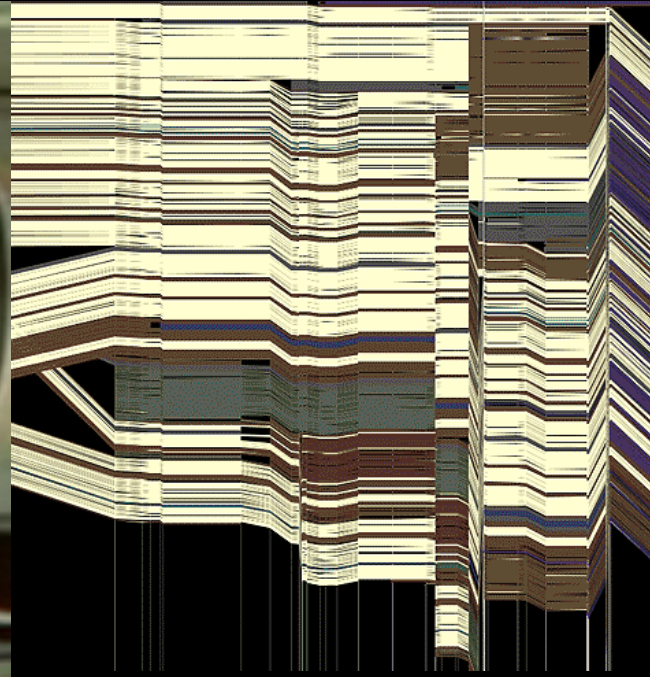
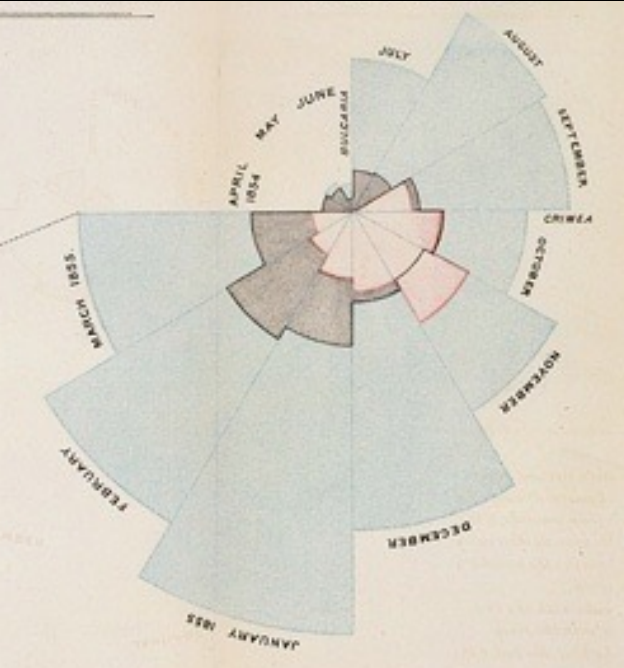


CSE 512 - Data Visualization

Multidimensional Vis



Jeffrey Heer University of Washington

Last Time:
Exploratory Data Analysis



Exposure, the effective laying open of the data to display the unanticipated, is to us a major portion of data analysis. Formal statistics has given almost no guidance to exposure; indeed, it is not clear how the **informality** and **flexibility** appropriate to the **exploratory character of exposure** can be fitted into any of the structures of formal statistics so far proposed.

Graph Viewer

Roll-up by:

All

Visualization:

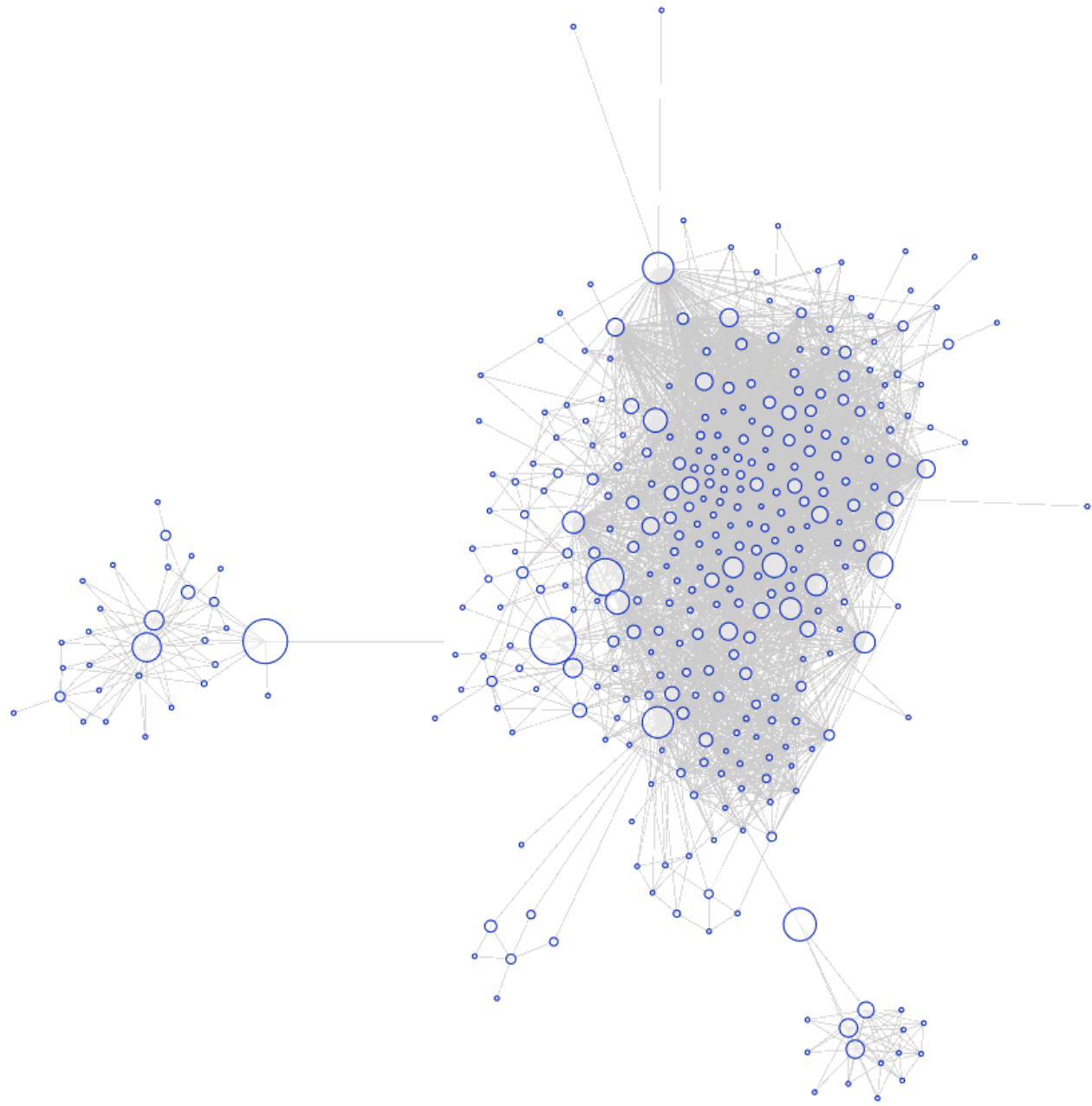
Node-Link

Sort by:

None

Edge centrality filters:

Two horizontal sliders for edge centrality filtering.



- Images
- Animate

Graph Viewer

Roll-up by:

All

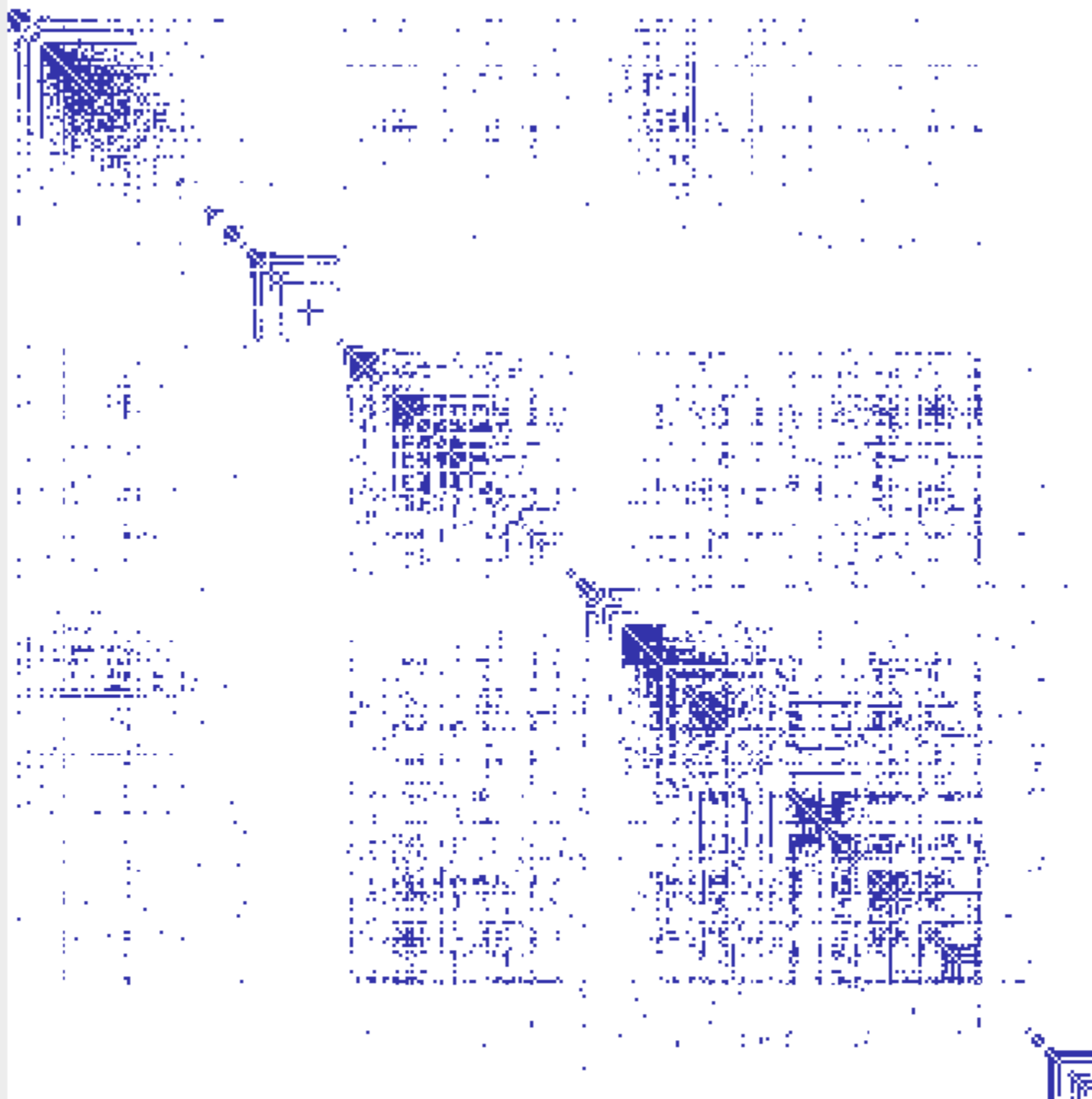
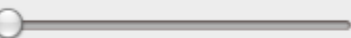
Visualization:

Matrix

Sort by:

Linkage

Edge centrality filters:



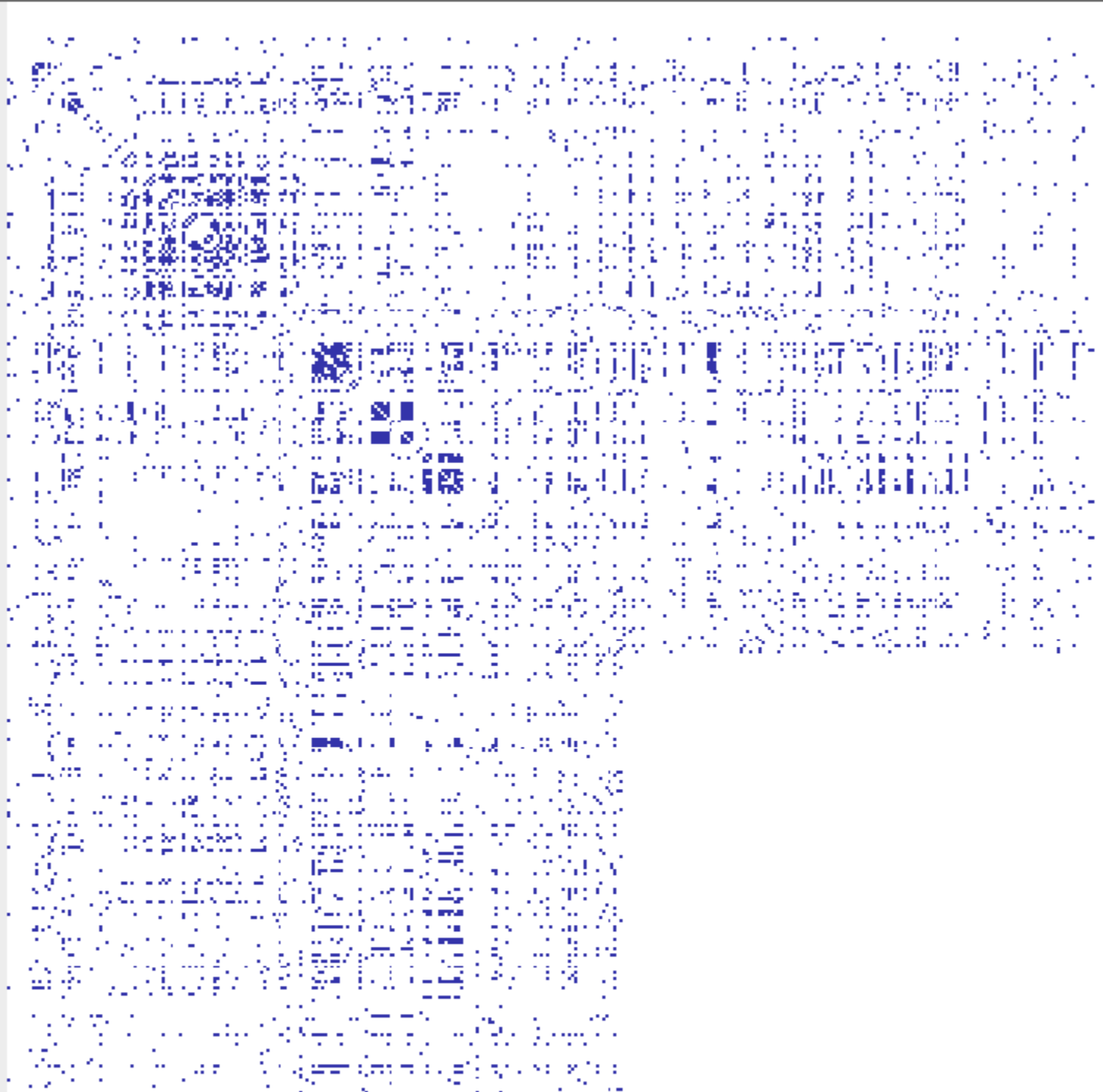
Graph Viewer

Roll-up by:

Visualization:

Sort by:

Edge centrality filters:

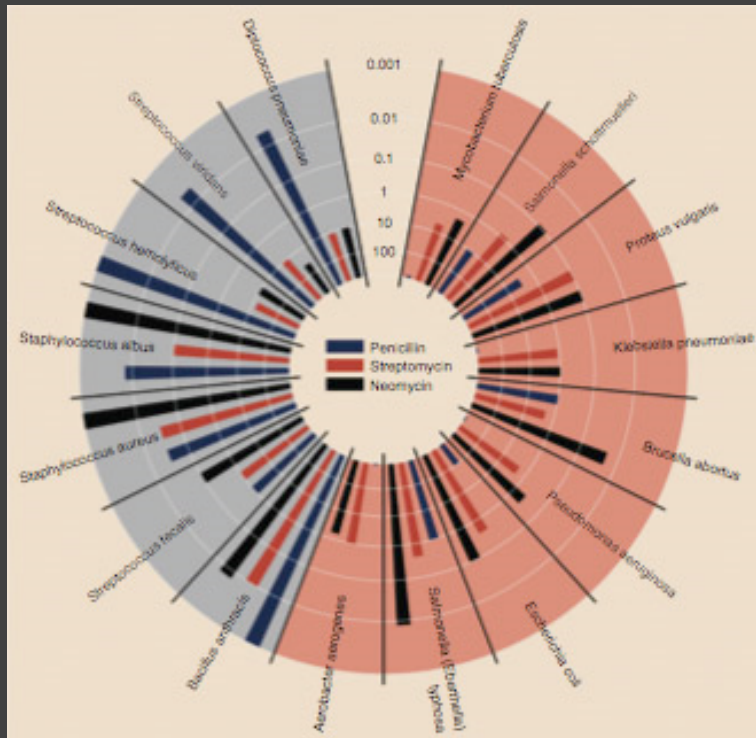


Antibiotic Effectiveness

Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

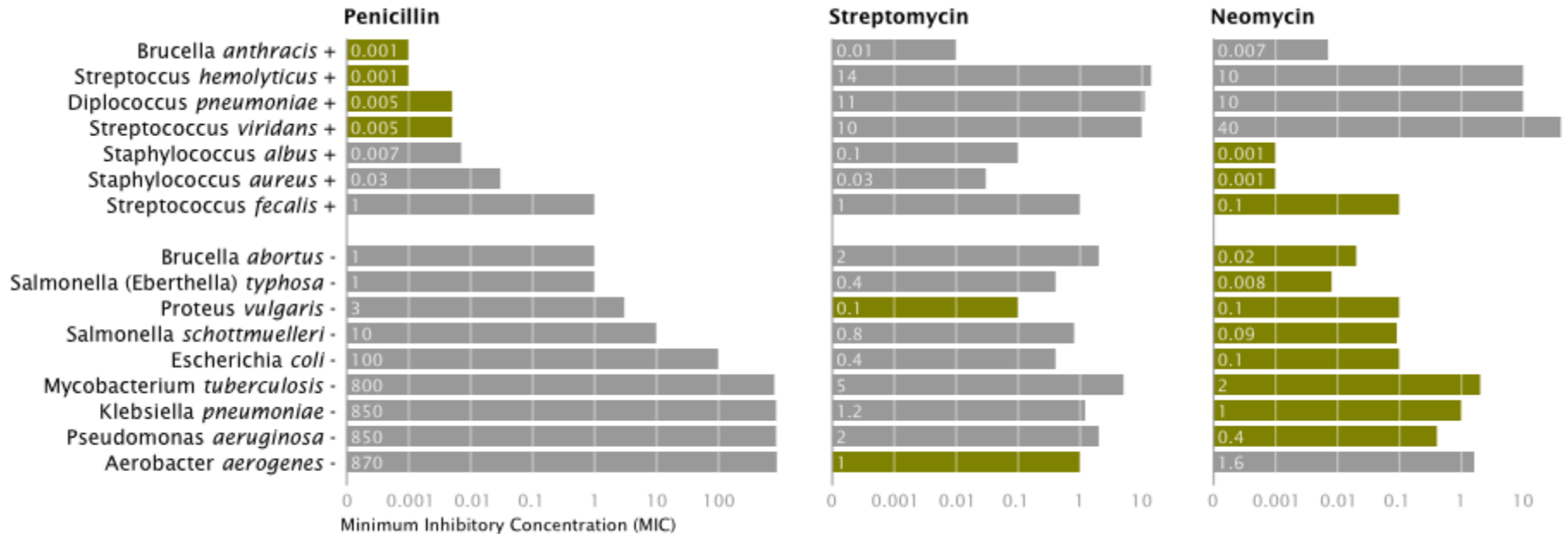
How do the drugs compare?

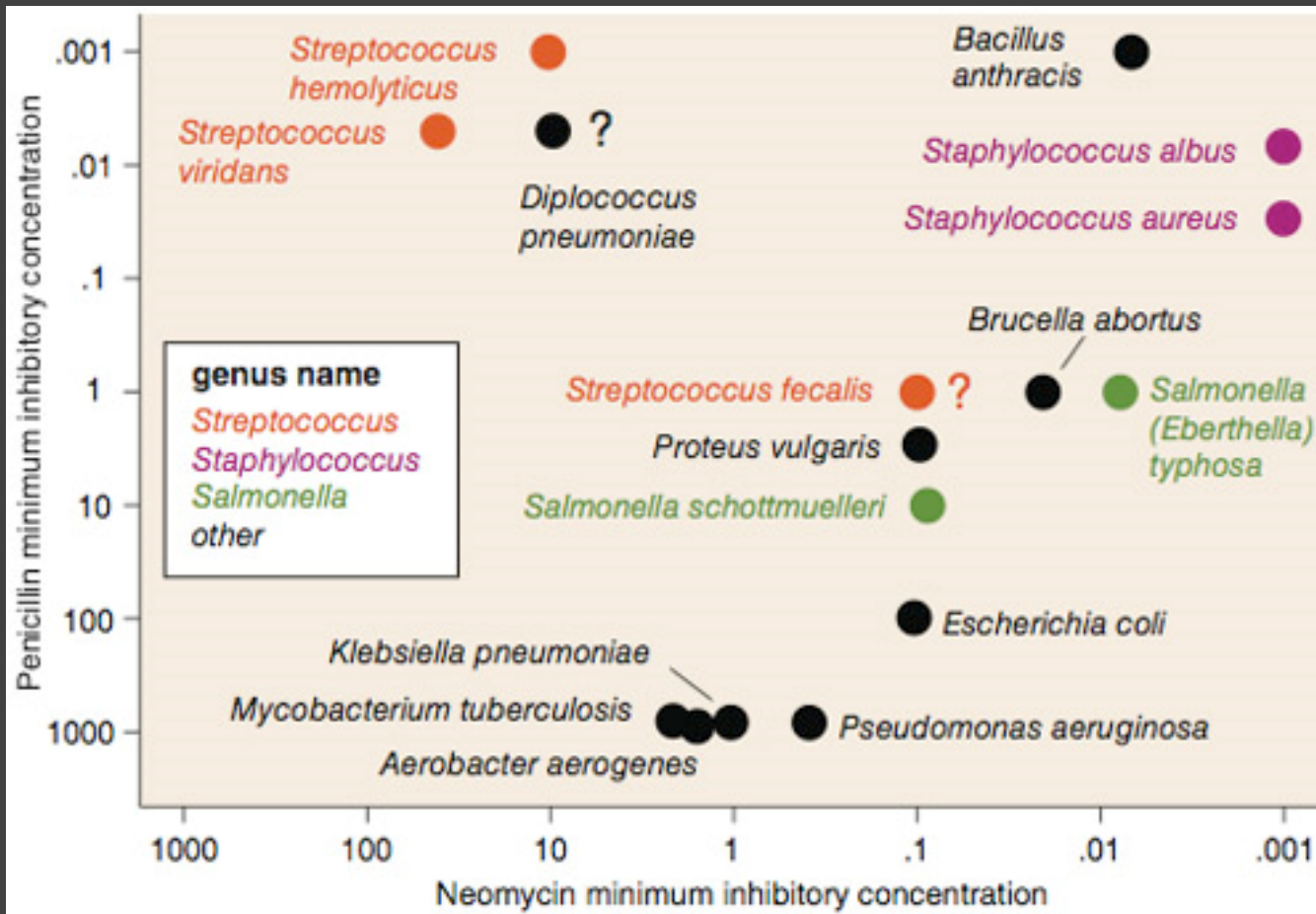


Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	1	-
<i>Mycobacterium tuberculosis</i>	800	5	2	-
<i>Proteus vulgaris</i>	3	0.1	0.1	-
<i>Pseudomonas aeruginosa</i>	850	2	0.4	-
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	-
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	-
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

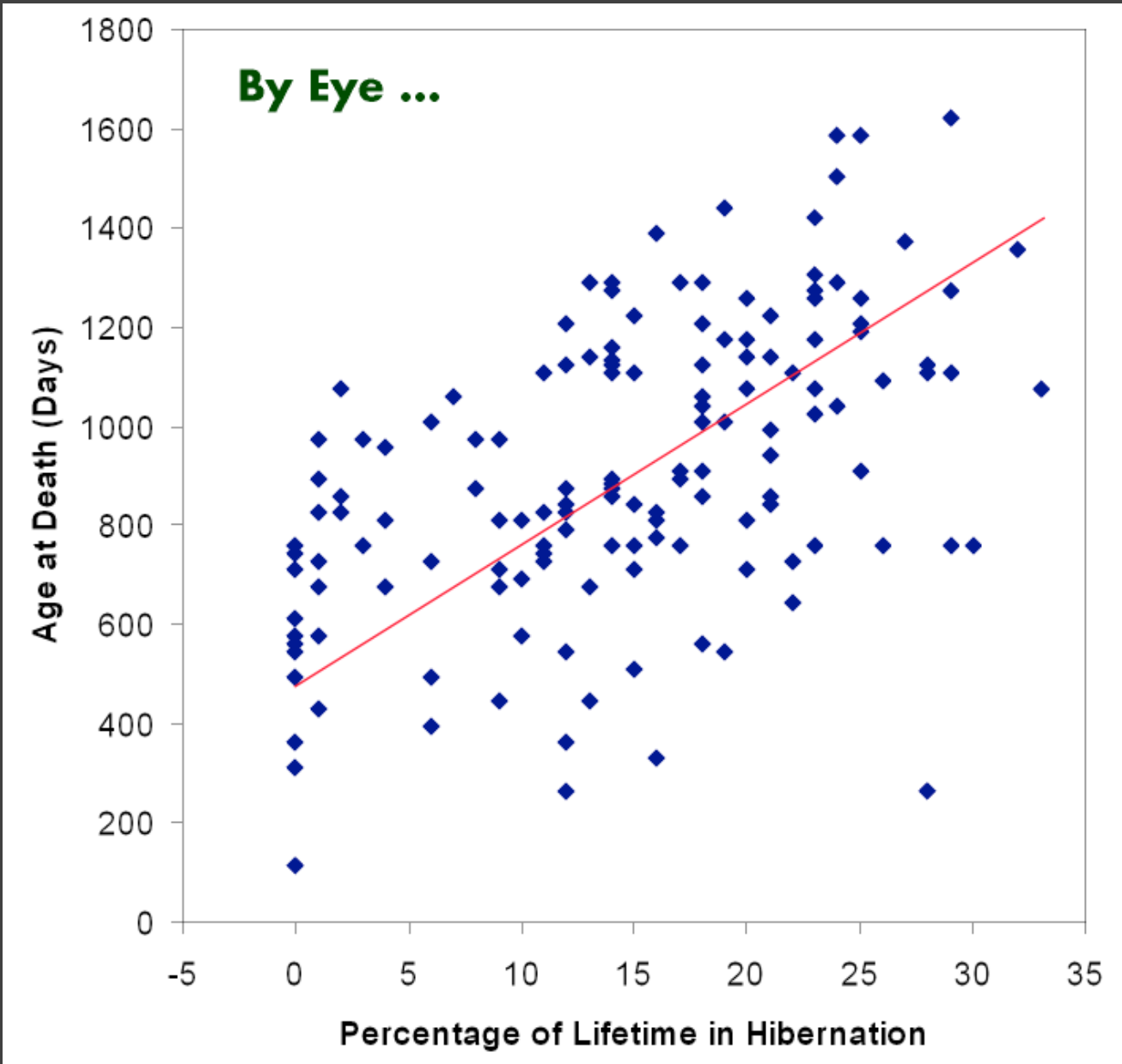
Original graphic by Will Burtin, 1951

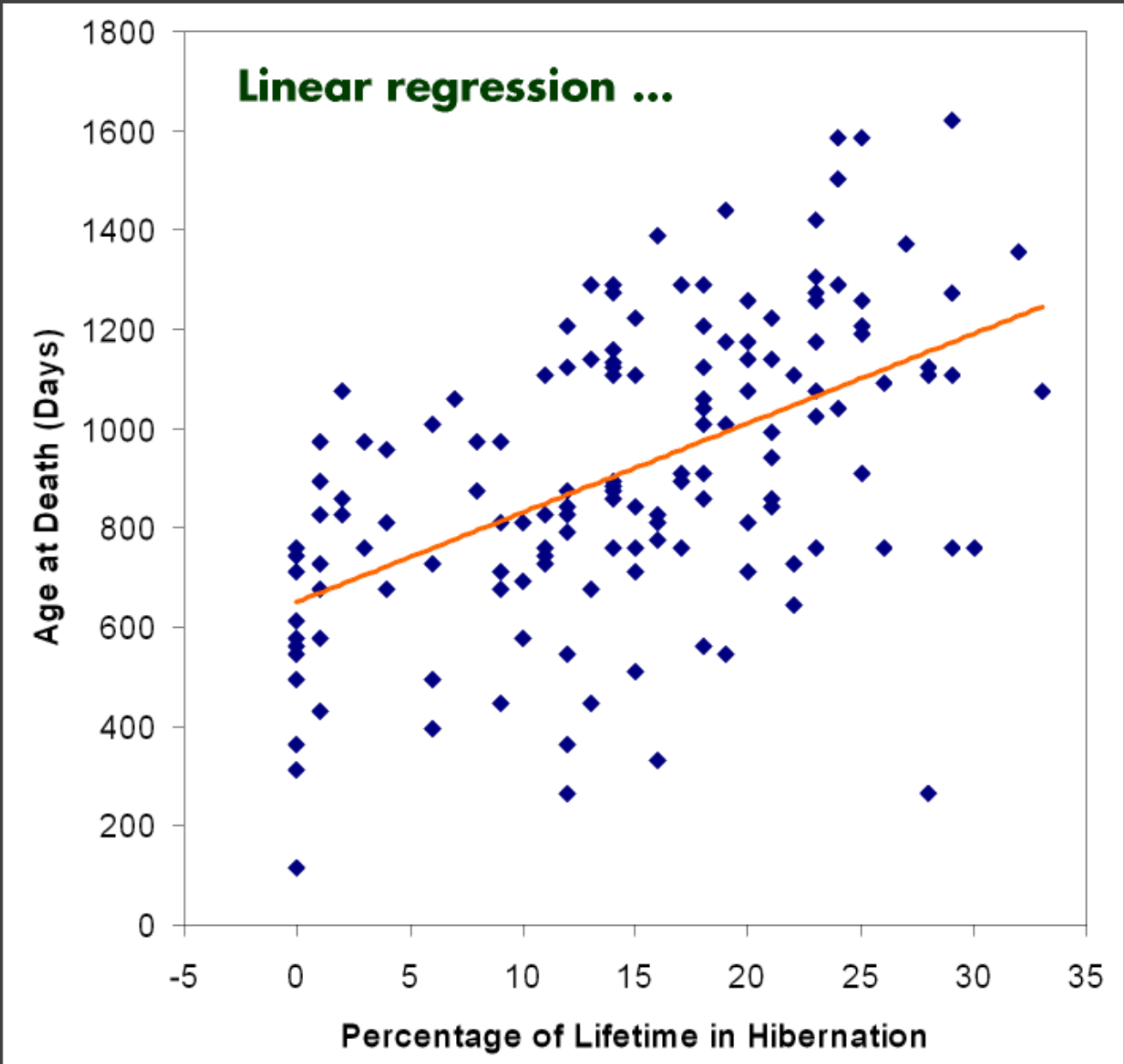
How do the drugs compare?



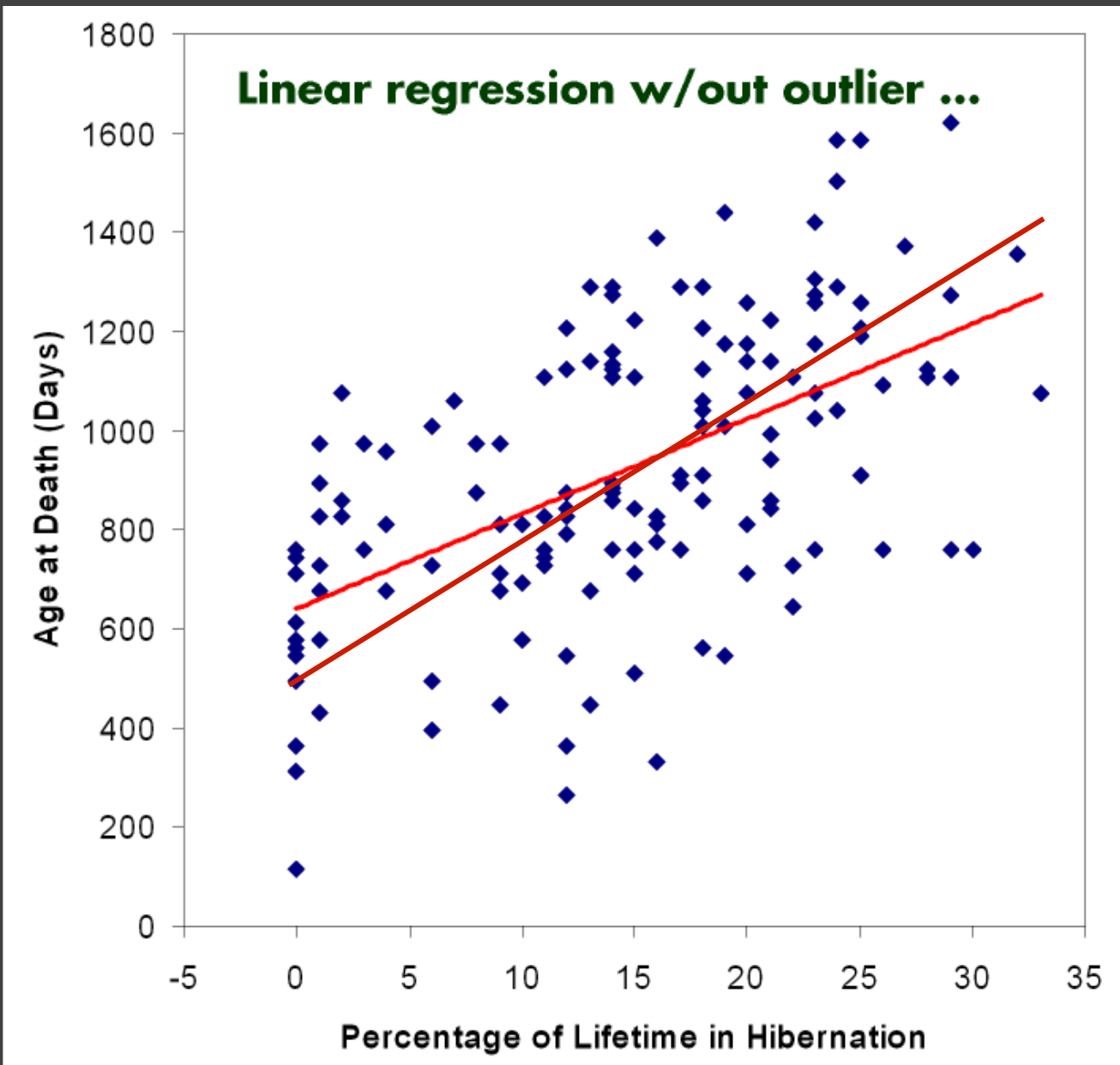


Do the bacteria group by resistance?
 Do different drugs correlate?





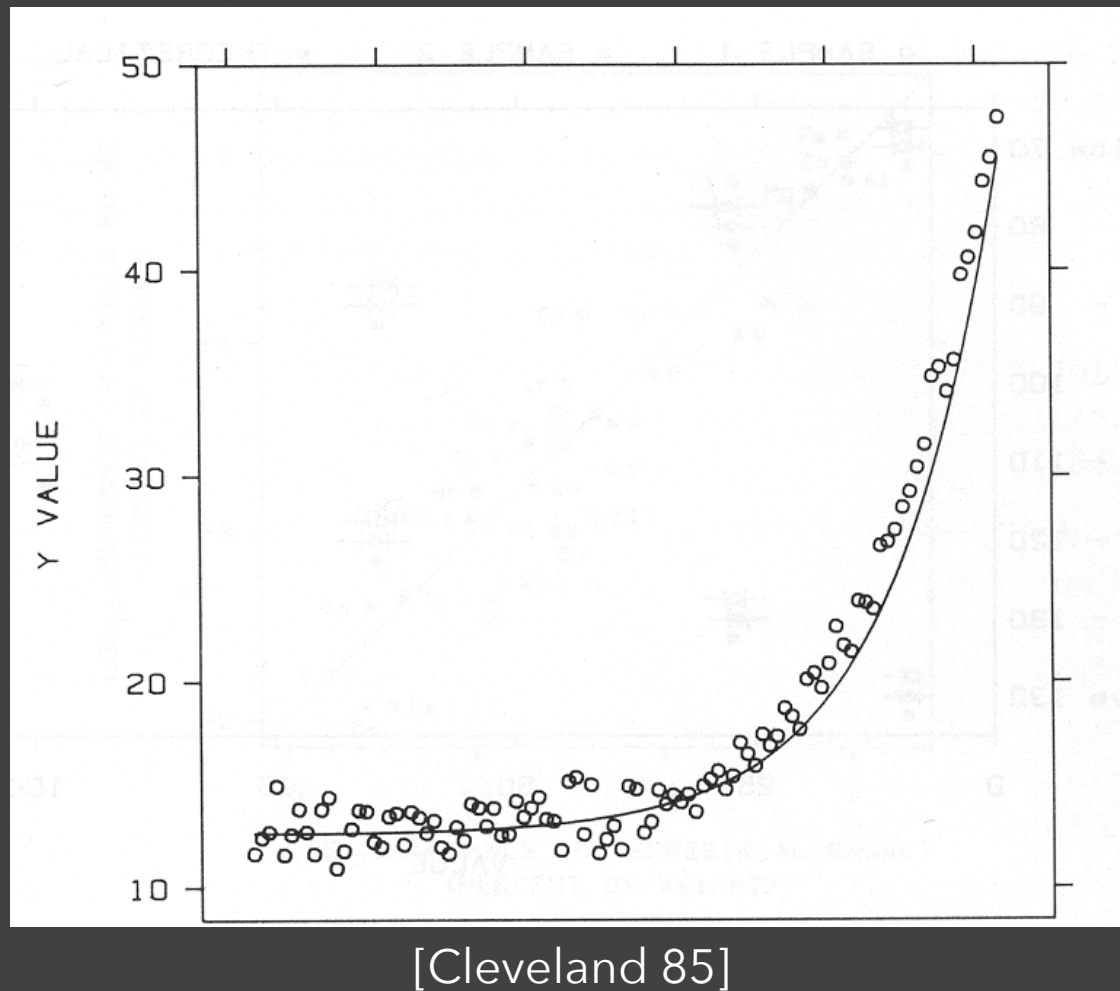
[The Elements of Graphing Data. Cleveland 94]



[The Elements of Graphing Data. Cleveland 94]

Transforming Data

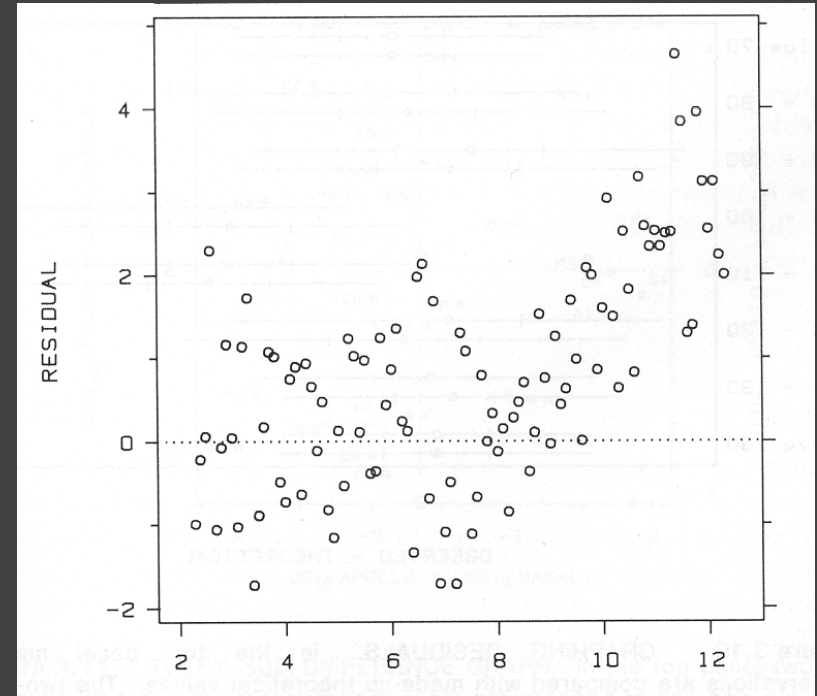
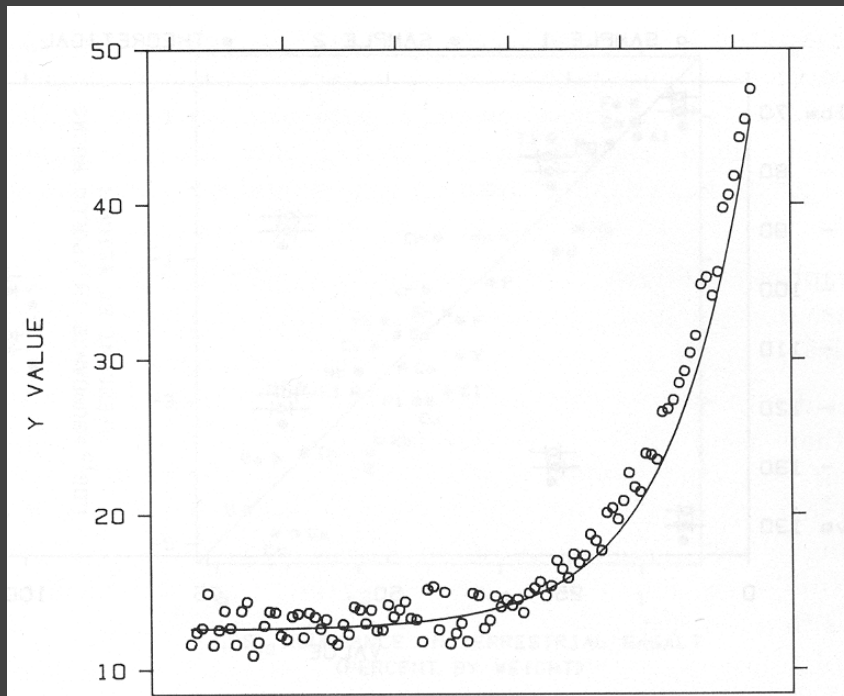
How well does the curve fit the data?



Plot the Residuals

Plot vertical distance from best fit curve

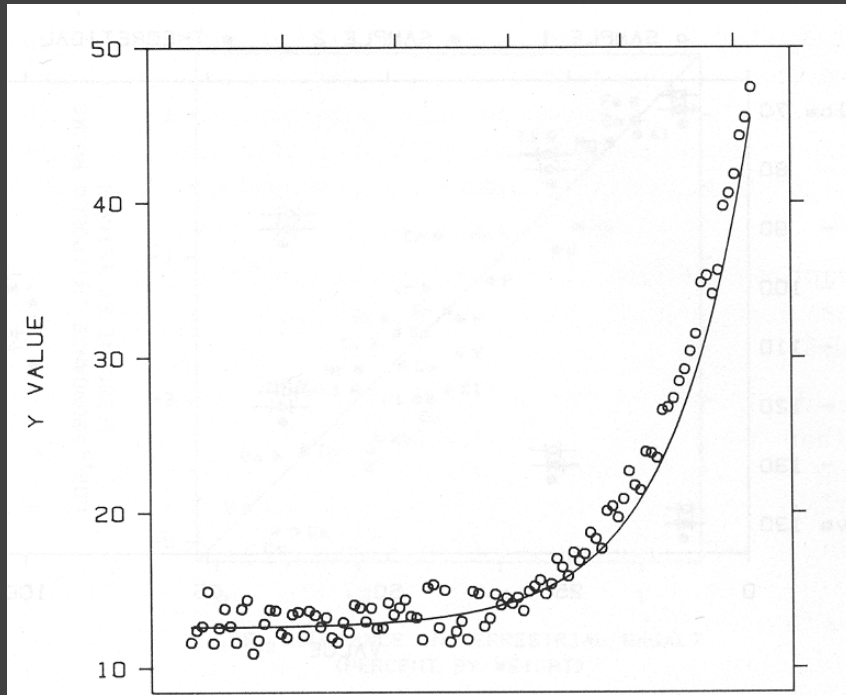
Residual graph shows accuracy of fit



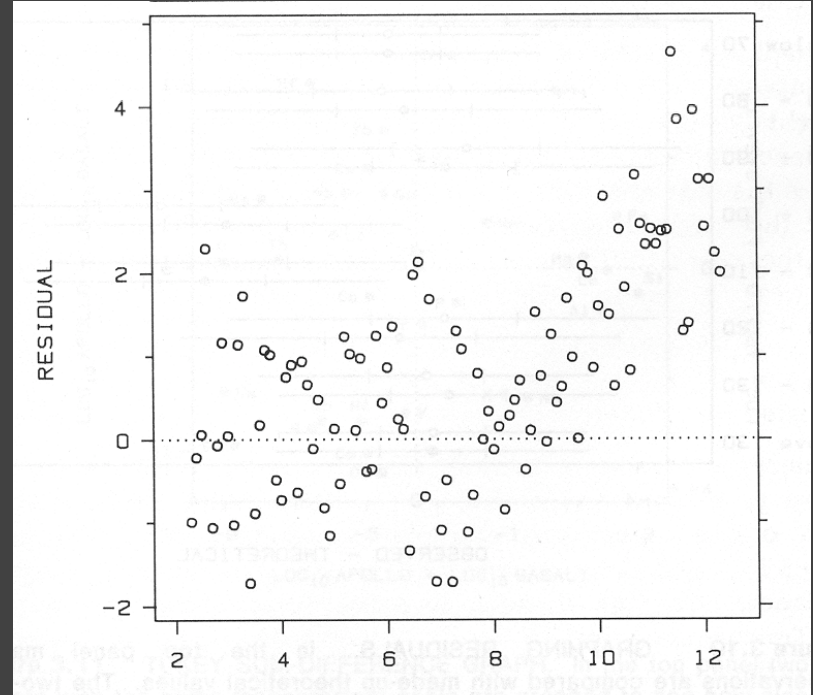
[Cleveland 85]

Multiple Plotting Options

Plot model in data space



Plot data in model space



[Cleveland 85]

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

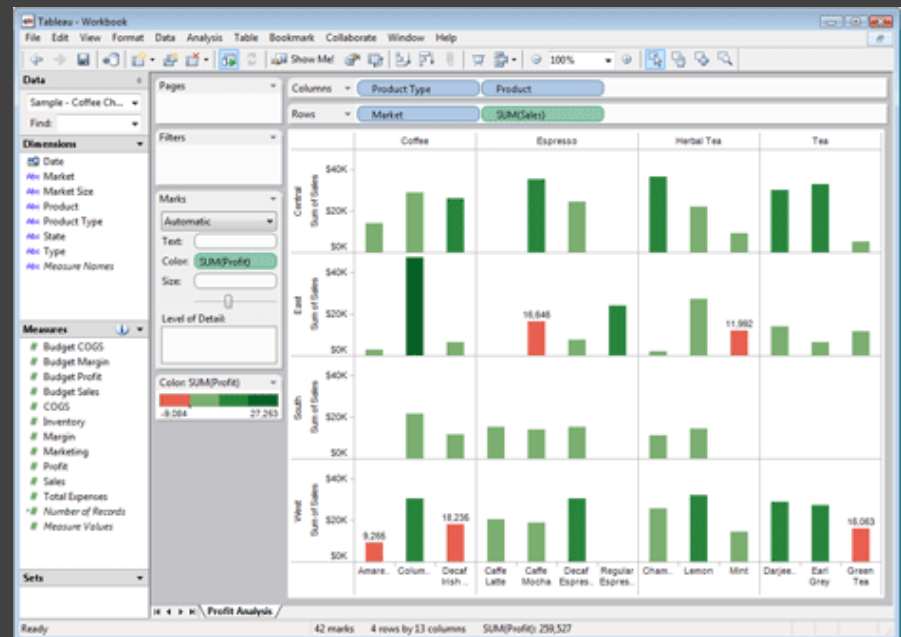
Interact with data

Refine your questions

Make a notebook

Keep record of your analysis

Prepare a final graphic and caption



Due by 5:00pm

Friday, April 15

Tutorials!

Visualization Tools

Tue 4/12, 3:00-4:20pm PAA 114A

Introduction to Tableau, plus a few others.

d3.js: Data-Driven Documents

Tue 4/19, 3:00-4:20pm PAA 114A

Focus on D3, touches on HTML/CSS/JS

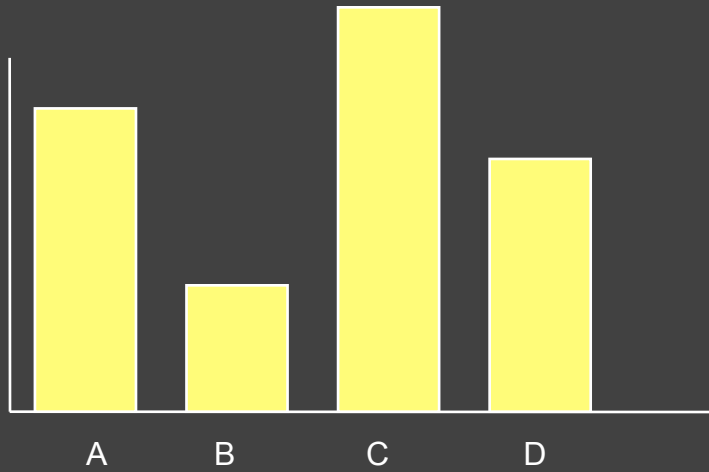
The Design Space of Visual Encodings

Univariate Data

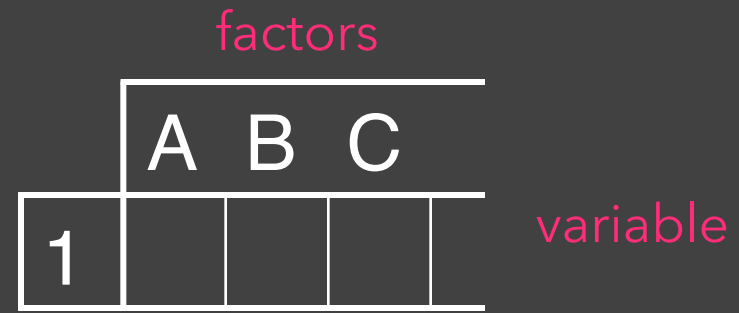
factors

	A	B	C	
1				

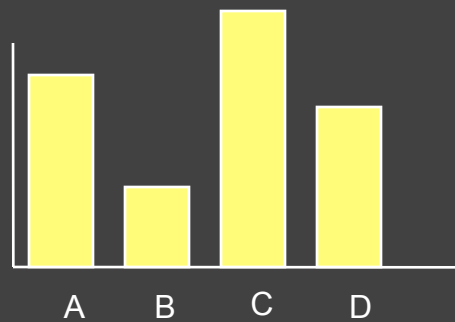
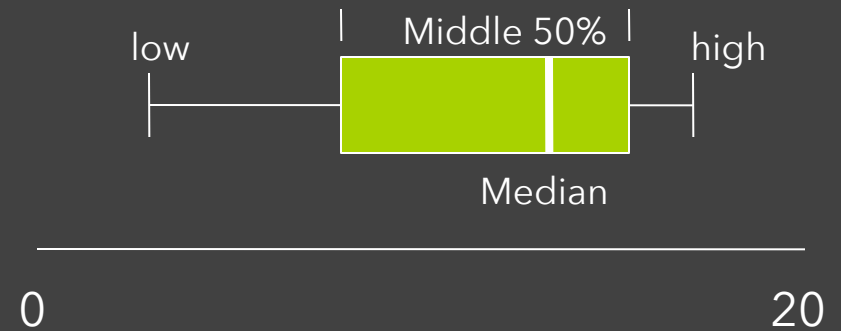
variable



Univariate Data

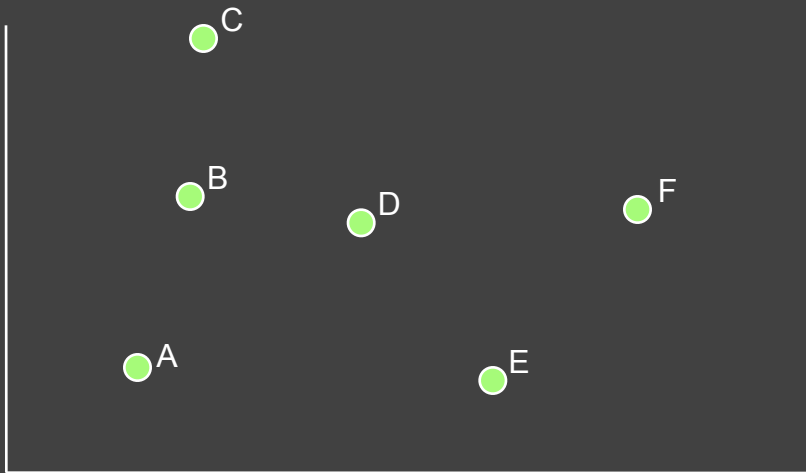


Tukey box plot



Bivariate Data

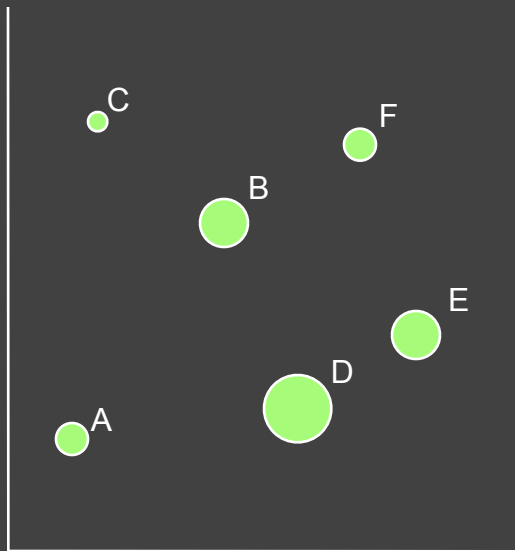
	A	B	C
1			
2			



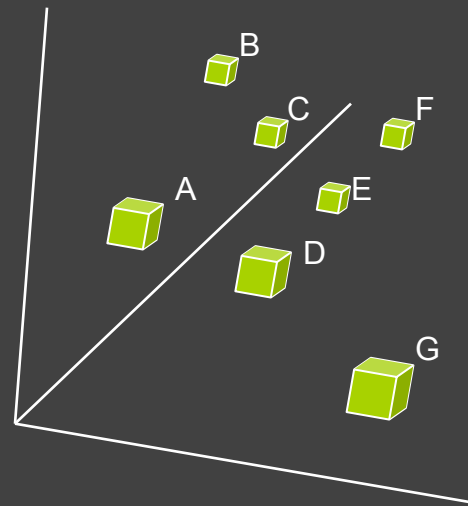
Scatter plot is common

Trivariate Data

	A	B	C	
1				
2				
3				



3D scatter plot is possible



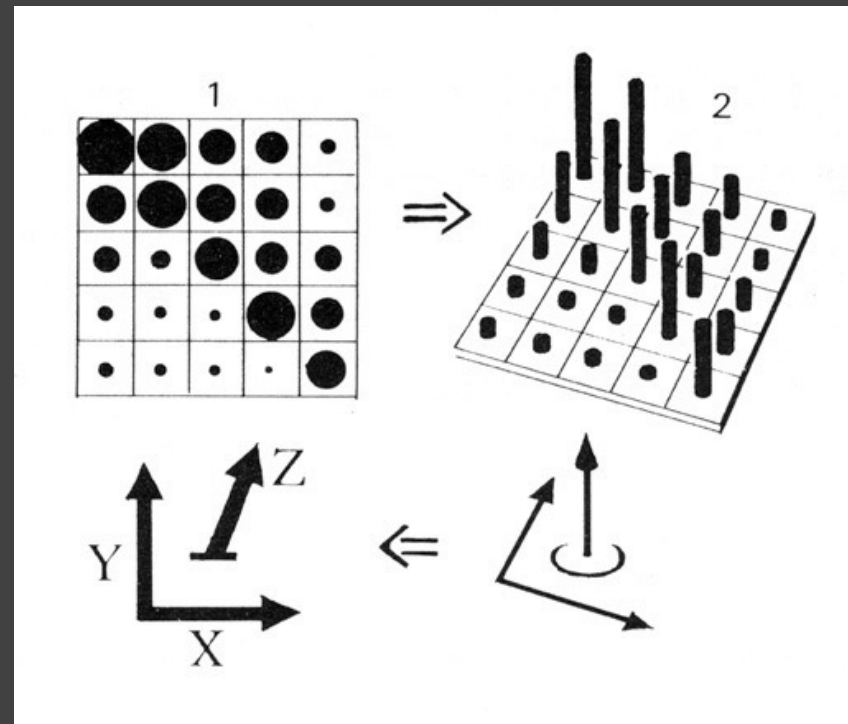
Three Variables

Two variables $[x,y]$ can map to points

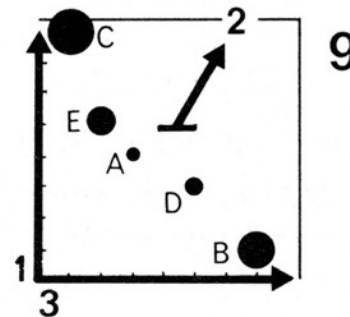
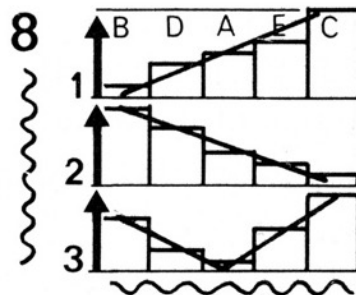
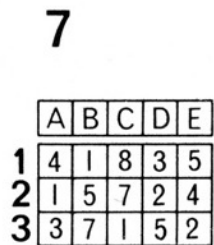
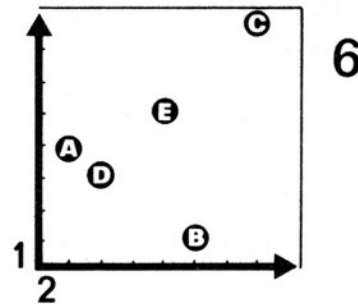
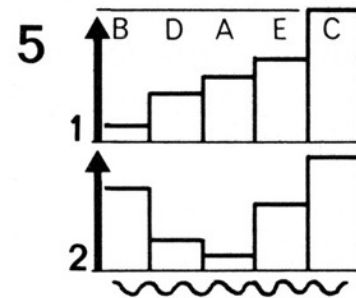
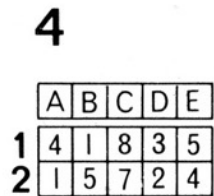
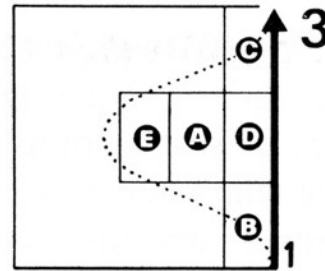
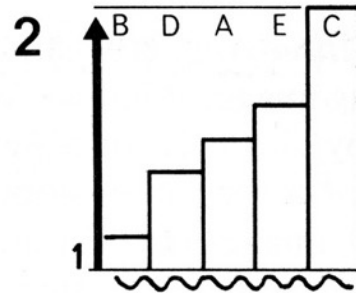
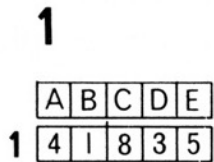
Scatterplots, maps, ...

Third variable $[z]$ must use

Color, size, shape, ...



Large Design Space



[Bertin, Graphics and Graphic Info. Processing, 1981]

Multidimensional Data

Visual Encoding Variables

Position (X)

Position (Y)

Size

Value

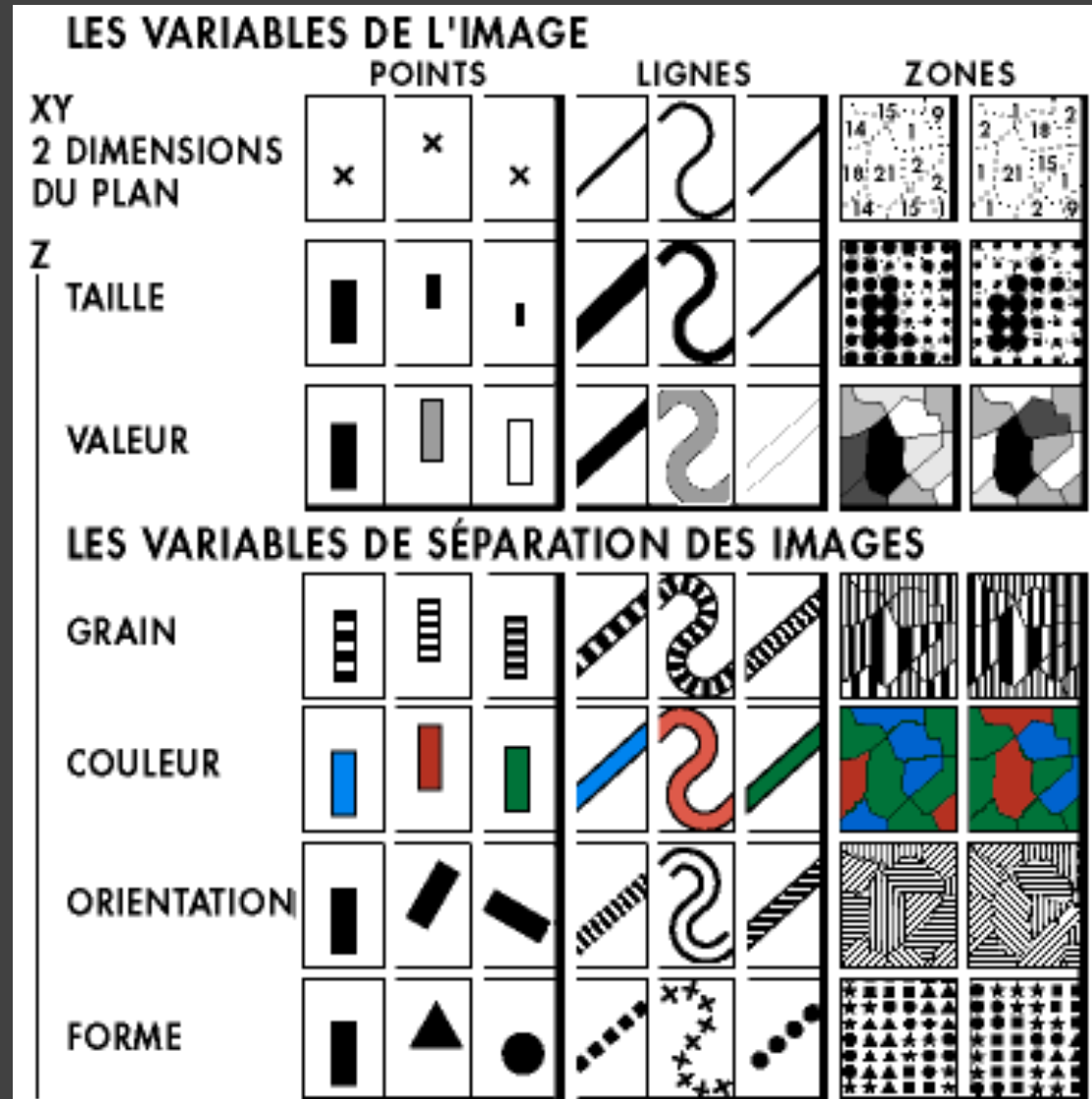
Texture

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

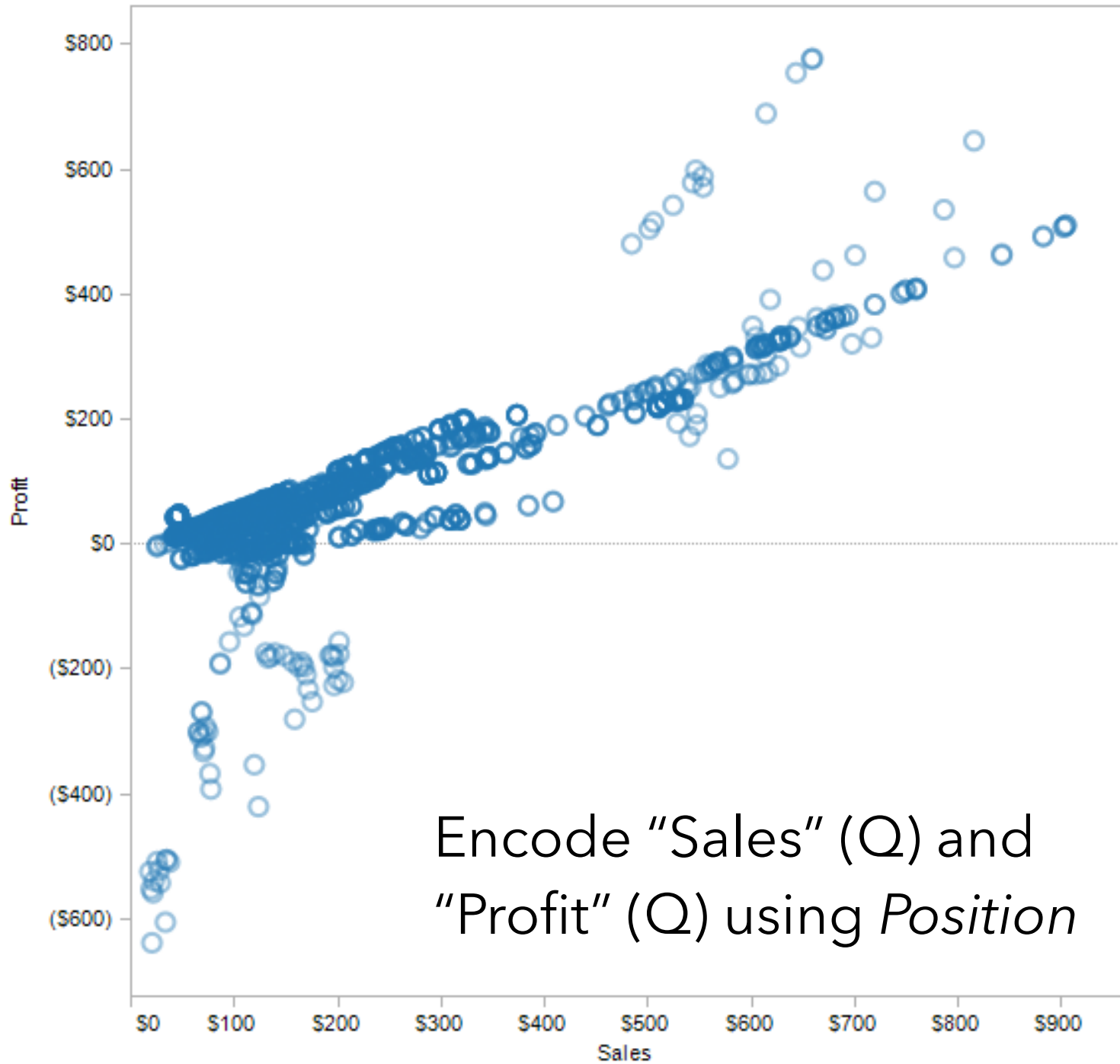
Shape ○

Label ▾

Color ▾

Size

Level of Detail



Filters

YEAR(Date): 2010

Marks

x+ Automatic

Shape

Label

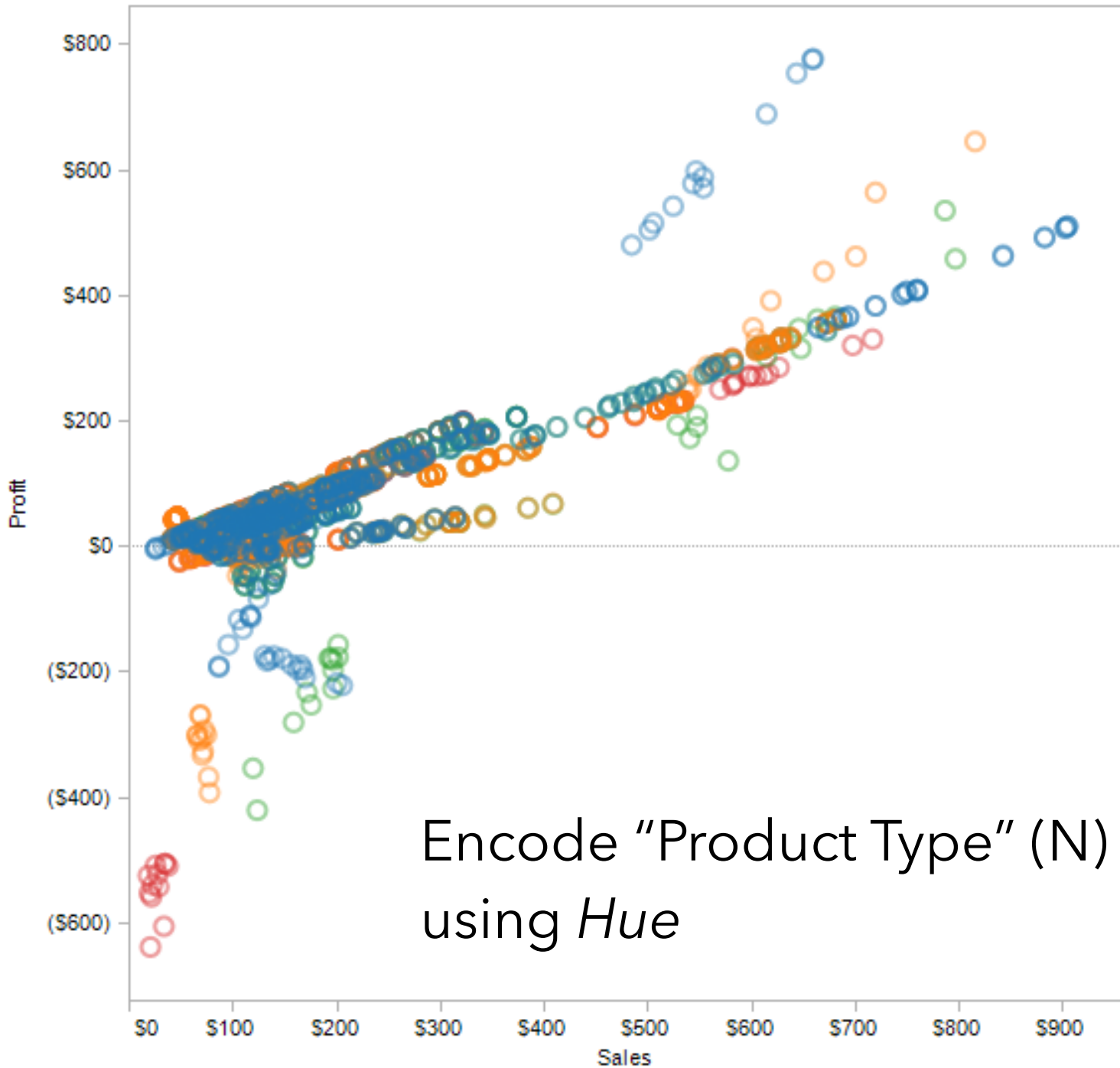
Color

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

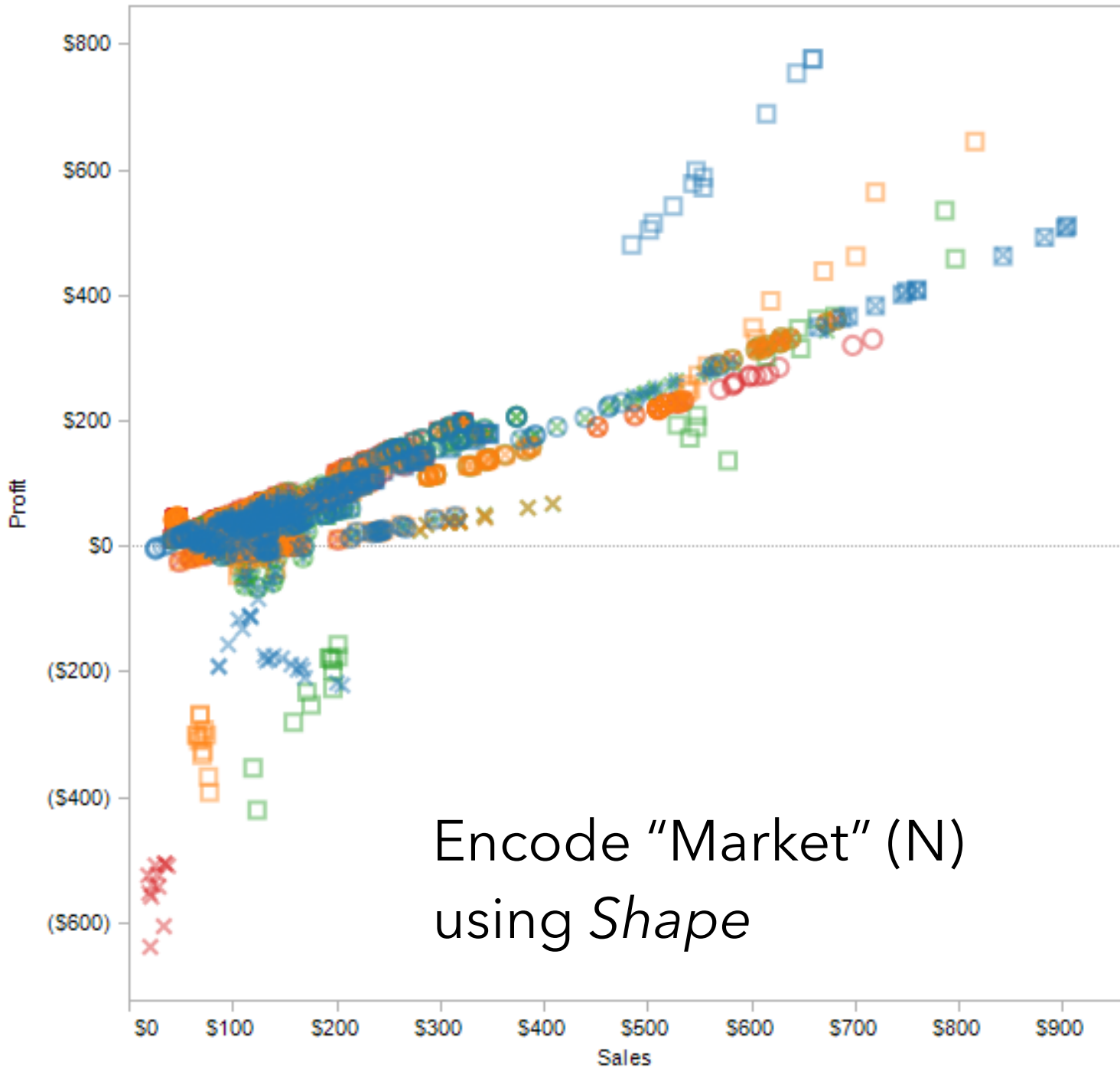
Level of Detail

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

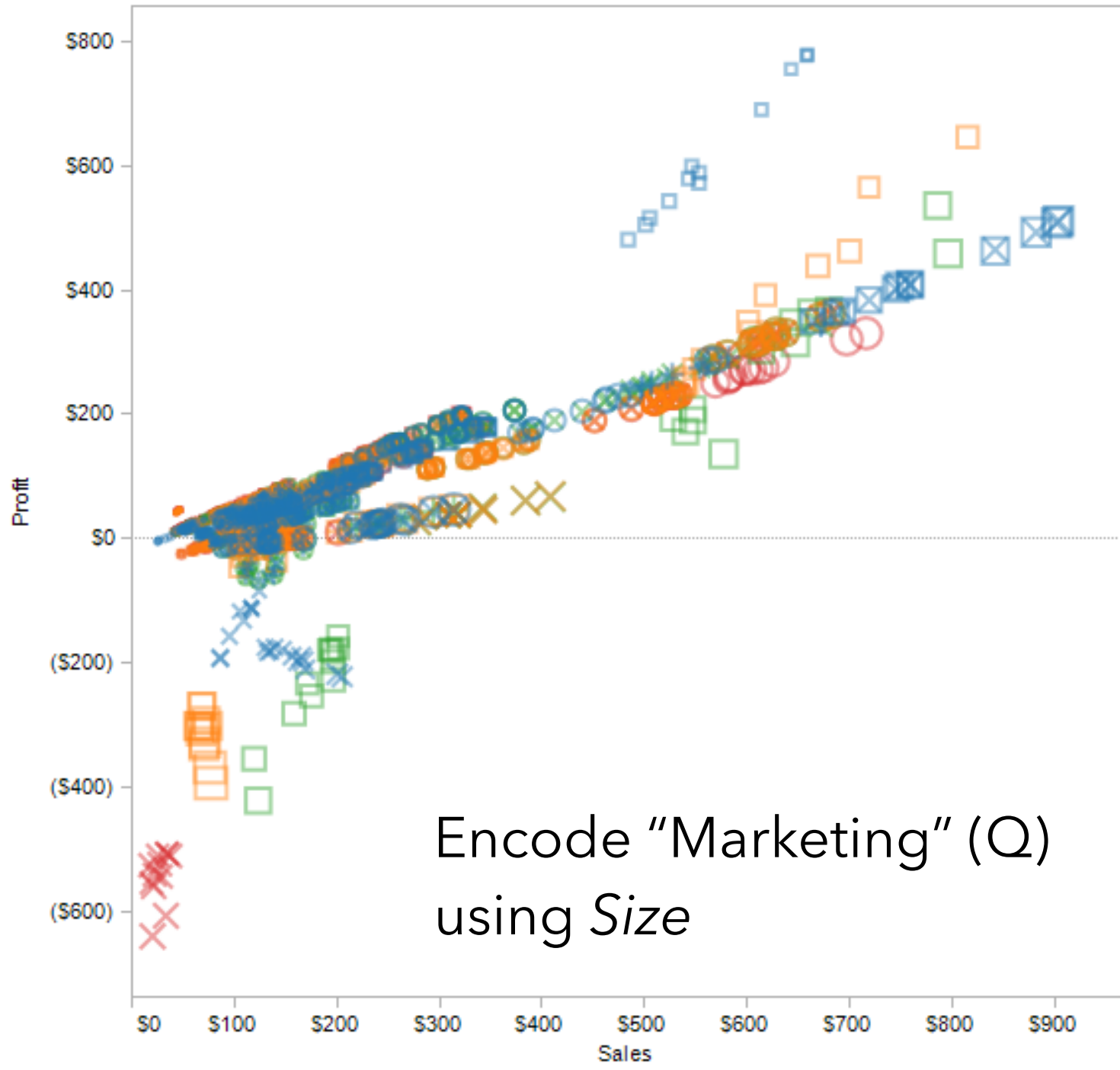
- Coffee
- Espresso
- Herbal Tea

Market

- Central
- East
- South

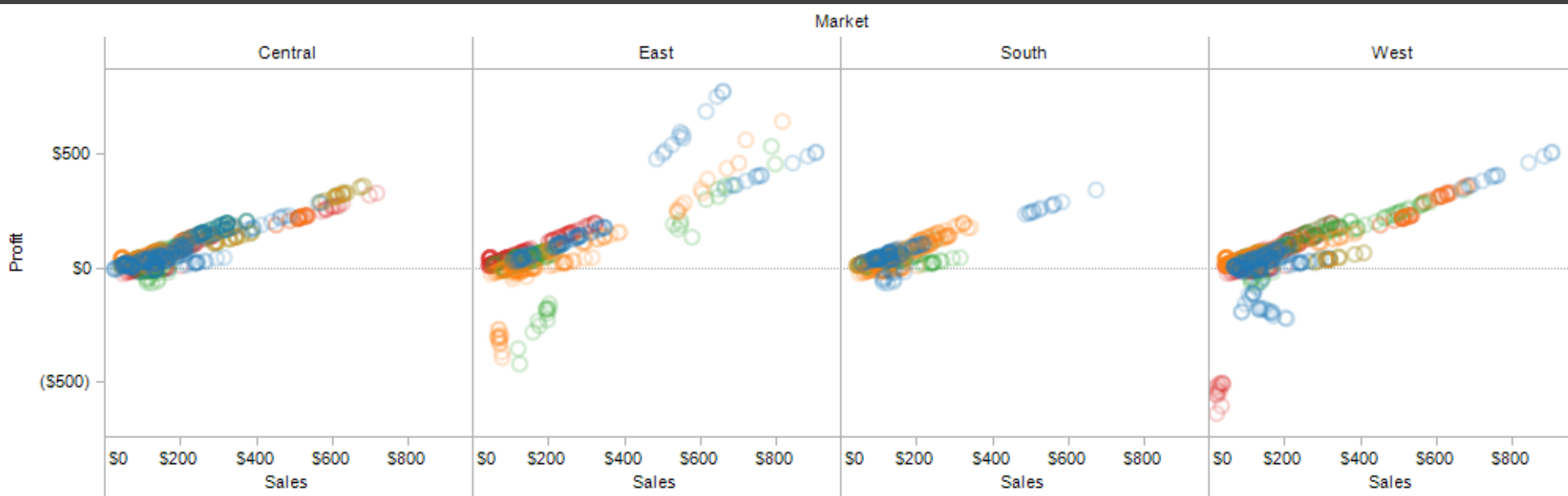
Marketing

- \$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.

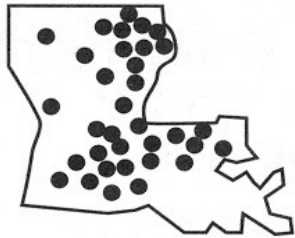
Small Multiples



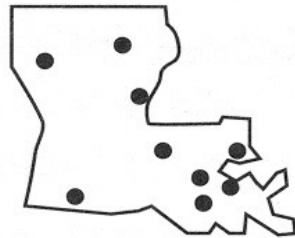
[MacEachren 95, Figure 2.11, p. 38]

Small Multiples

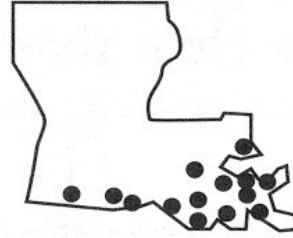
alfisol



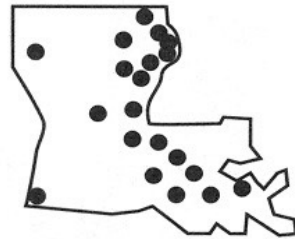
entisol



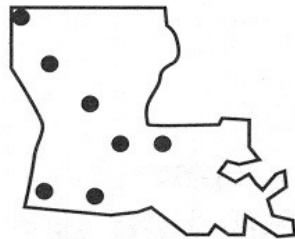
histosol



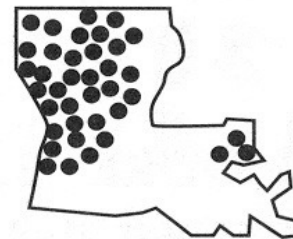
inceptisol



mollisol

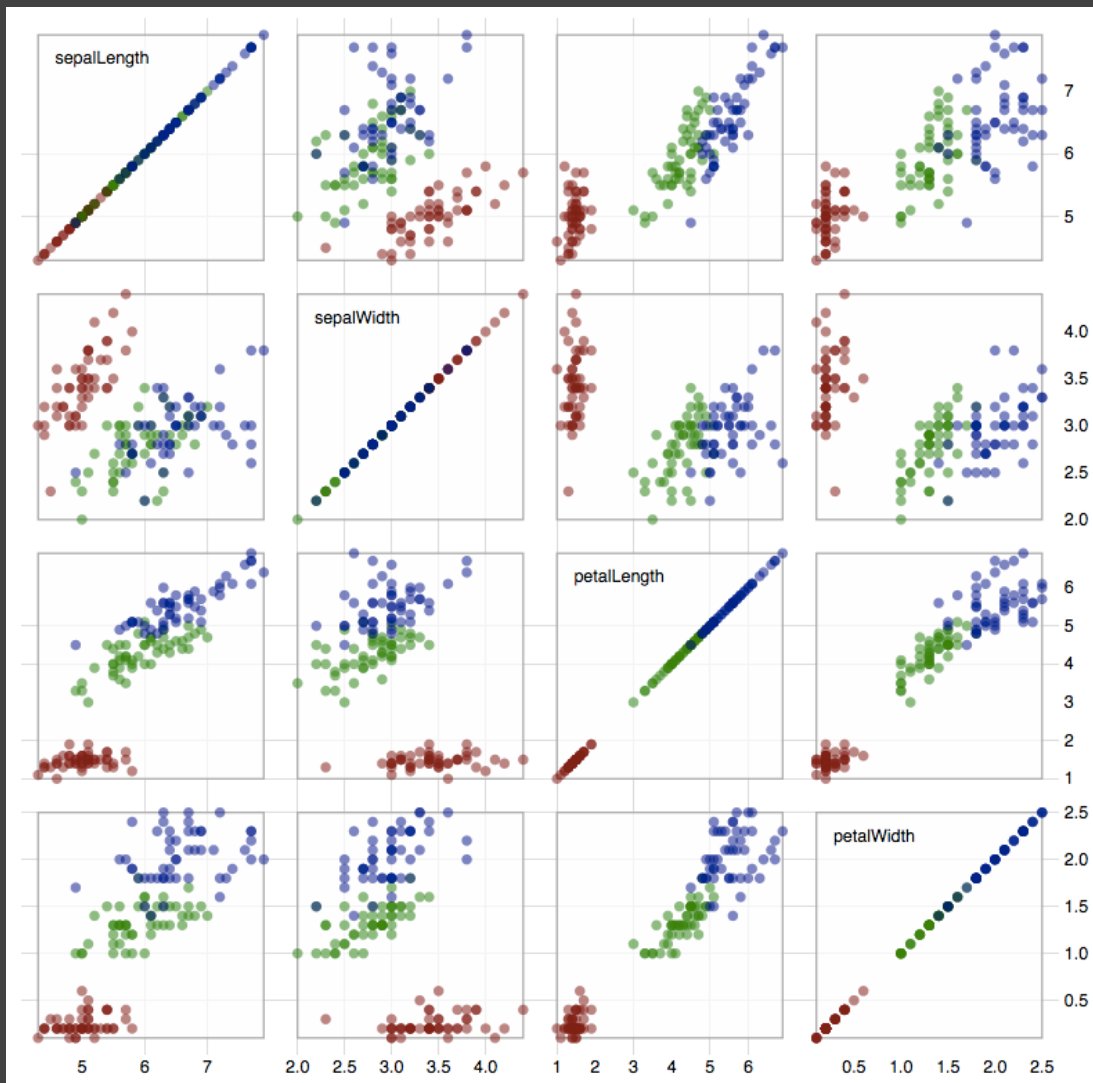


ultisol

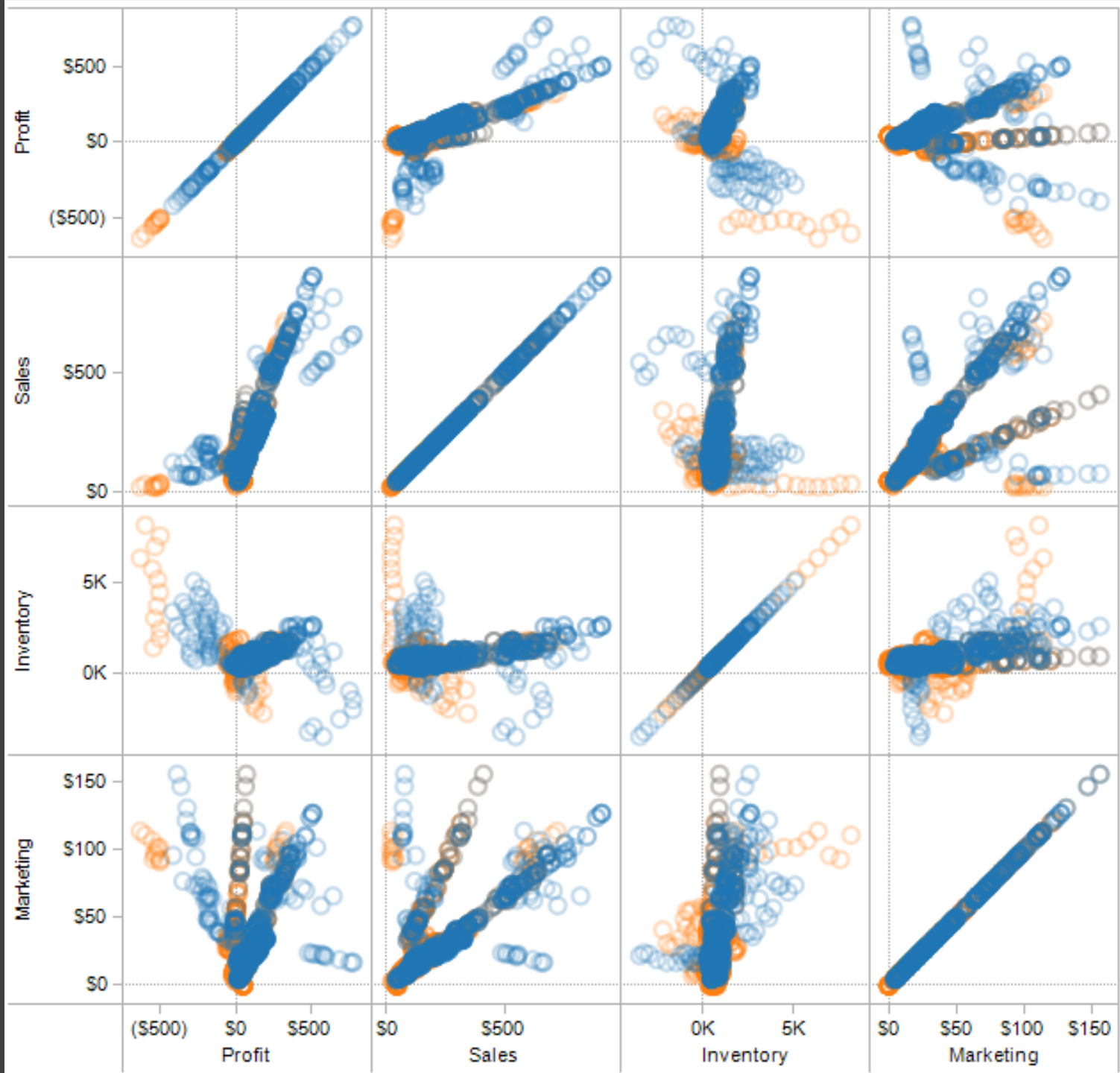


[MacEachren 95, Figure 2.11, p. 38]

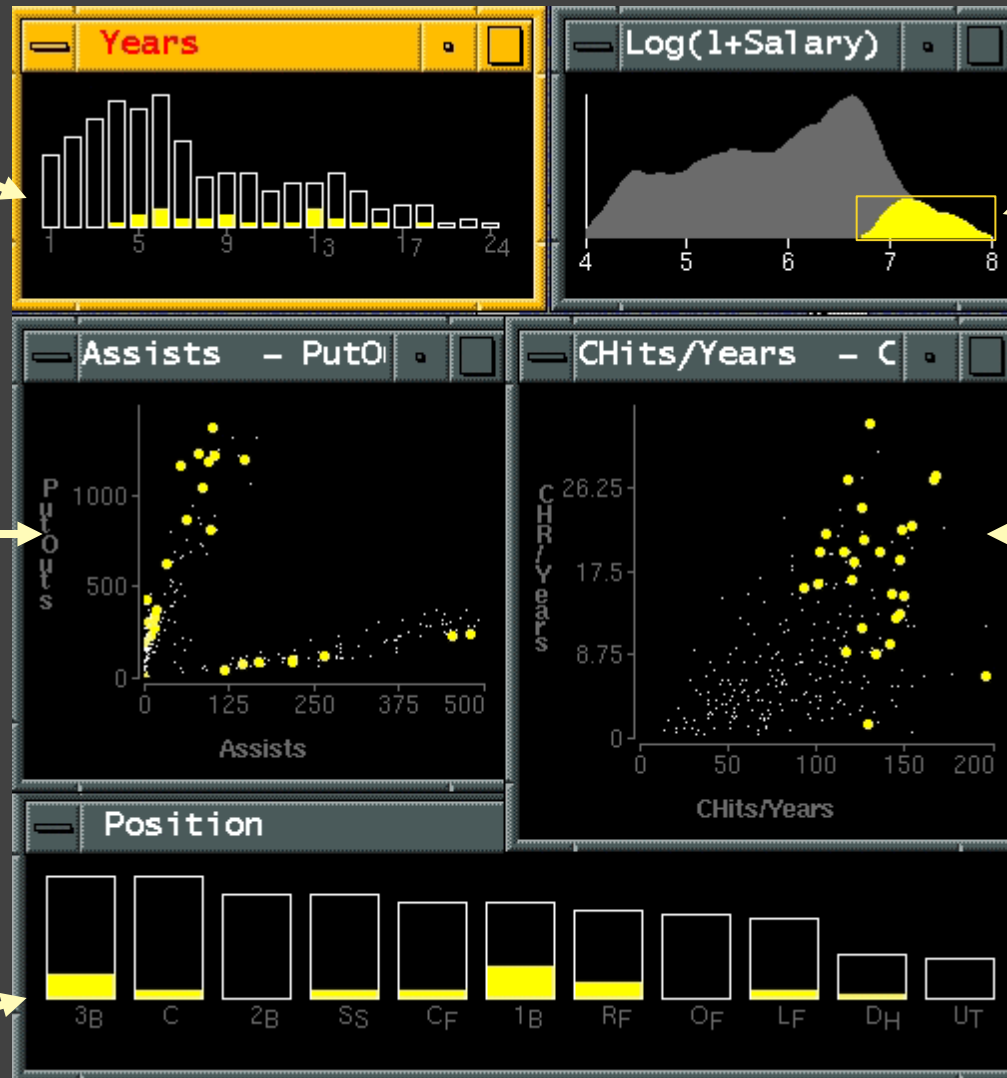
Scatterplot Matrix (SPLOM)



Scatter plots for pairwise comparison of each data dimension.



Multiple Coordinated Views



how long
in majors

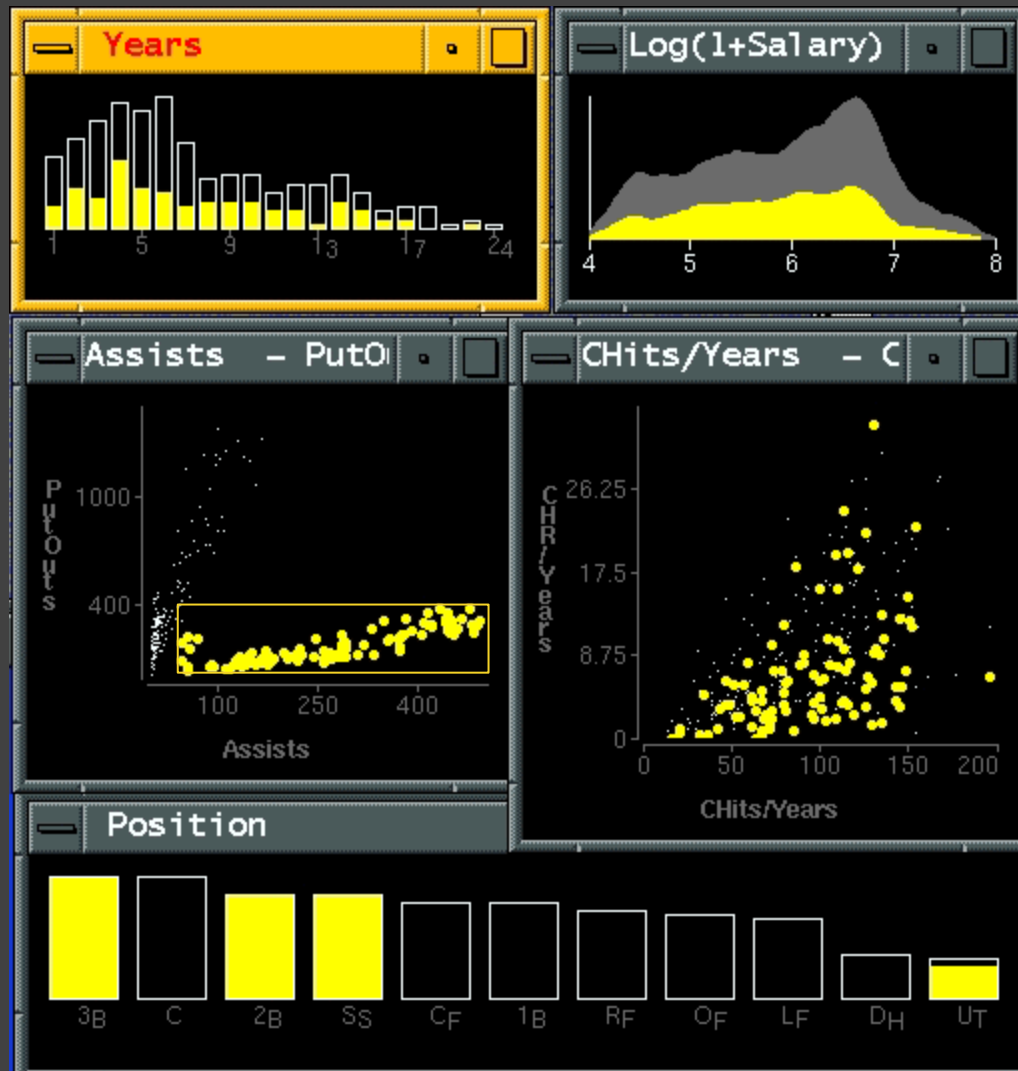
select high
salaries

avg assists vs
avg putouts
(fielding ability)

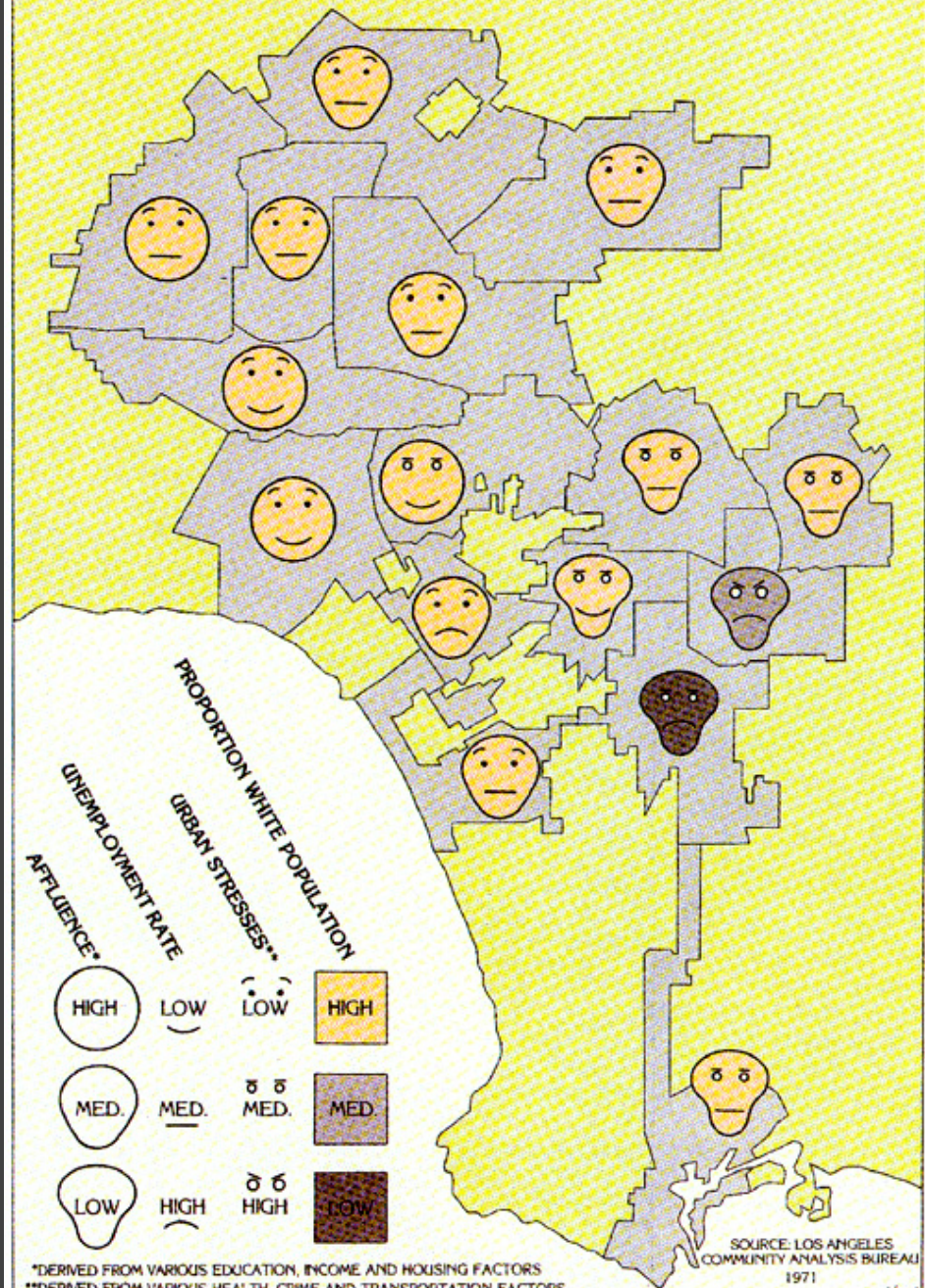
avg career
HRs vs avg
career hits
(batting ability)

distribution
of positions
played

Linking Assists to Position



Life in Los Angeles

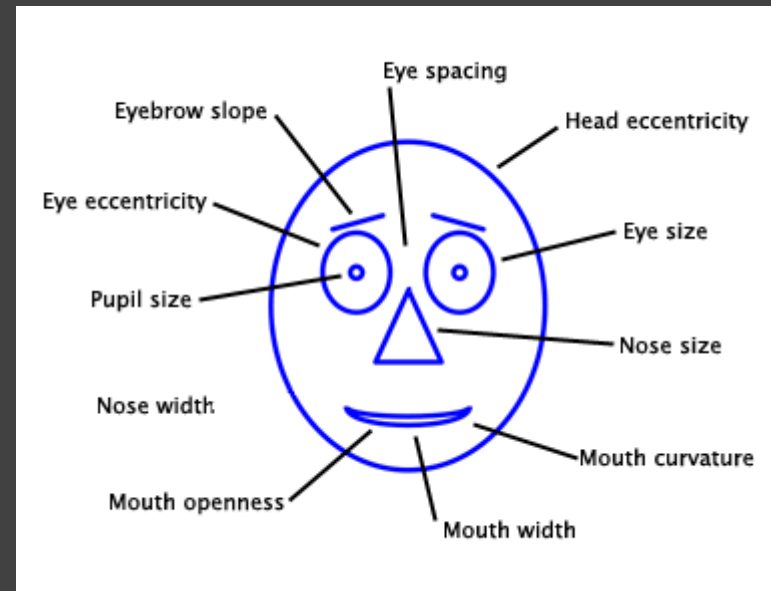


*DERIVED FROM VARIOUS EDUCATION, INCOME AND HOUSING FACTORS
 **DERIVED FROM VARIOUS HEALTH, CRIME AND TRANSPORTATION FACTORS

Chernoff Faces

Observation: We have evolved a sophisticated ability to interpret faces.

Idea: Map data variables to facial features.



Question: Do we process facial features in an uncorrelated way? (i.e., are they *separable*?)

This is just one example of nD "glyphs"

Visualizing Multiple Dimensions

Strategies:

Avoid "over-encoding"

Use space and small multiples intelligently

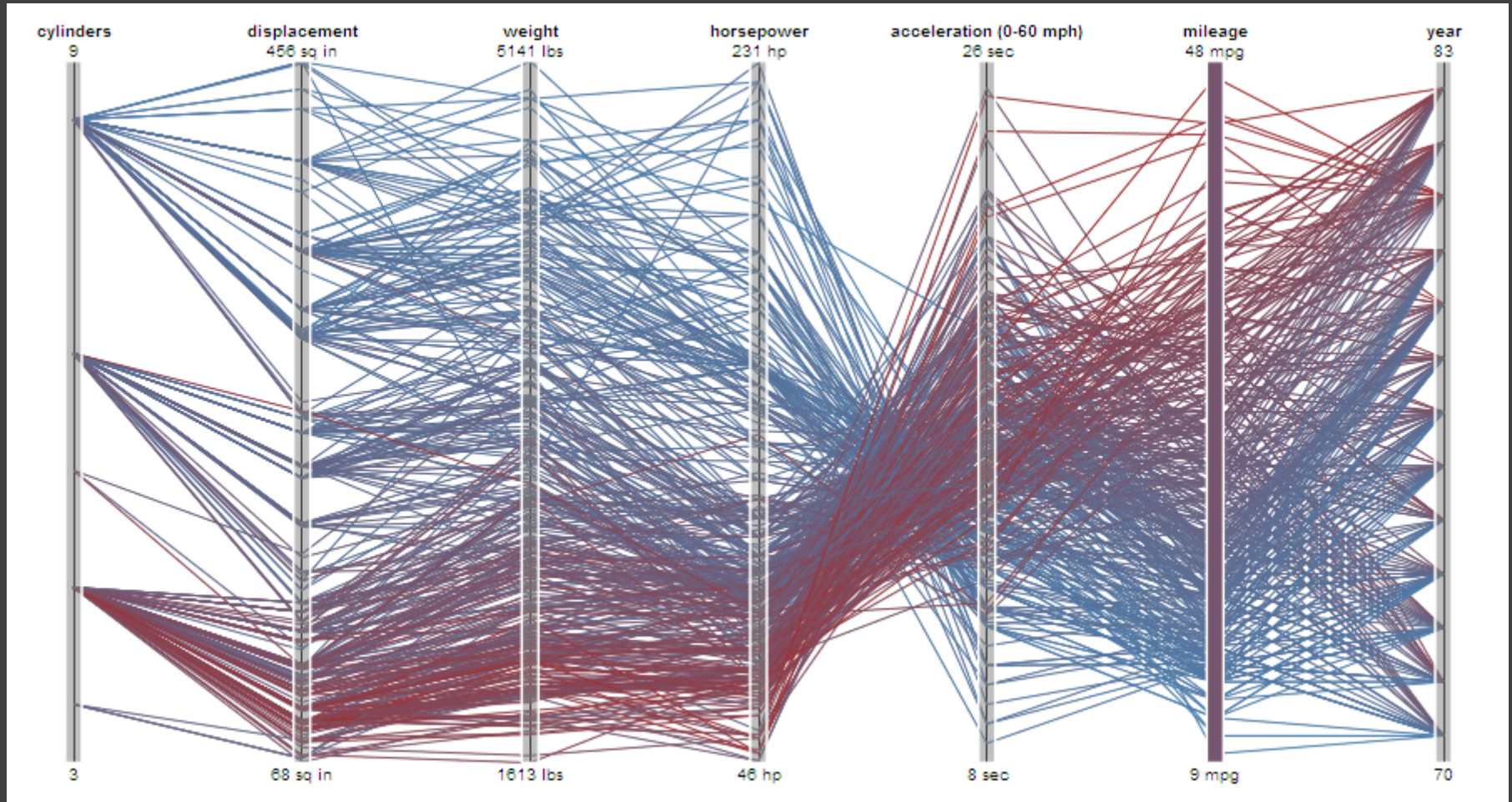
Reduce the problem space

Use interaction to generate *relevant* views

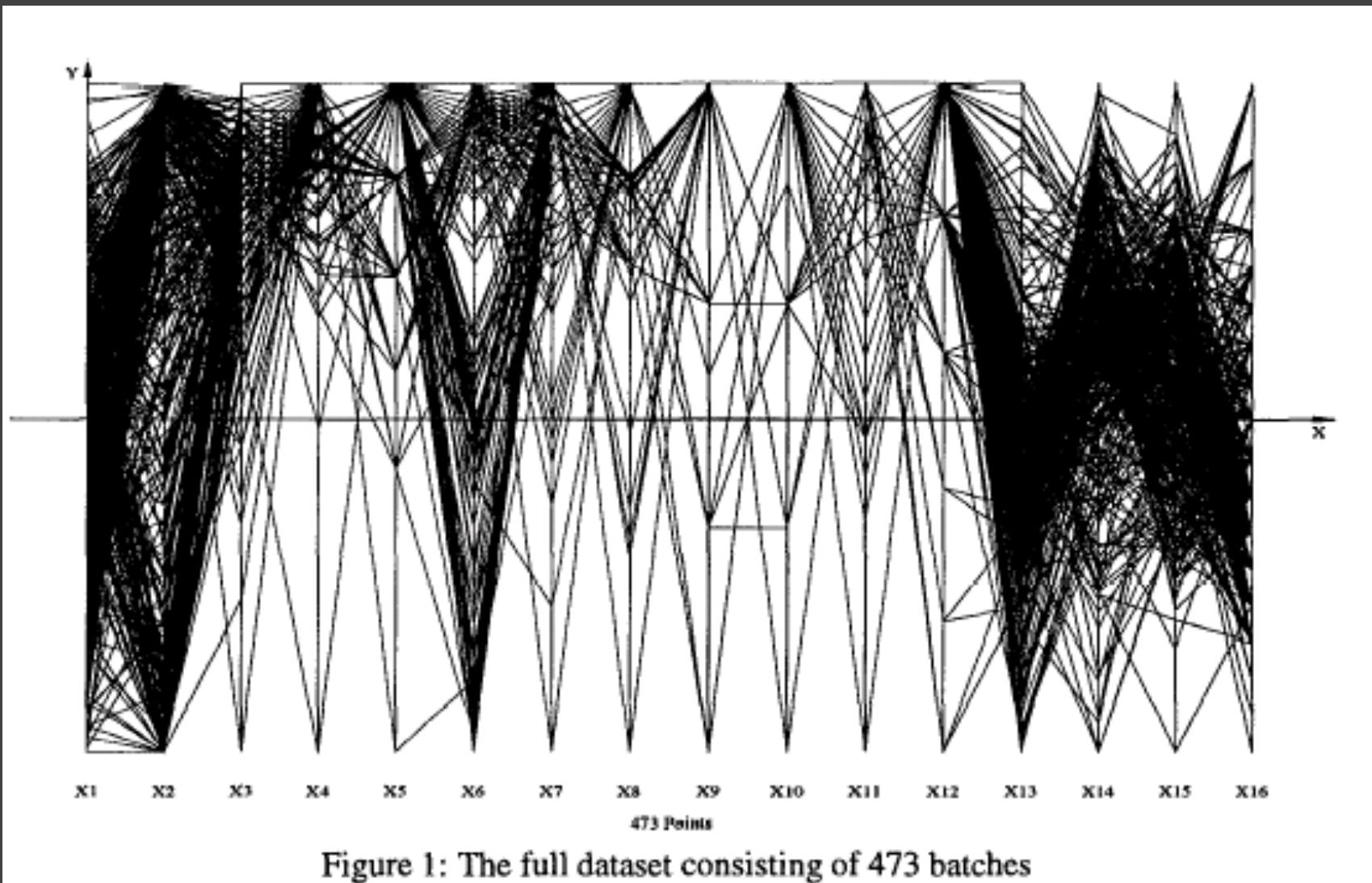
Rarely does a single visualization answer all questions. Instead, the ability to generate appropriate visualizations quickly is key.

Parallel Coordinates

Parallel Coordinates [Inselberg]



Parallel Coordinates [Inselberg]



The Multidimensional Detective

Production data for 473 batches of a VLSI chip

16 process parameters

$X1$: The yield: % of produced chips that are useful

$X2$: The quality of the produced chips (speed)

$X3-12$: 10 types of defects (0 defects shown at top)

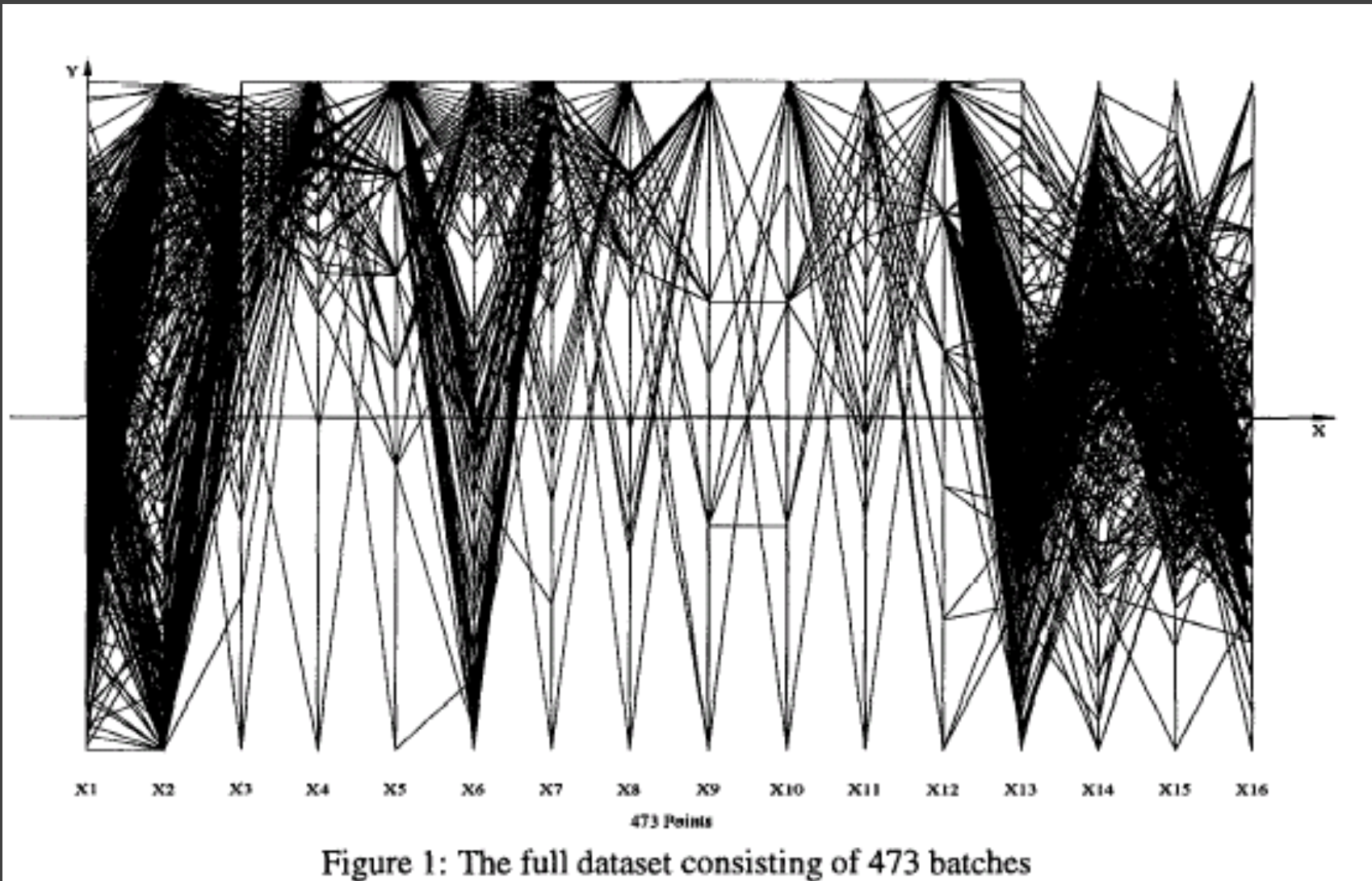
$X13-16$: 4 physical parameters

Objective:

Raise the yield ($X1$) and maintain high quality ($X2$)

A. Inselberg, Multidimensional Detective, Proc. IEEE InfoVis, 1997

Parallel Coordinates [Inselberg]



Inselberg's Principles

1. Do not let the picture scare you.
2. Understand your objectives. Use them to obtain visual cues.
3. Carefully scrutinize the picture.
4. Test your assumptions, especially the "I am really sure of's".
5. You can't be unlucky all the time!

Each line represents a tuple (e.g., VLSI batch)
Filtered below for high values of $X1$ and $X2$

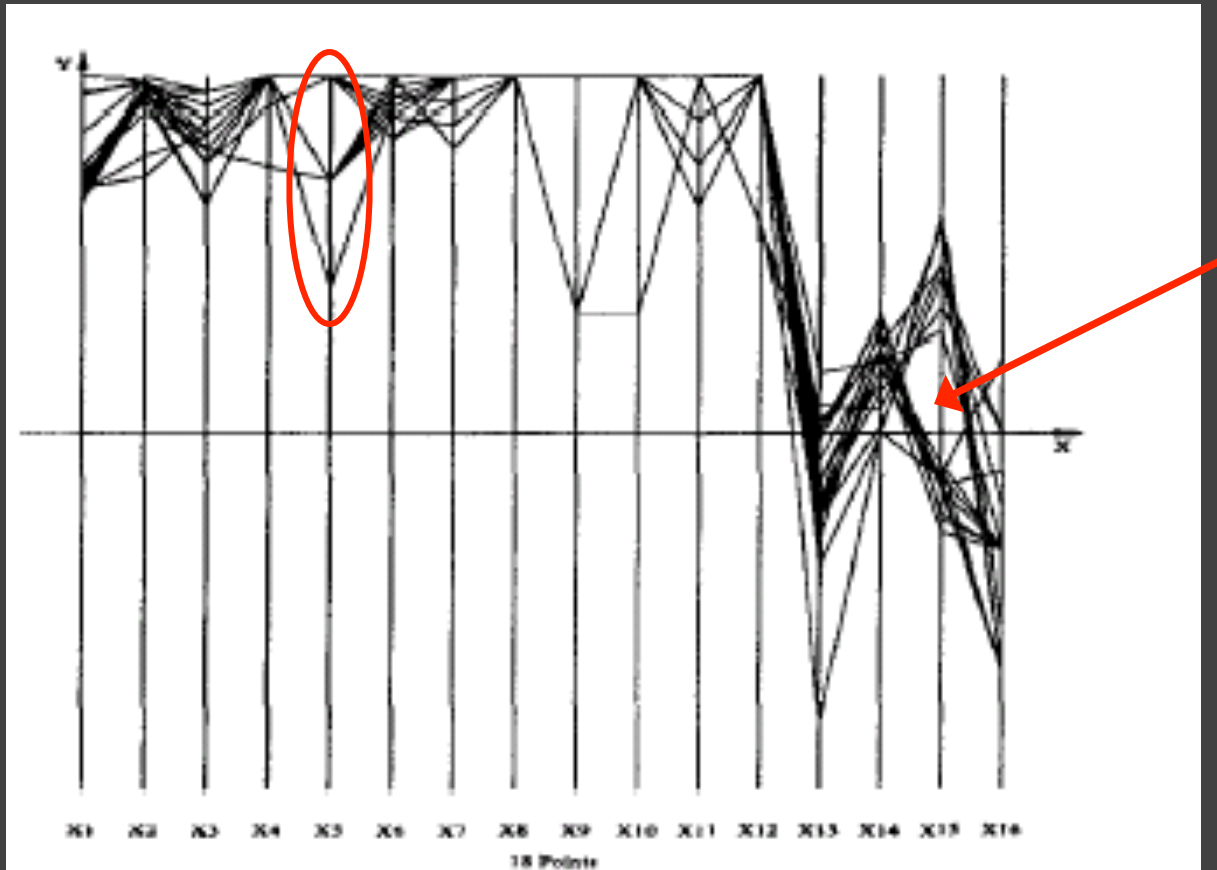
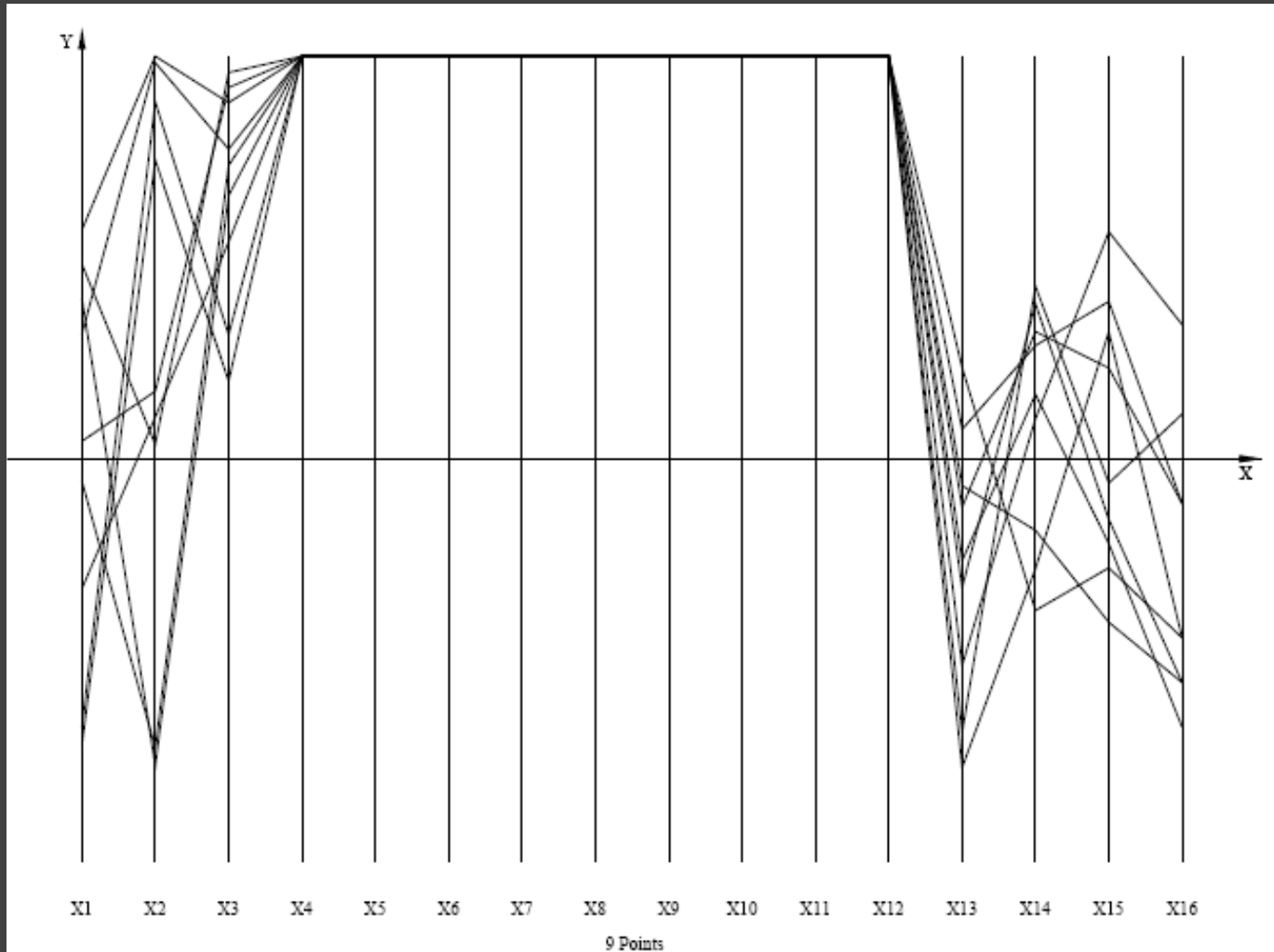


Figure 2: The batches high in Yield, $X1$, and Quality, $X2$.

Look for batches with *nearly* zero defects (9/10)
Most of these have low yields -> defects OK.



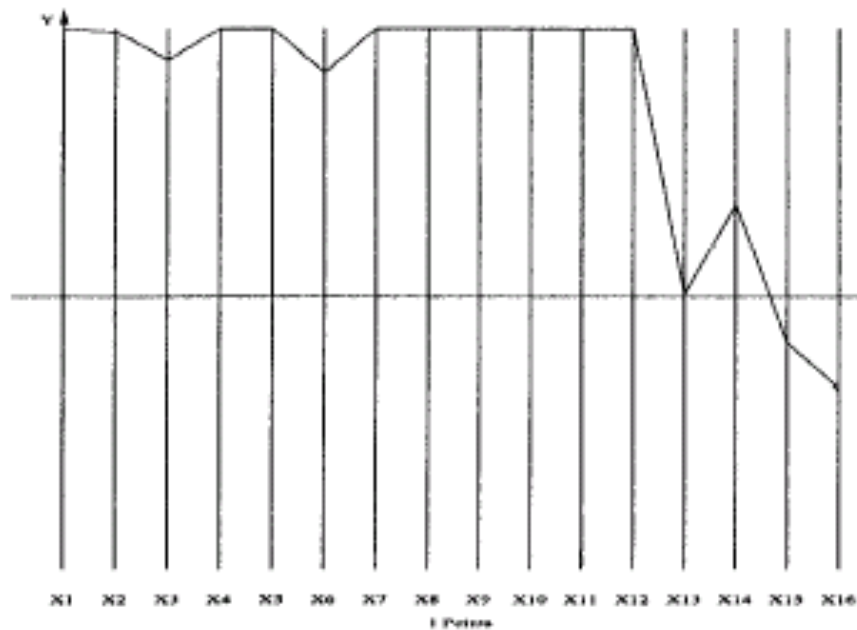


Figure 5: The best batch. Highest in Yield, X_1 , and very high in Quality, X_2 .

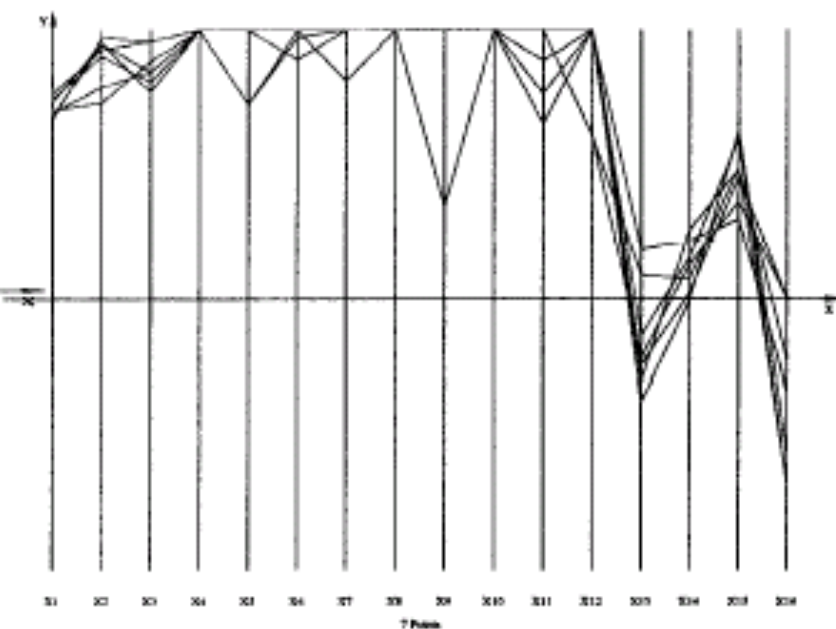
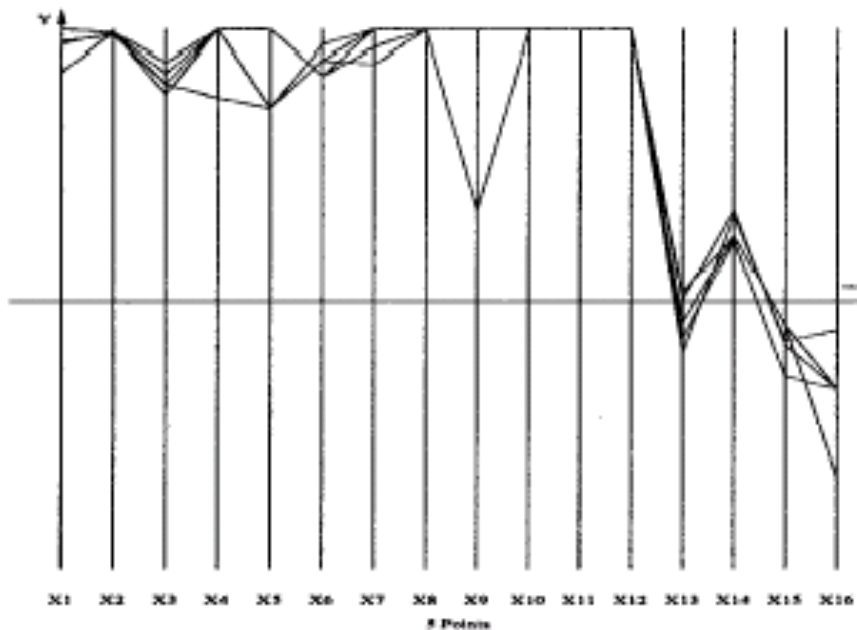
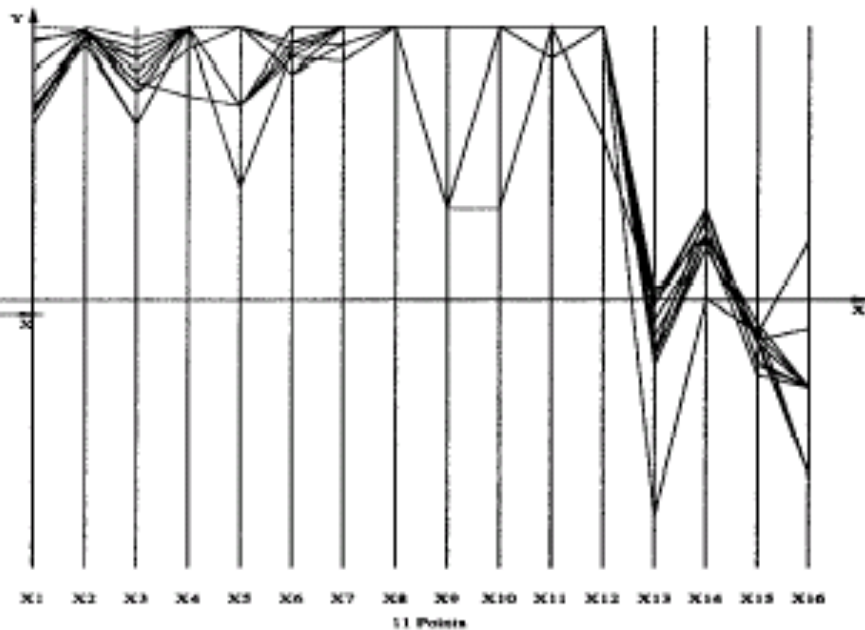


Figure 7: Upper range of split in X_{15}



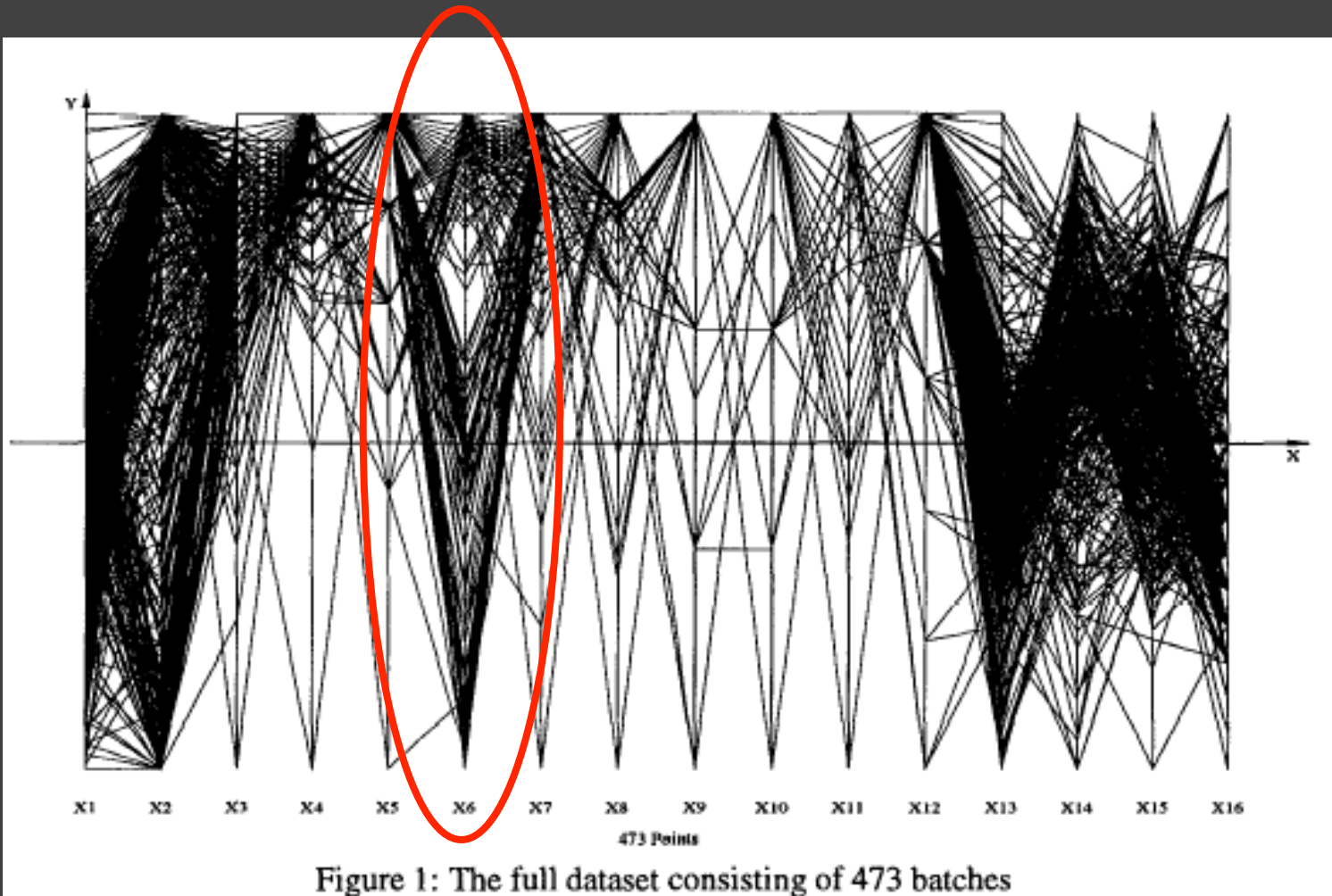
8 Points



11 Points

Notice that $X6$ behaves differently.

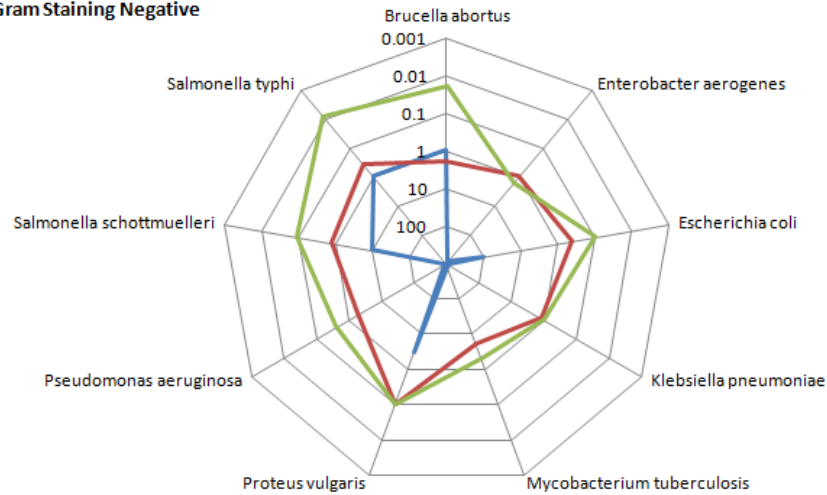
Allow 2 defects, including $X6$ -> best batches



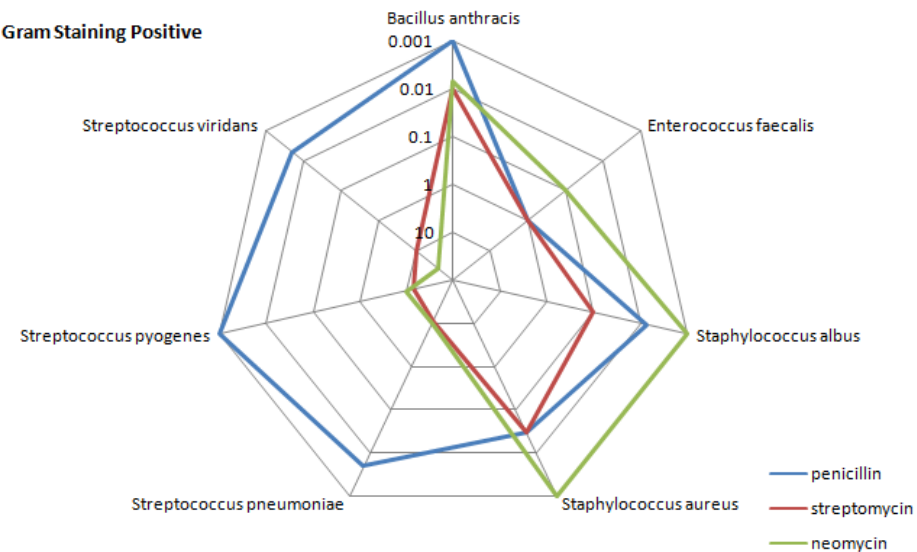
Radar Plot / Star Graph

Antibiotics MIC Concentrations

Gram Staining Negative



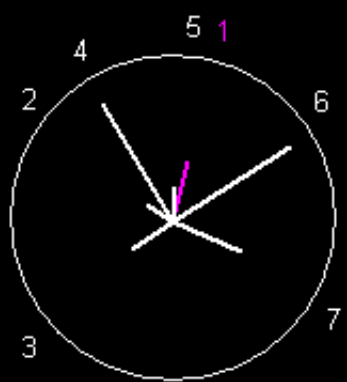
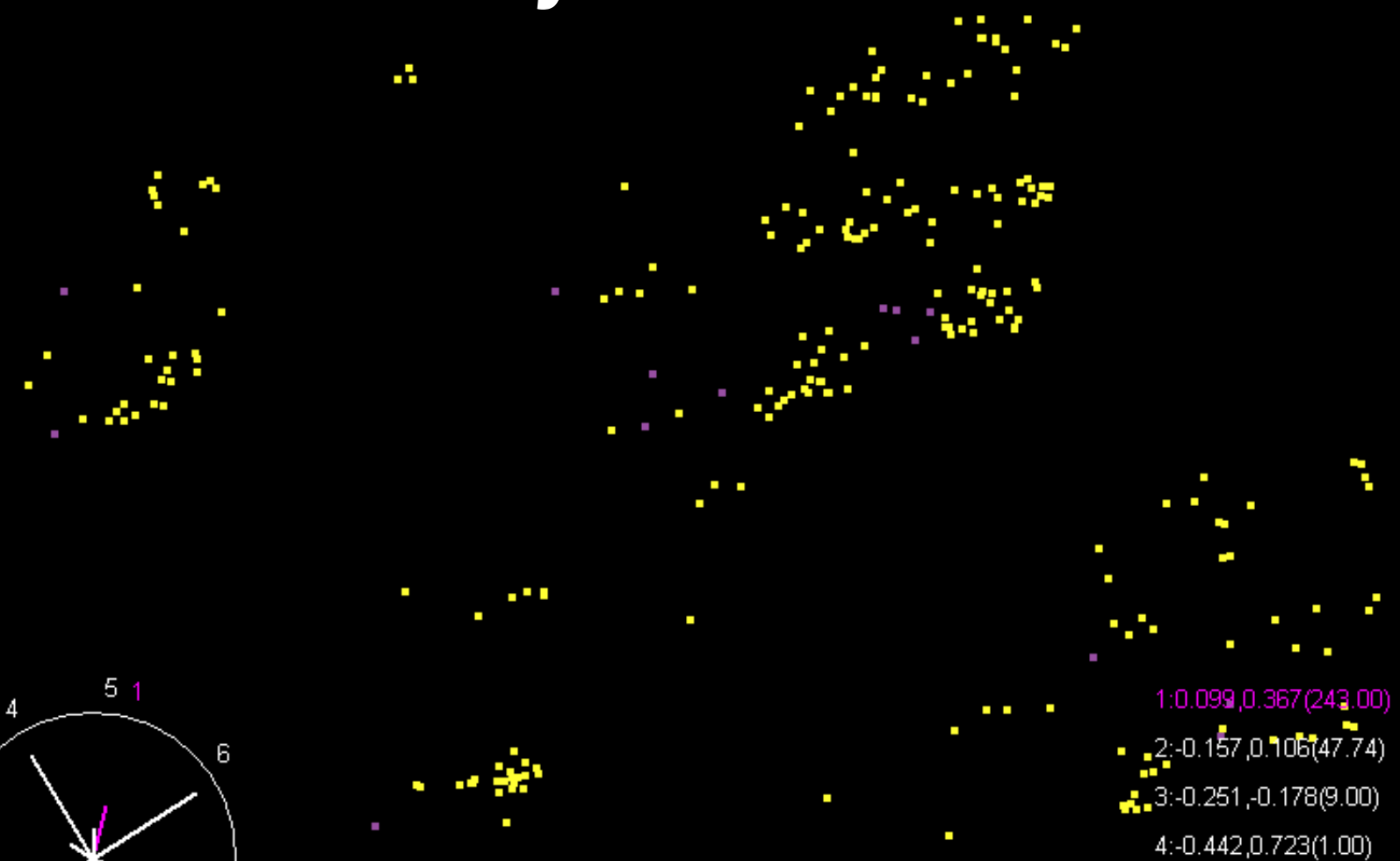
Gram Staining Positive



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

Dimensionality Reduction

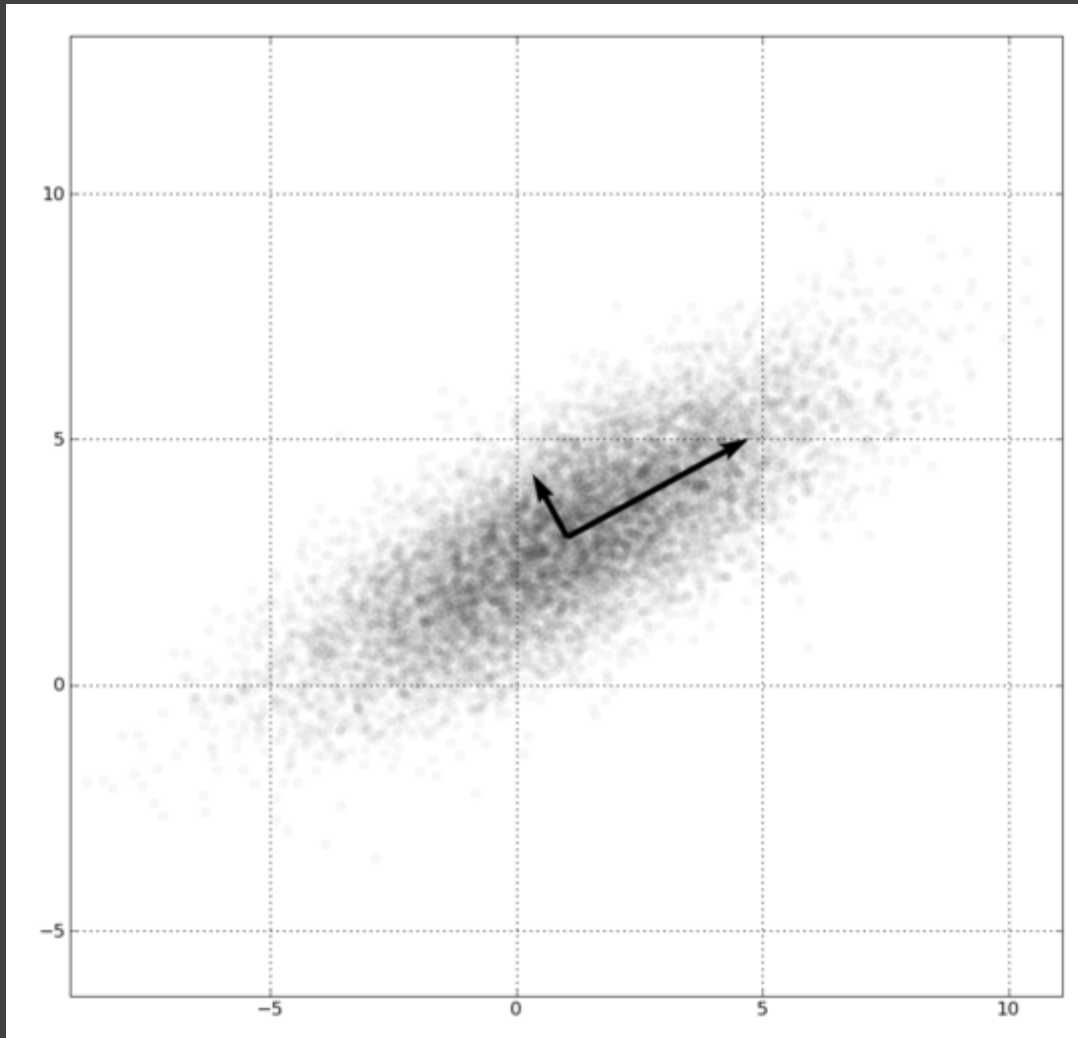
Dimensionality Reduction



- 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)

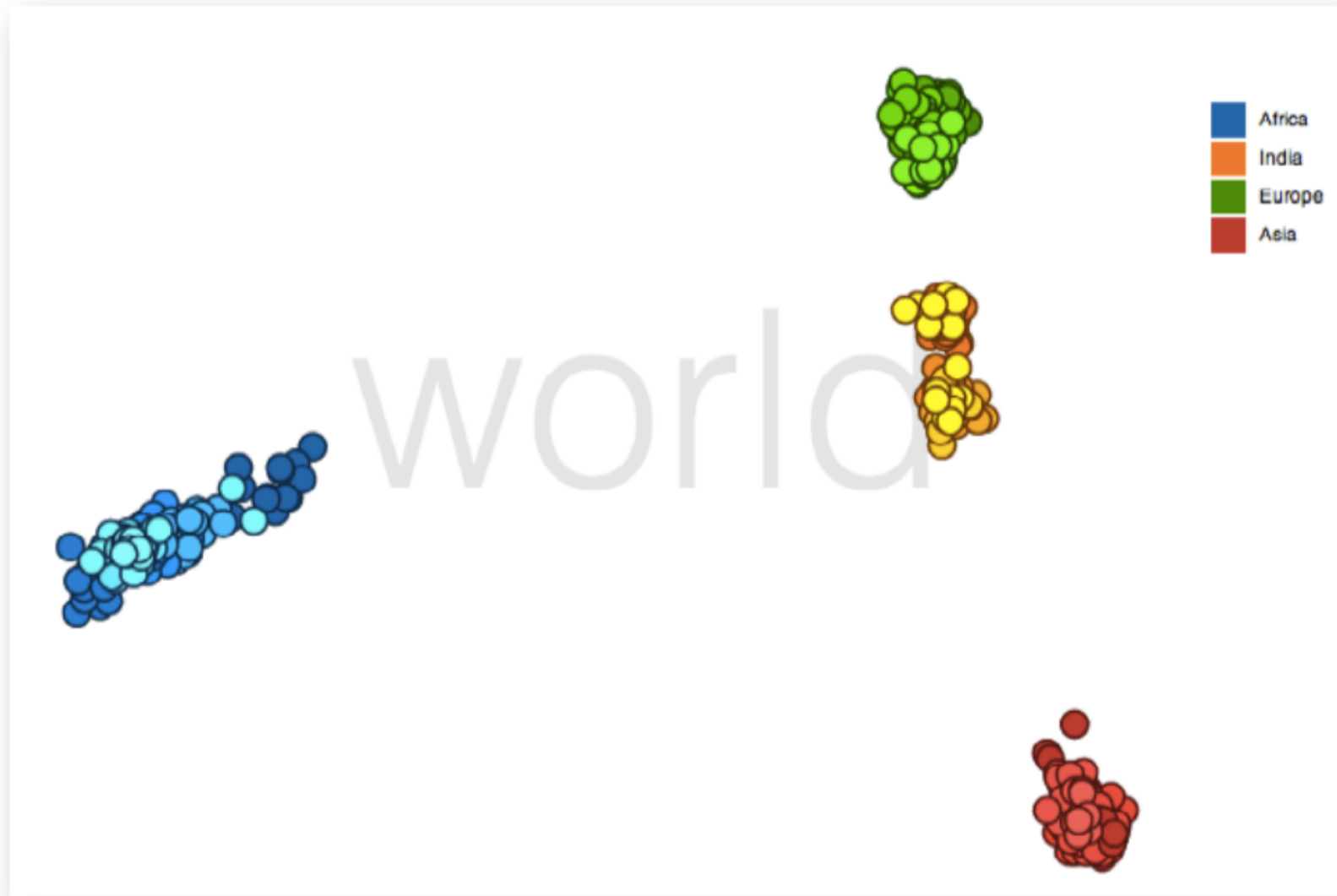
<http://www.ggobi.org/>

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]



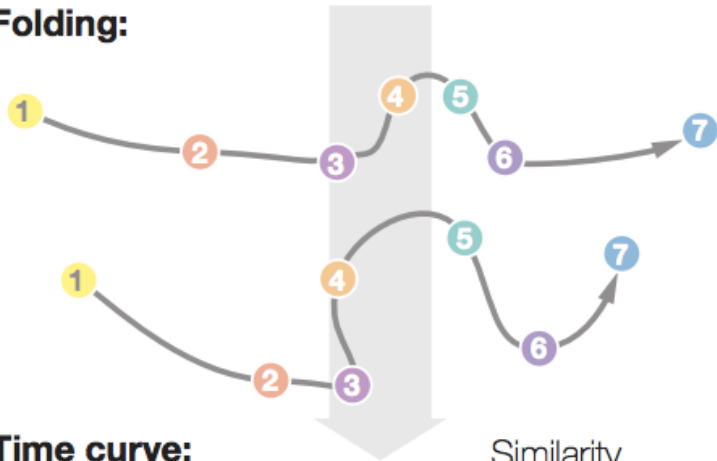
Time Curves [Bach et al. '16]

Timeline:



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

Folding:

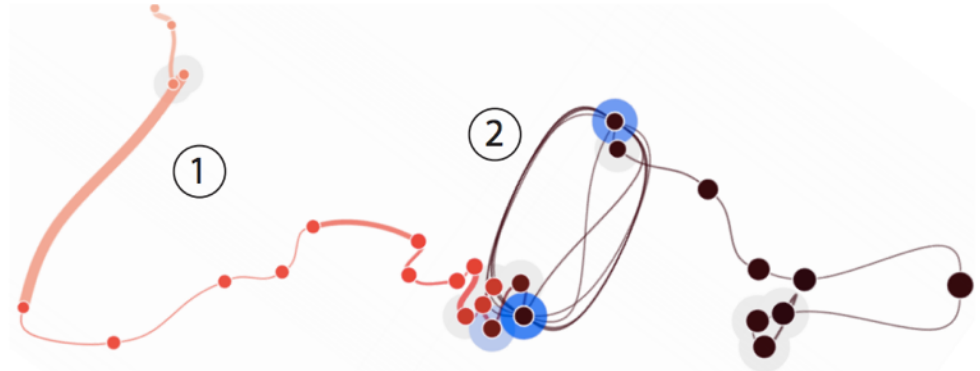


Time curve:

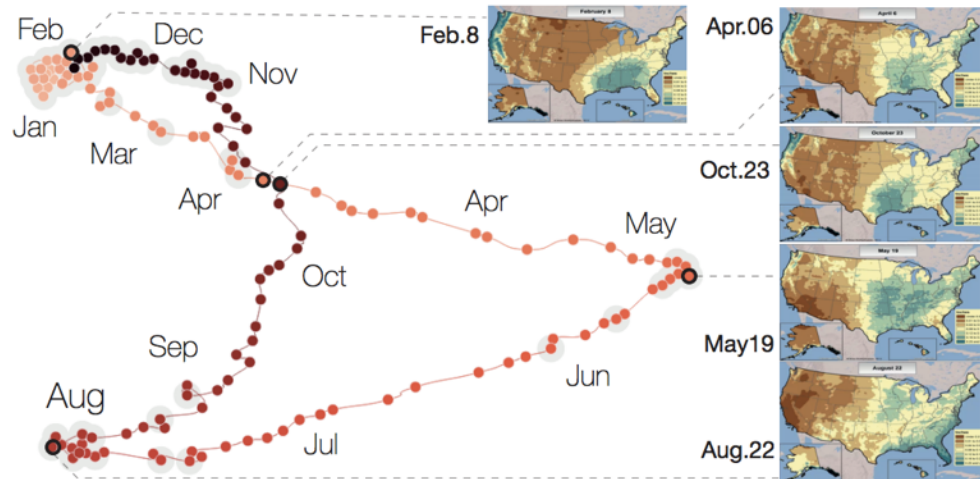


The temporal ordering of data cases is preserved.
Spatial proximity now indicates similarity.

(a) Folding time



Wikipedia "Chocolate" Article



U.S. Precipitation over 1 Year

Many Reduction Techniques!

Principal Components Analysis (PCA)

Multidimensional Scaling (MDS)

Locally Linear Embedding (LLE)

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

Isomap

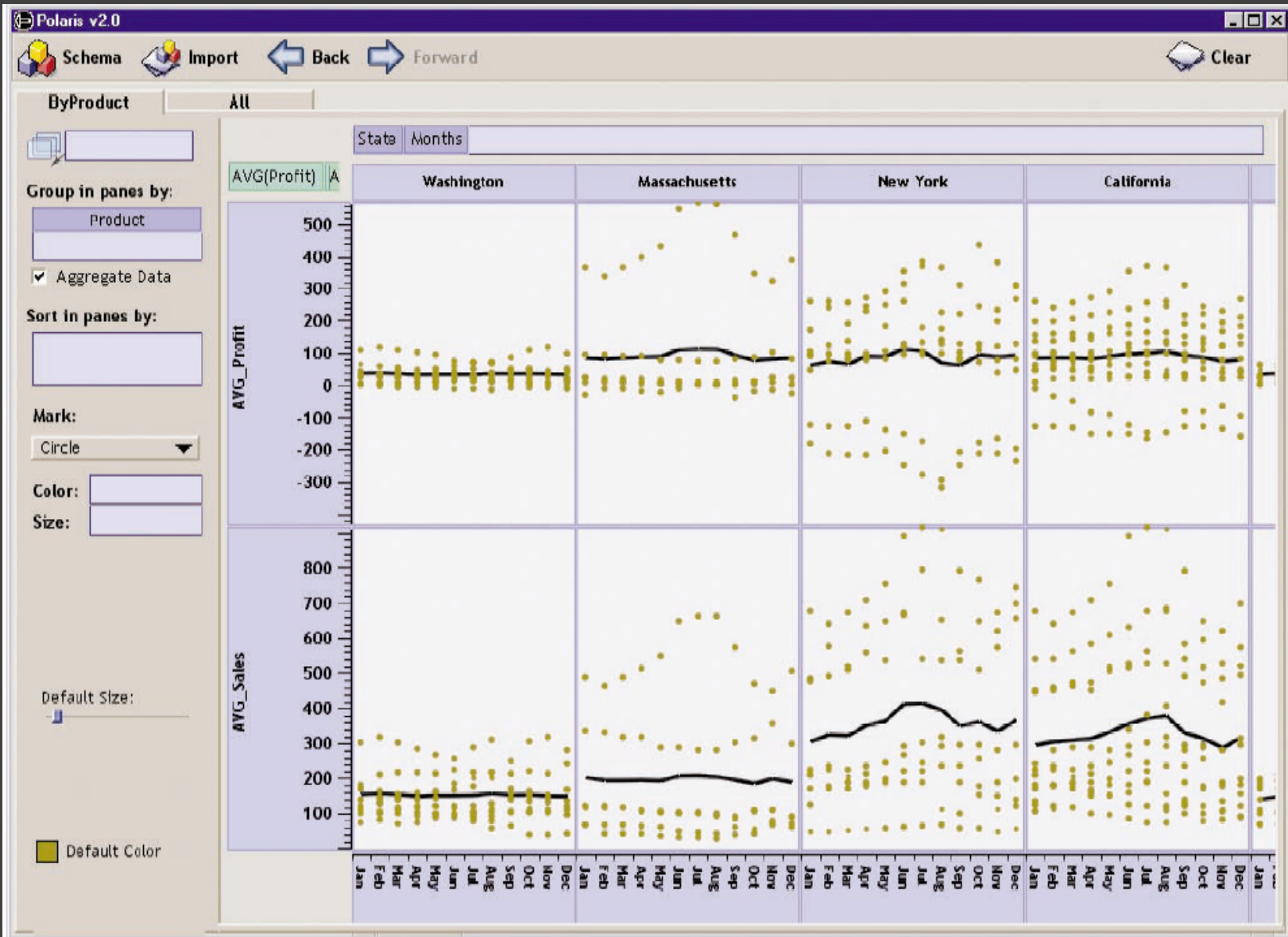
Auto-Encoder Neural Networks

Topological methods

...

Tableau / Polaris

Polaris [Stolte et al.]



Tableau

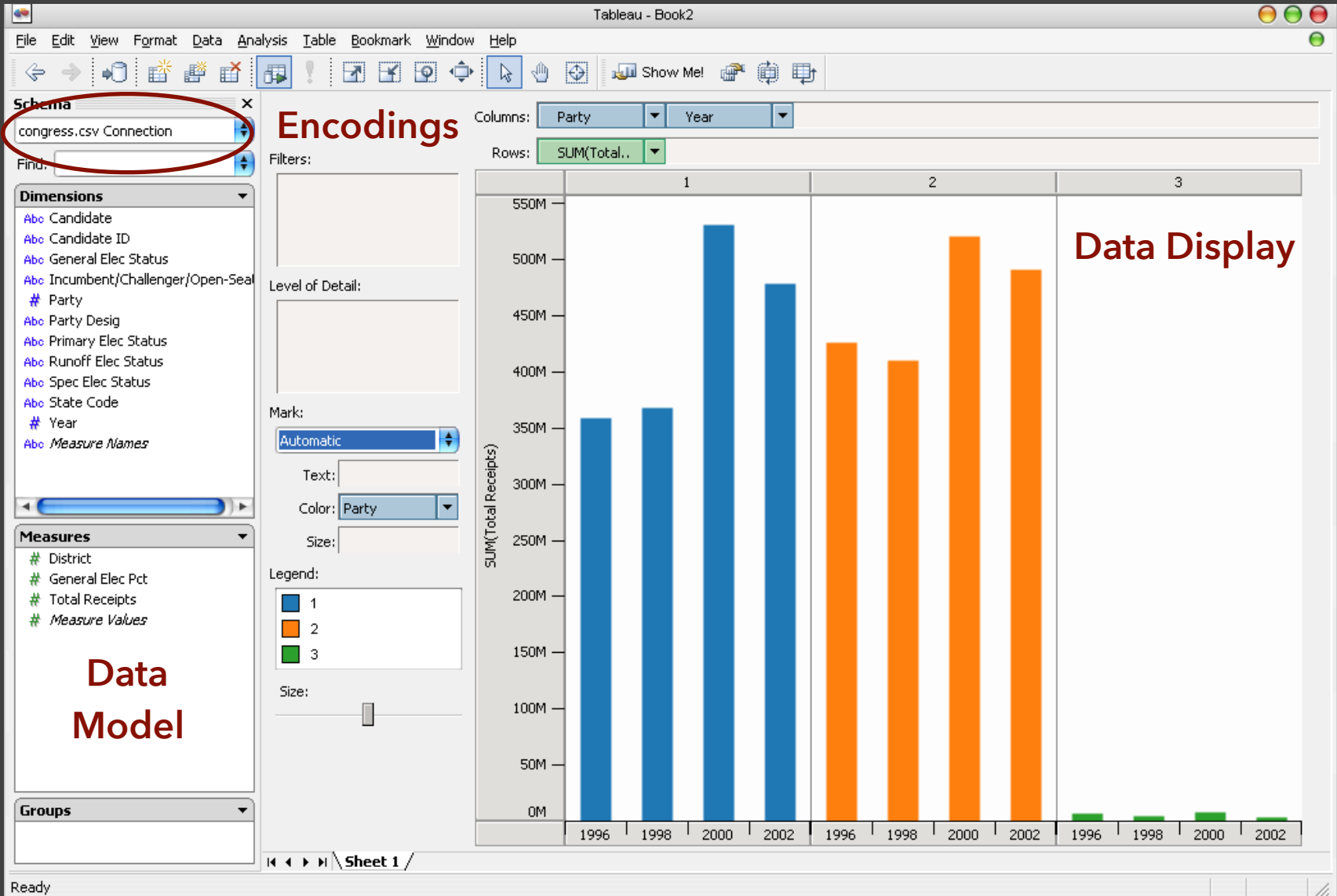


Tableau Demo

The dataset:

Federal Elections Commission Receipts

Every Congressional Candidate from 1996 to 2002

4 Election Cycles

9216 Candidacies

Dataset Schema

Year (Qi)

Candidate Code (N)

Candidate Name (N)

Incumbent / Challenger / Open-Seat (N)

Party Code (N) [1=Dem,2=Rep,3=Other]

Party Name (N)

Total Receipts (Qr)

State (N)

District (N)

This is a subset of the larger data set available from the FEC.

Hypotheses?

What might we learn from this data?

Hypotheses?

What might we learn from this data?

Correlation between receipts and winners?

Do receipts increase over time?

Which states spend the most?

Which party spends the most?

Margin of victory vs. amount spent?

Amount spent between competitors?

Tableau Demo

Tableau/Polaris Approach

Insight: can simultaneously specify both database queries and visualization

Choose data, then visualization, not vice versa

Use smart defaults for visual encodings

More recently: automate visualization design

Specifying Table Configurations

Operands are the database fields

Each operand interpreted as a set {...}

Quantitative and Ordinal fields treated differently

Three operators:

concatenation (+)

cross product (x)

nest (/)

Data | Analytics

Sample - Superstore

Dimensions

- Customer
 - Customer Name
 - Segment
- Order
- Location
- Product
 - Category
 - Sub-Category
 - Manufacturer
 - Product Name
- Profit (bin)
- Region
- Measure Names

Measures

- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
- Measure Values

Pages

Filters

Marks

Automatic

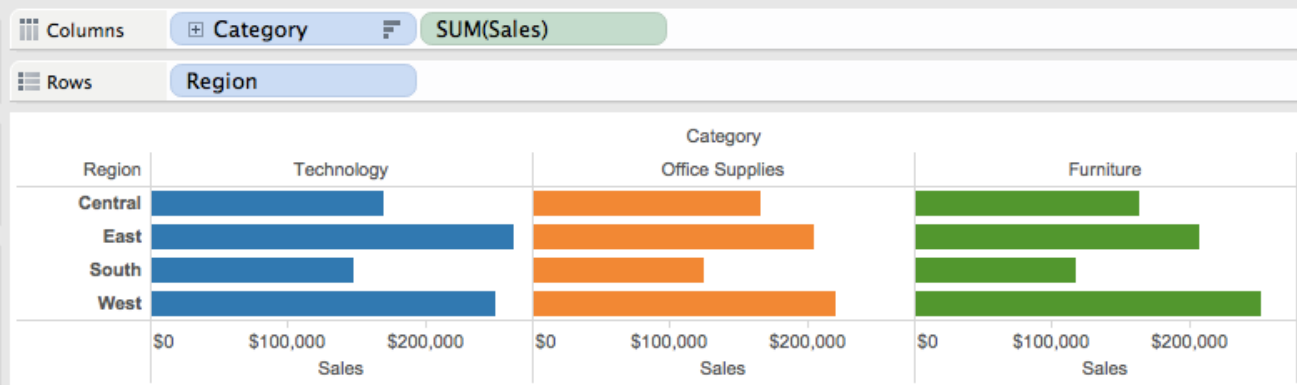
Color Size Label

Detail Tooltip

Category

Category

- Technology
- Office Supplies
- Furniture



Data | Analytics

Sample - Superstore

Dimensions

- Customer
 - Customer Name
 - Segment
- Order
 - Location
- Product
 - Category
 - Sub-Category
 - Manufacturer
 - Product Name
- Profit (bin)
- Region
- Measure Names

Measures

- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
- Measure Values

Pages

Filters

Marks

Automatic

Color Size Label

Detail Tooltip

Category

Category

- Technology
- Office Supplies
- Furniture



Data | Analytics

Sample - Superstore

Dimensions

- Customer
 - Customer Name
 - Segment
- Order
- Location
- Product
 - Category
 - Sub-Category
 - Manufacturer
 - Product Name
- Profit (bin)
- Region
- Measure Names

Measures

- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
- Measure Values

Pages

Filters

Marks

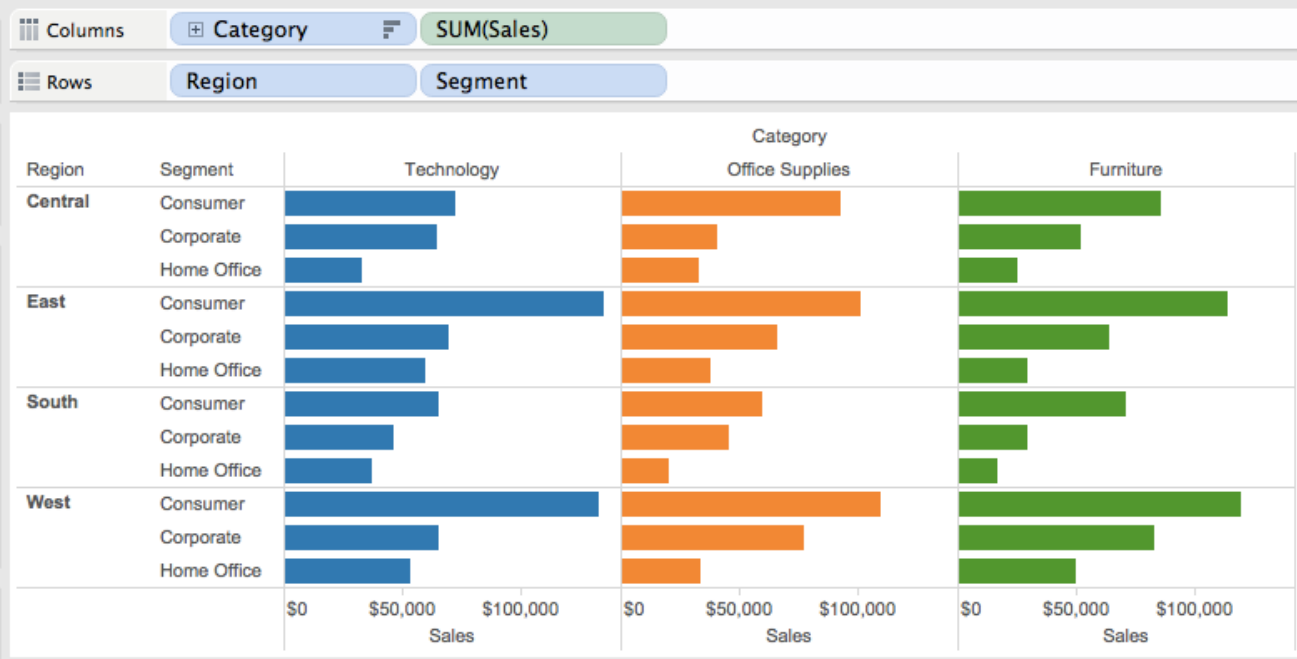
Automatic

Color Size Label

Detail Tooltip

Category

Technology
Office Supplies
Furniture



Data | Analytics

Sample - Superstore

Dimensions

- Customer
 - Customer Name
 - Segment
- Order
- Location
- Product
 - Category
 - Sub-Category
 - Manufacturer
 - Product Name
- Profit (bin)
- Region
- Measure Names

Measures

- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
- Measure Values

Pages

Filters

Marks

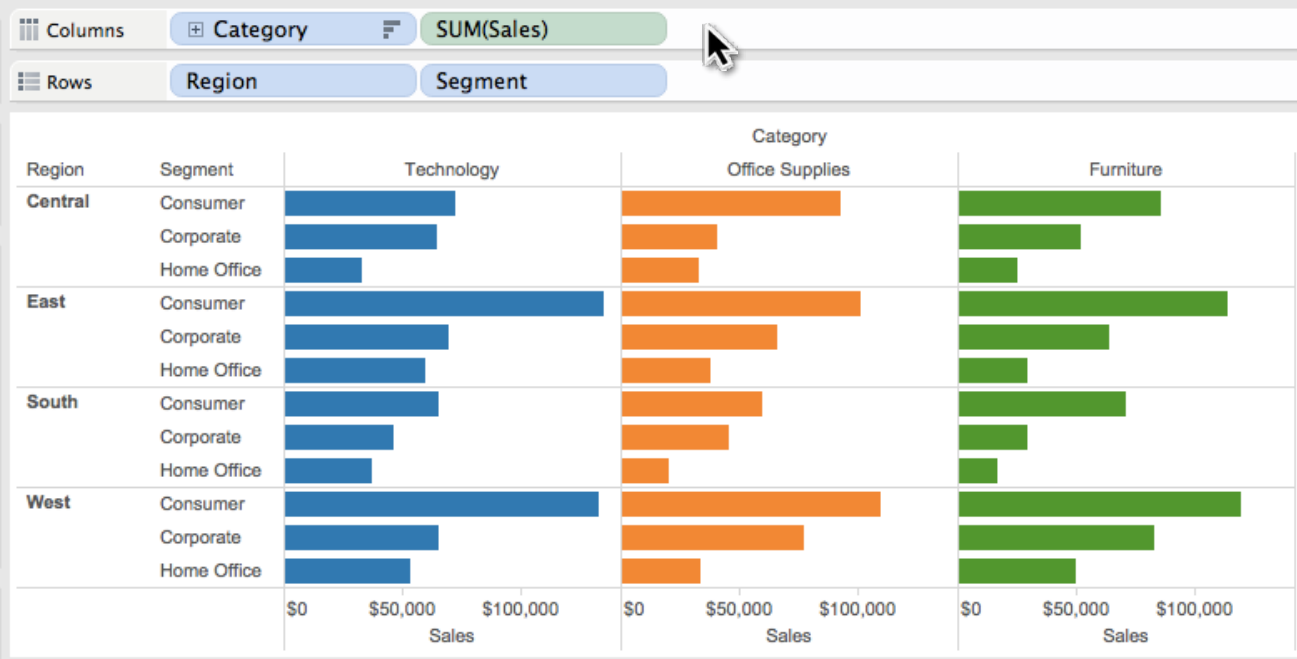
Automatic

Color Size Label

Detail Tooltip

Category

Technology
Office Supplies
Furniture



Data | Analytics

Sample - Superstore

Dimensions

- Customer
 - Customer Name
 - Segment
- Order
- Location
- Product
 - Category
 - Sub-Category
 - Manufacturer
 - Product Name
- Profit (bin)
- Region
- Measure Names

Measures

- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
- Measure Values

Pages

Filters

Marks

All

Automatic

Color Size Label

Detail Tooltip

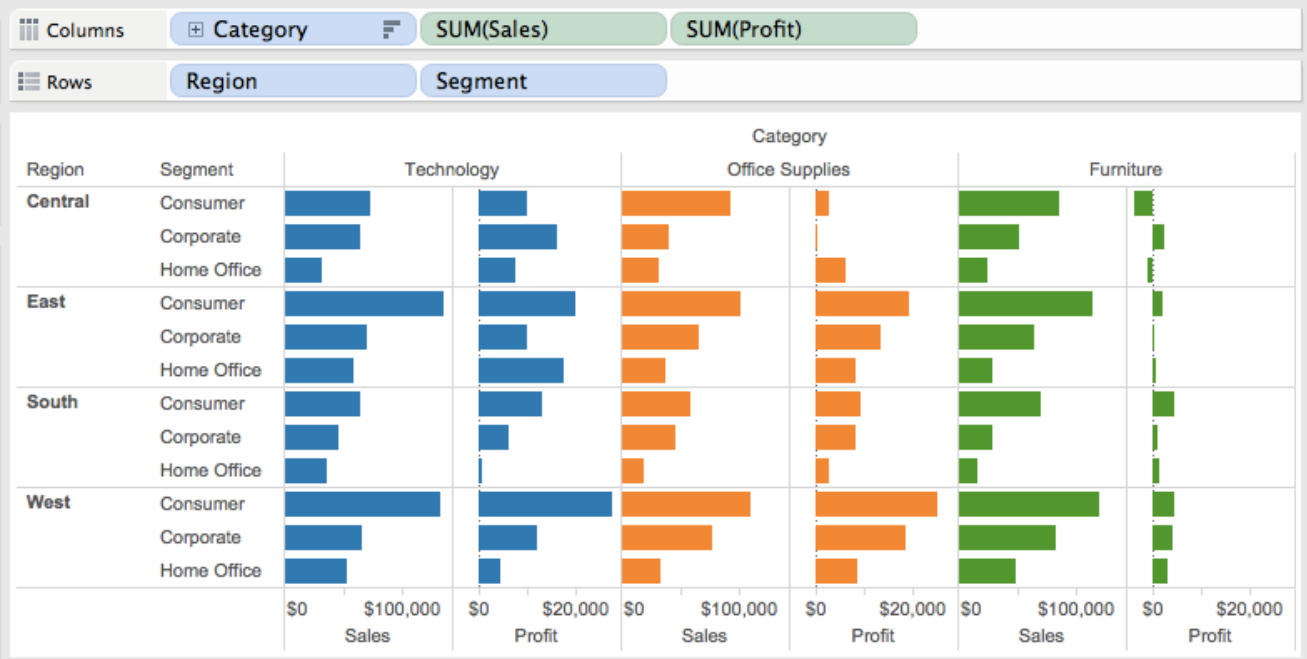
Category

SUM(Sales)

SUM(Profit)

Category

- Technology
- Office Supplies
- Furniture



Columns: Category, ~~SUM(Sales)~~, SUM(Profit)
 Rows: Region, Segment
 -> GROUP BY Category, Region, Segment

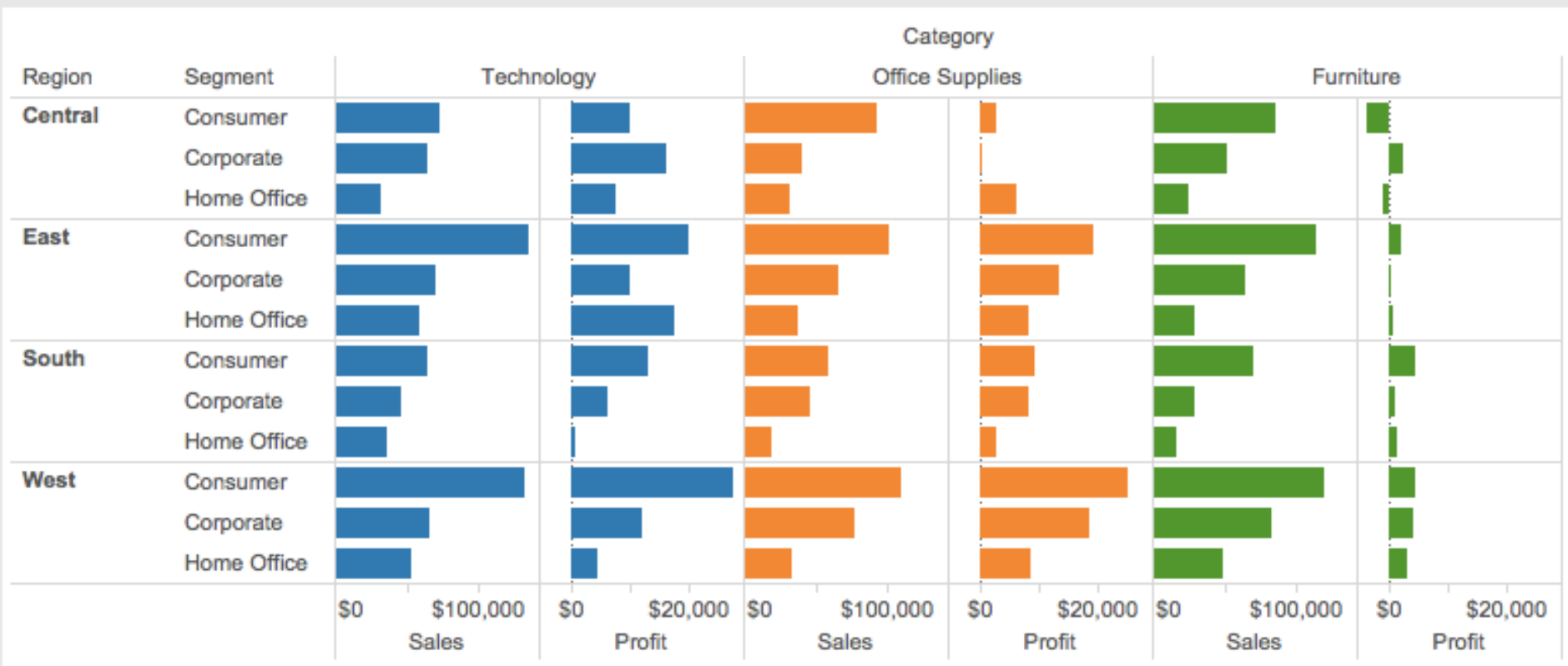


Table Algebra: Operands

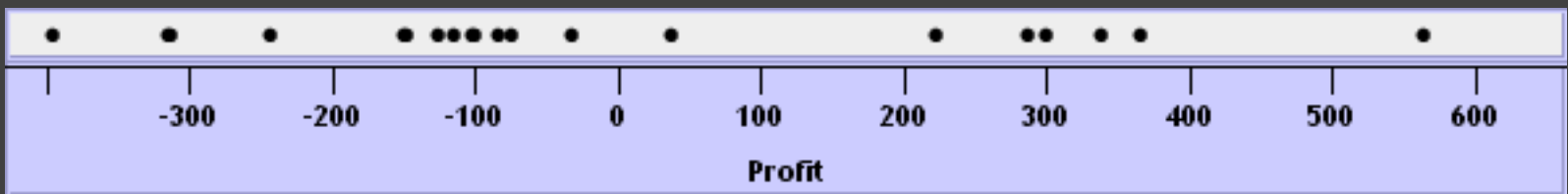
Ordinal fields: interpret domain as a set that partitions table into rows and columns.

Quarter = {(Qtr1),(Qtr2),(Qtr3),(Qtr4)} ->

Qtr1	Qtr2	Qtr3	Qtr4
95892	101760	105282	98225

Quantitative fields: treat domain as single element set and encode spatially as axes.

Profit = {(Profit[-410,650])} ->



Concatenation (+) Operator

Ordered union of set interpretations

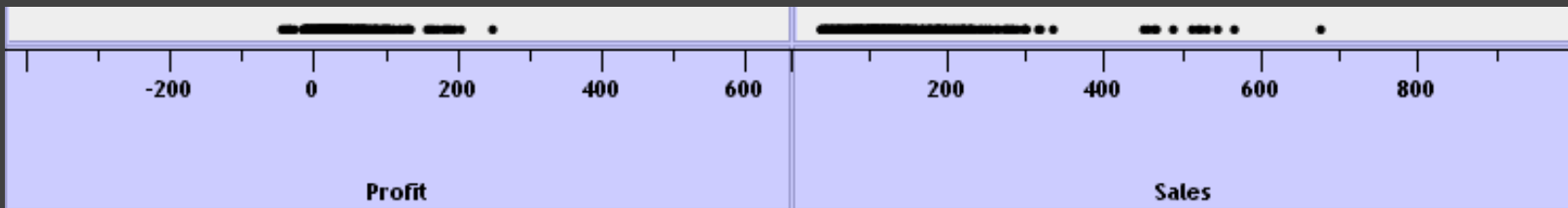
Quarter + Product Type

= {(Qtr1),(Qtr2),(Qtr3),(Qtr4)} + {(Coffee), (Espresso)}

= {(Qtr1),(Qtr2),(Qtr3),(Qtr4),(Coffee),(Espresso)}

Qtr1	Qtr2	Qtr3	Qtr4	Coffee	Espresso
48	59	57	53	151	21

Profit + Sales = {(Profit[-310,620]),(Sales[0,1000])}



Cross (x) Operator

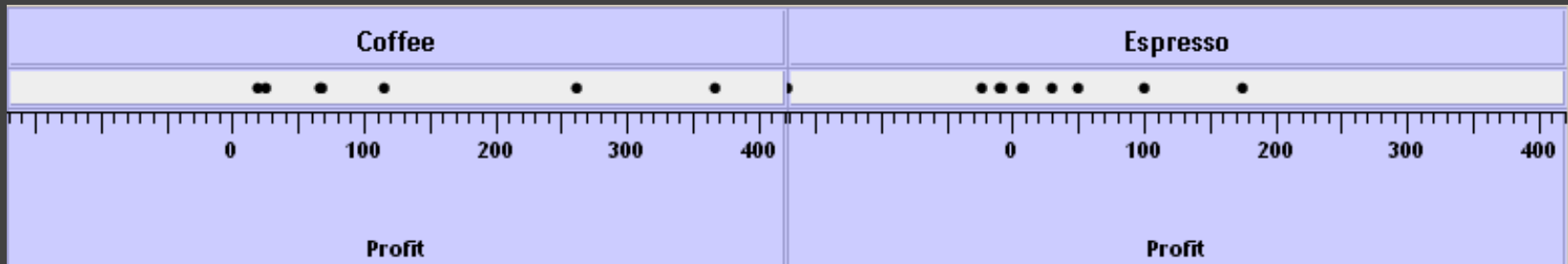
Cross-product of set interpretations

Quarter x Product Type =

{(Qtr1, Coffee), (Qtr1, Tea), (Qtr2, Coffee), (Qtr2, Tea), (Qtr3, Coffee), (Qtr3, Tea), (Qtr4, Coffee), (Qtr4, Tea)}

Qtr1		Qtr2		Qtr3		Qtr4	
Coffee	Espresso	Coffee	Espresso	Coffee	Espresso	Coffee	Espresso
131	19	160	20	178	12	134	33

Product Type x Profit =



Nest (/) Operator

Cross-product filtered by existing records

Quarter x Month ->

creates twelve entries for each quarter. i.e.,
(Qtr1, December)

Quarter / Month ->

creates three entries per quarter based on
tuples in database (not semantics)

Table Algebra

The operators (+, x, /) and operands (O, Q) provide an *algebra* for tabular visualization.

Algebraic statements are then mapped to:

Visualizations - trellis plot partitions, visual encodings

Queries - selection, projection, group-by aggregation

In Tableau, users make statements via drag-and-drop

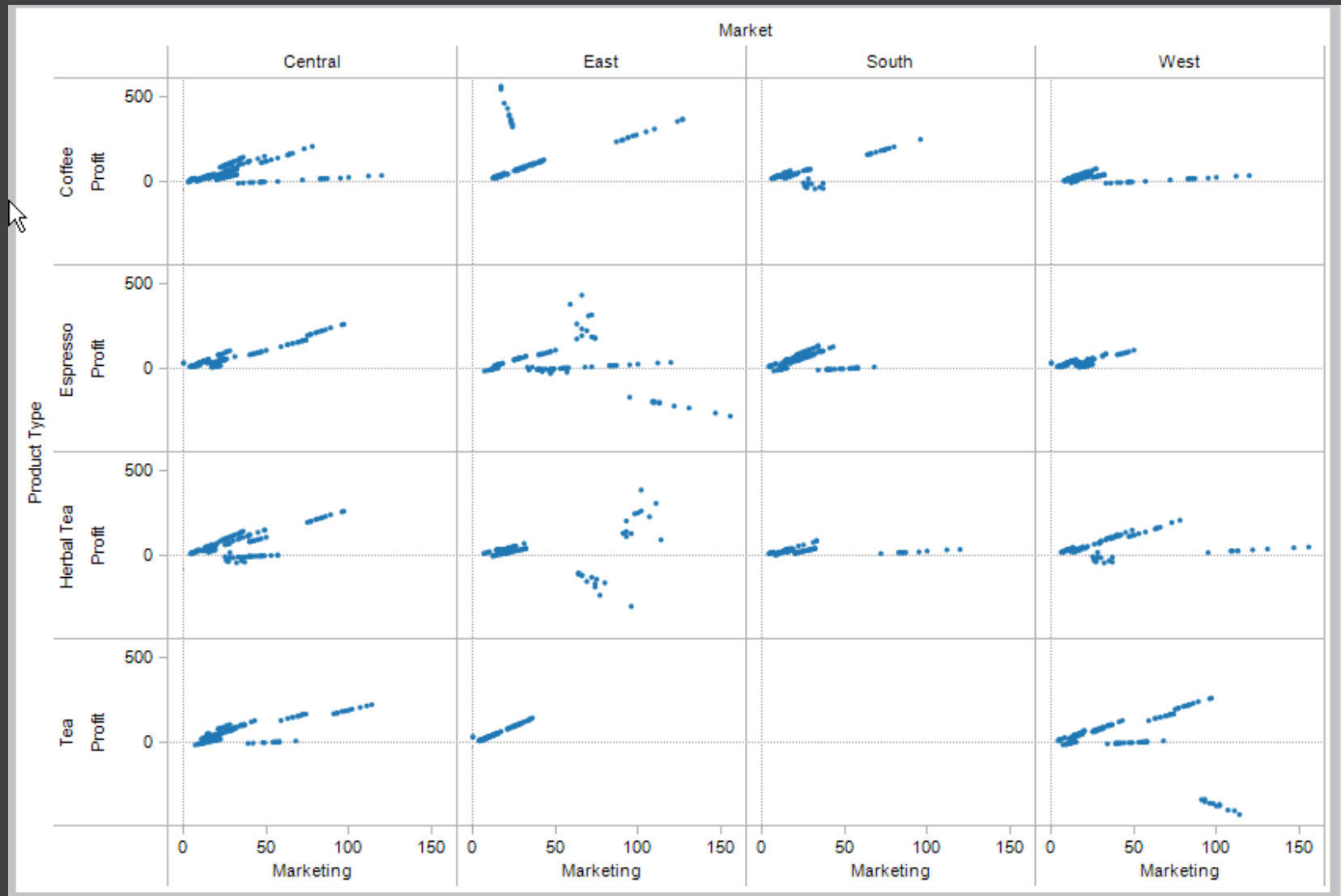
Note that this specifies operands *NOT* operators!

Operators are inferred by data type (O, Q)

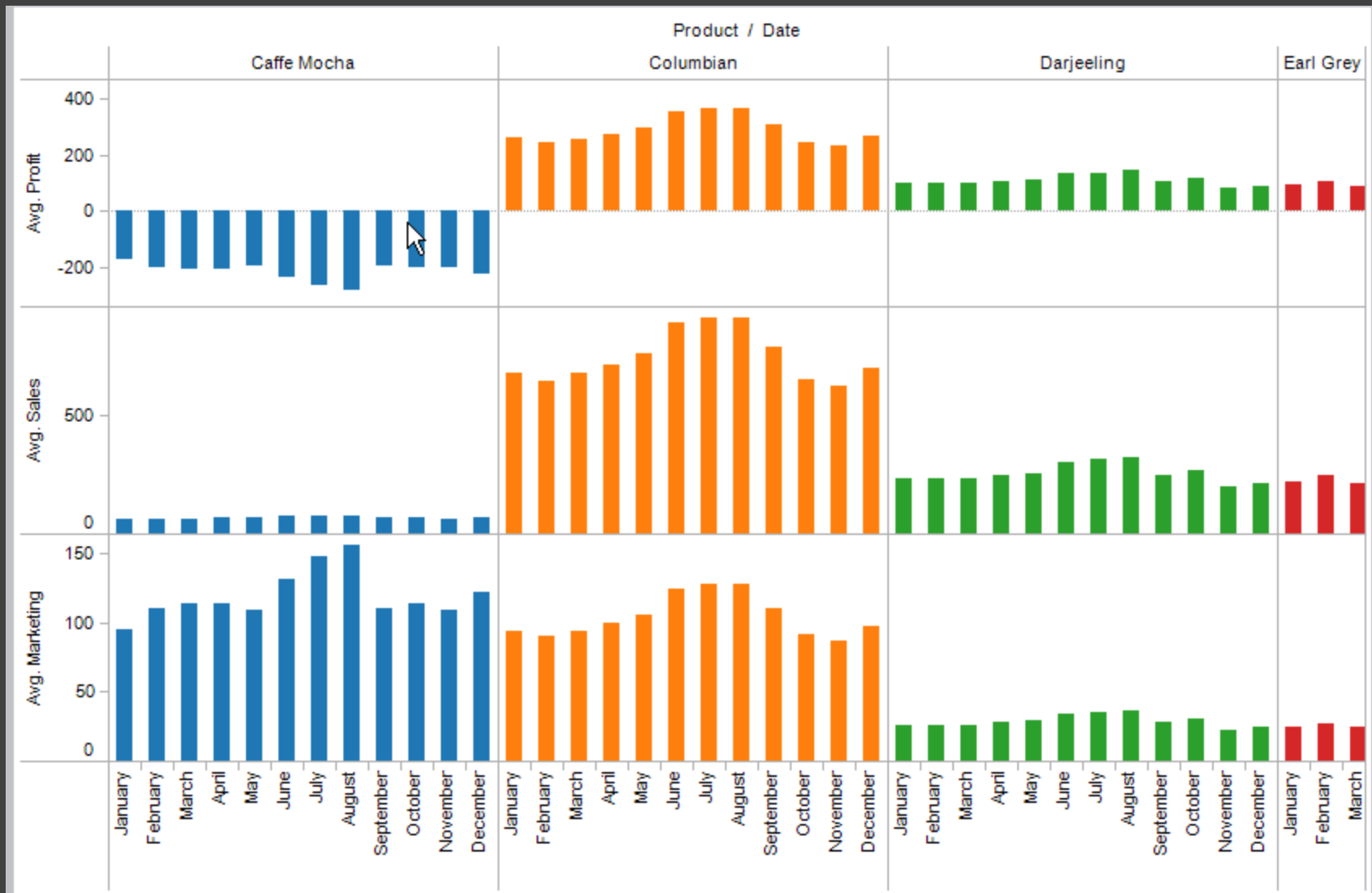
Ordinal-Ordinal

State	Product Type			
	Coffee	Espresso	Herbal Tea	Tea
Colorado	●	●	●	●
Connecticut	●	●	●	●
Florida	●	●	●	●
Illinois	●	●	●	●
Iowa	●	●	●	●
Louisiana	●	●	●	●
Massachusetts	●	●	●	●
Missouri	●	●	●	●
Nevada	●	●	●	●
New Hampshire	●	●	●	●
New Mexico	●	●	●	●
New York	●	●	●	●
Ohio	●	●	●	●
Oklahoma	●	●	●	●
Oregon	●	●	●	●
Texas	●	●	●	●
Utah	●	●	●	●
Washington	●	●	●	●
Wisconsin	●	●	●	●

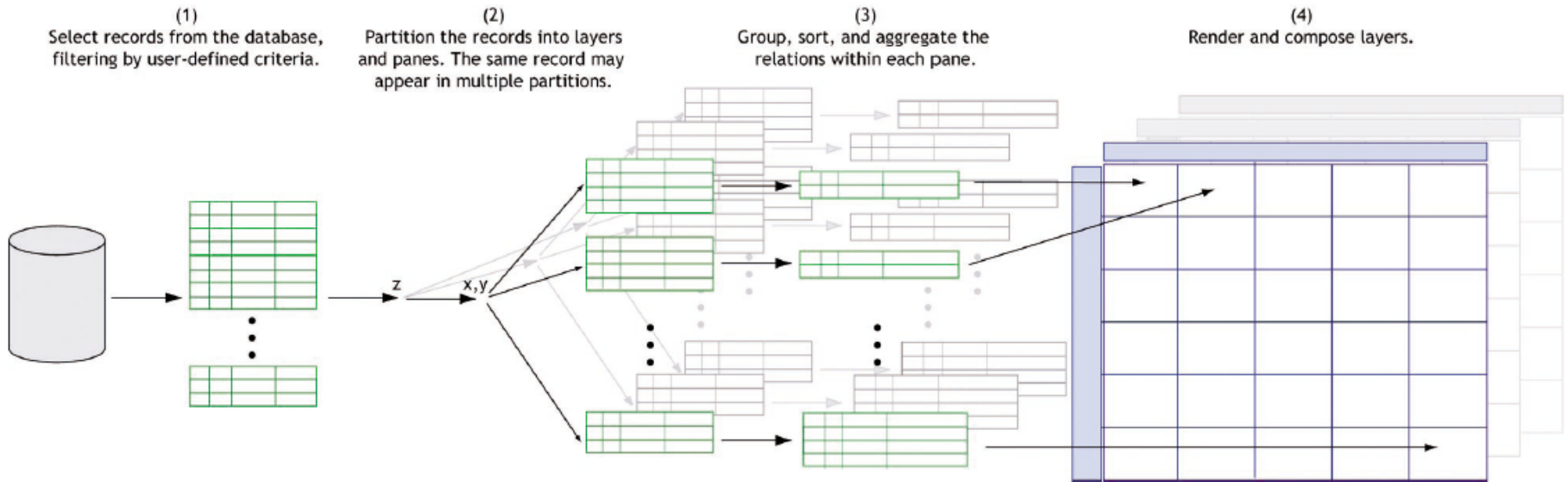
Quantitative-Quantitative



Ordinal-Quantitative



Querying the Database



Visualizing Multiple Dimensions

Strategies:

Avoid “over-encoding”

Use space and small multiples intelligently

Reduce the problem space

Use interaction to generate *relevant* views

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