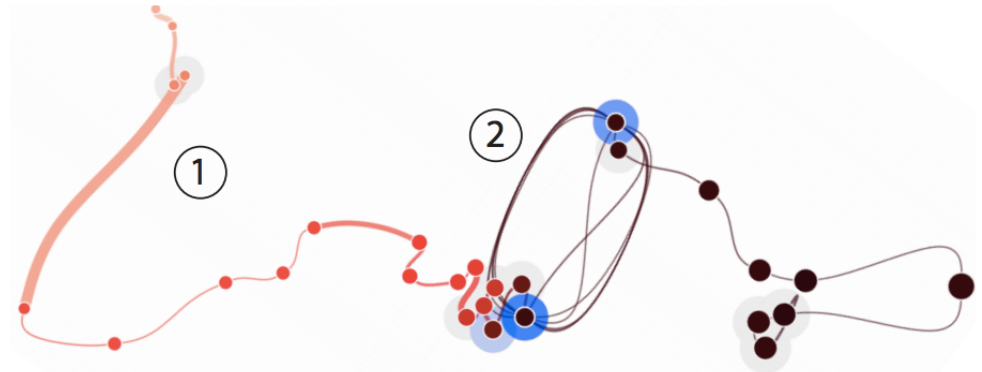


Time curve:

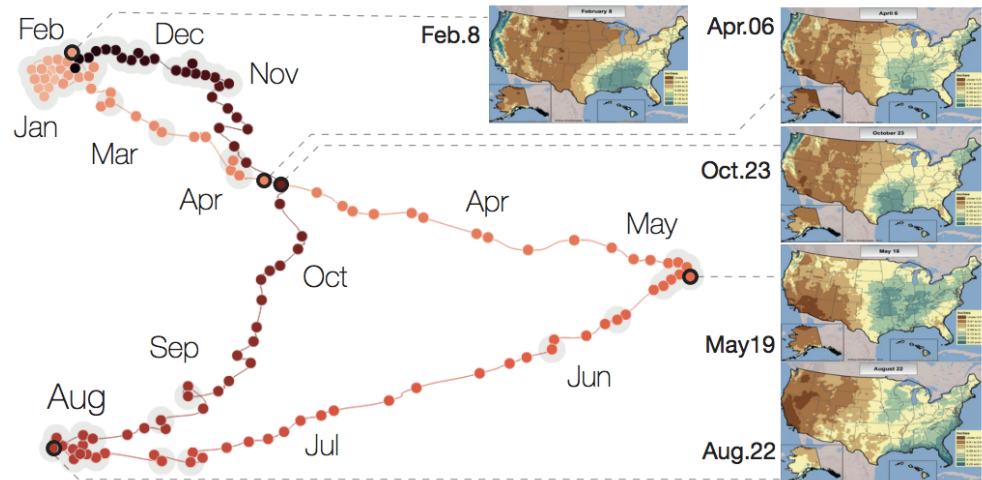


The temporal ordering of data cases is preserved.
Spatial proximity now indicates similarity.

(a) Folding time



Wikipedia "Chocolate" Article



U.S. Precipitation over 1 Year

Visual Encoding Design

Use **expressive** and **effective** encodings

Avoid **over-encoding**

Reduce the problem space

Use **space** and **small multiples** intelligently

Use **interaction** to generate *relevant* views

Rarely does a single visualization answer all questions. Instead, the ability to generate appropriate visualizations quickly is critical!