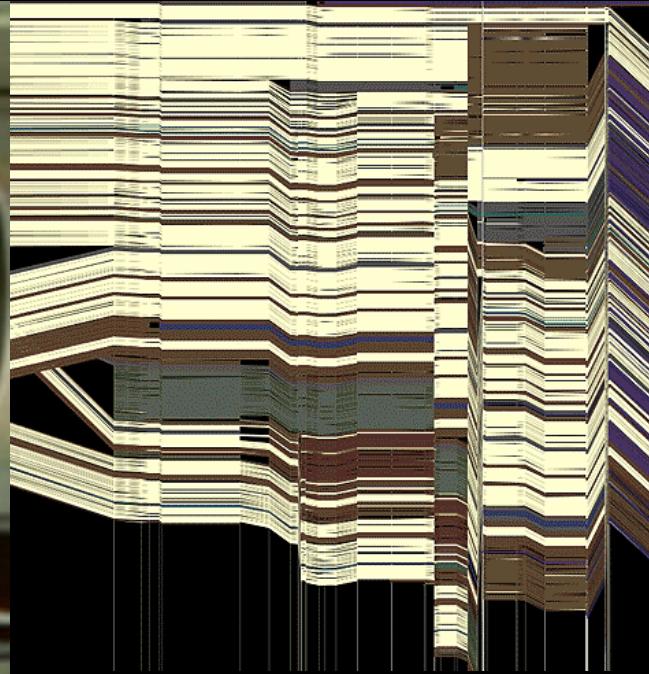
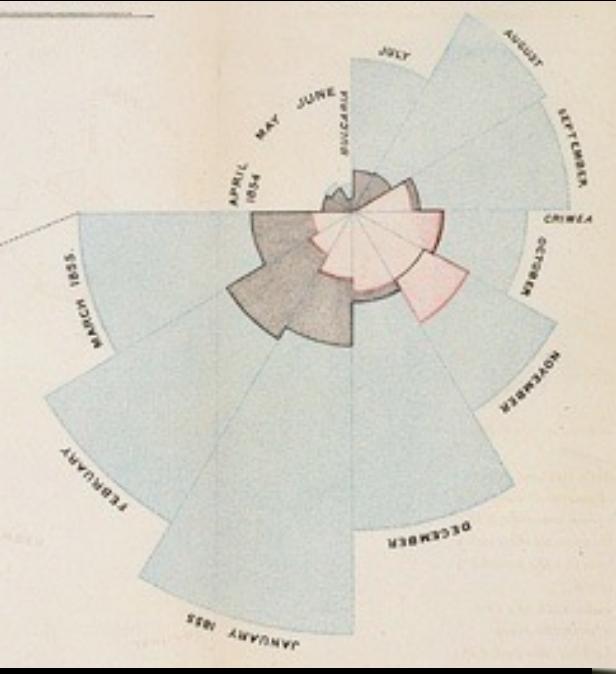


CSE 512 - Data Visualization

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

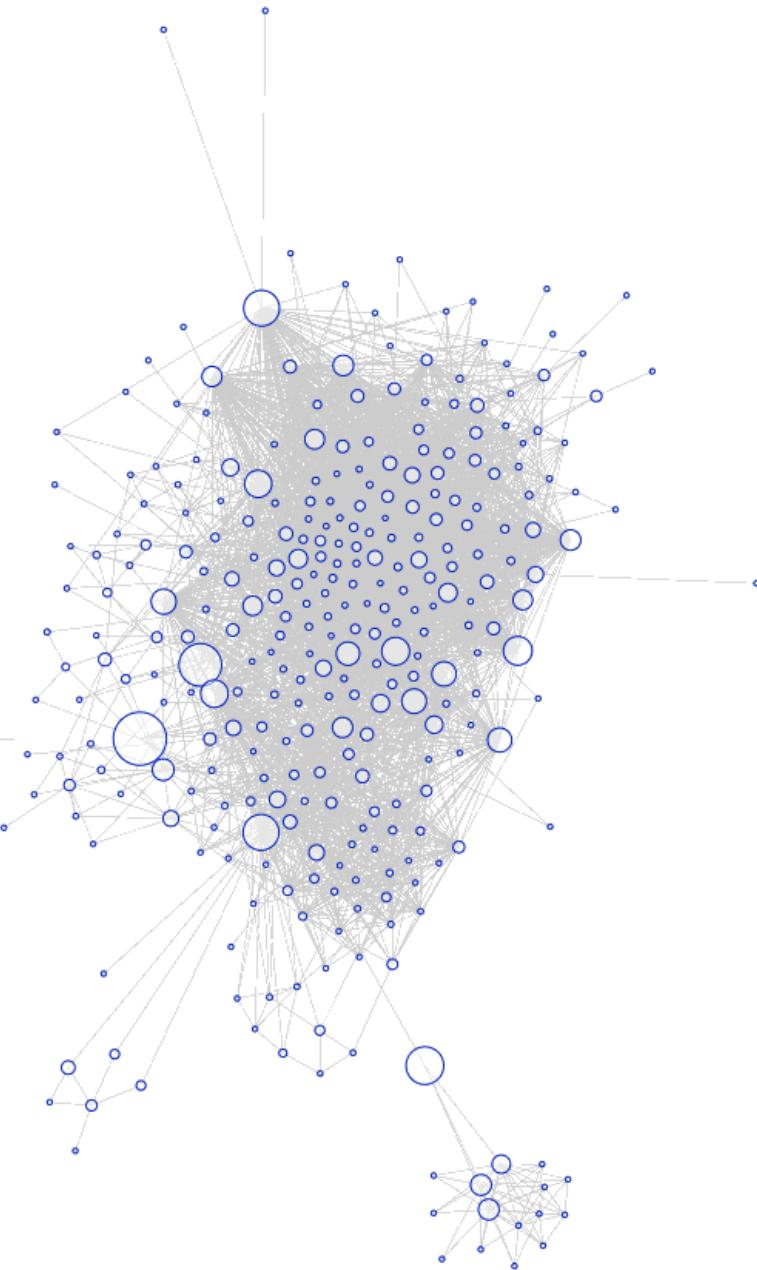
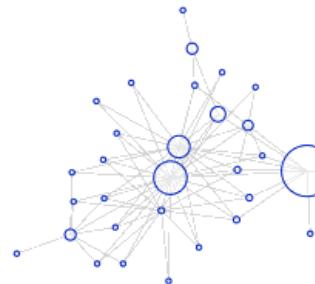
Graph Viewer

Roll-up by:

Visualization:

Sort by:

Edge centrality filters:



Images

Animate

Graph Viewer

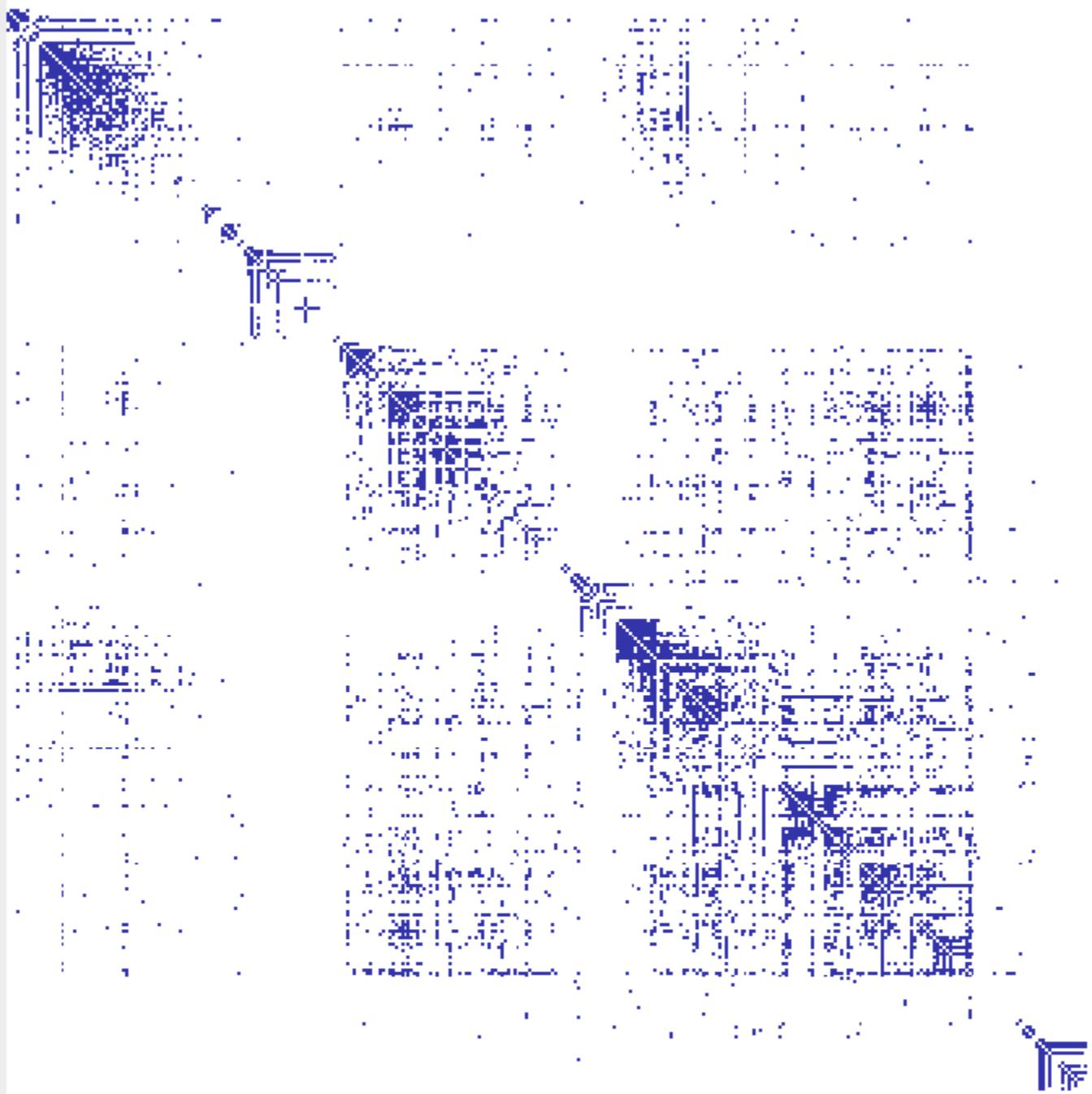
Graph Viewer

Roll-up by:

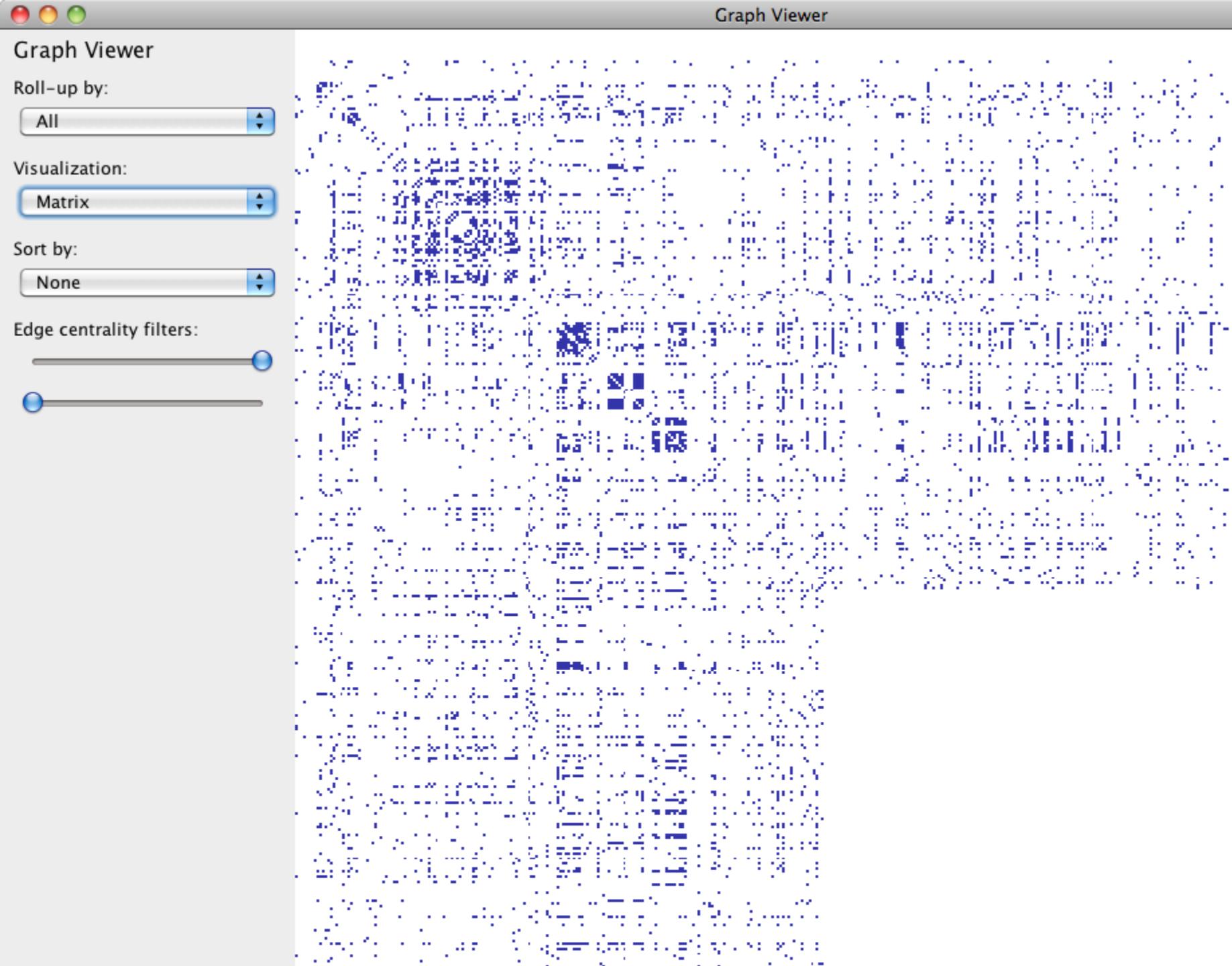
Visualization:

Sort by:

Edge centrality filters:



Graph Viewer

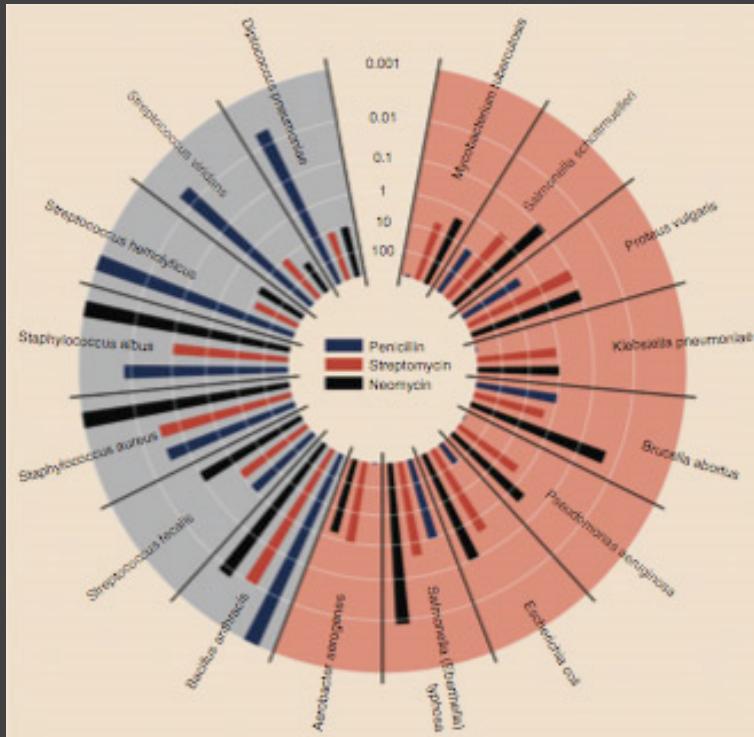


Antibiotic Effectiveness

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

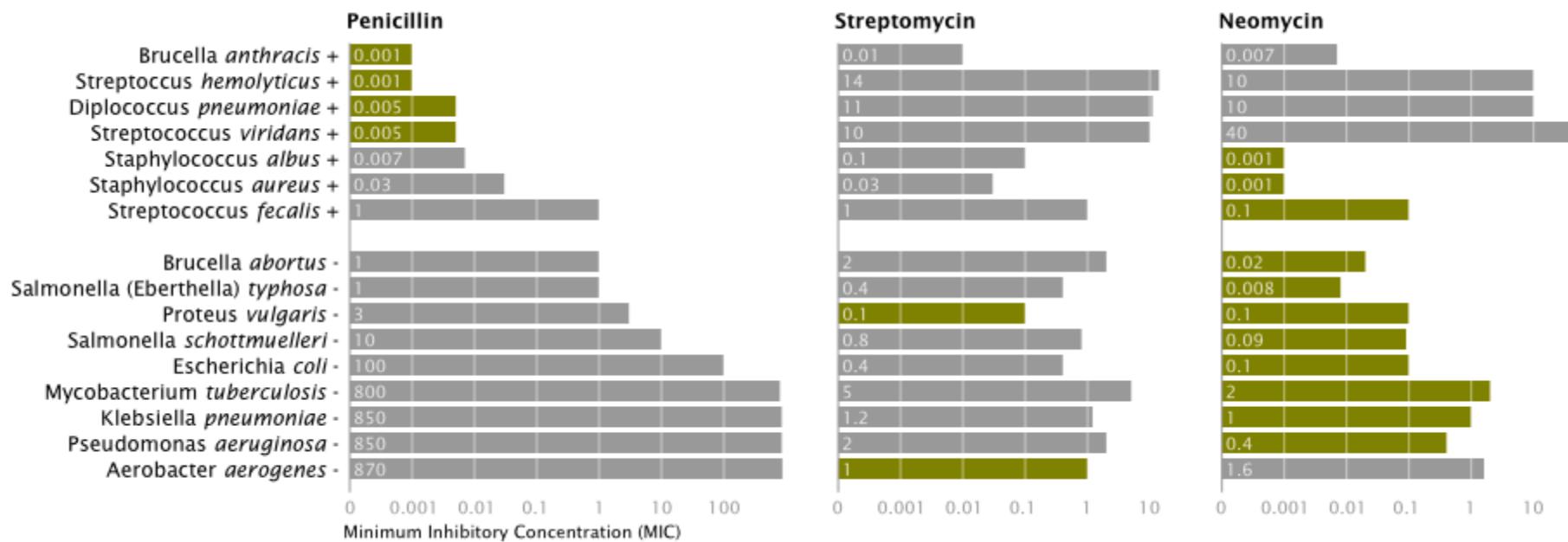
How do the drugs compare?



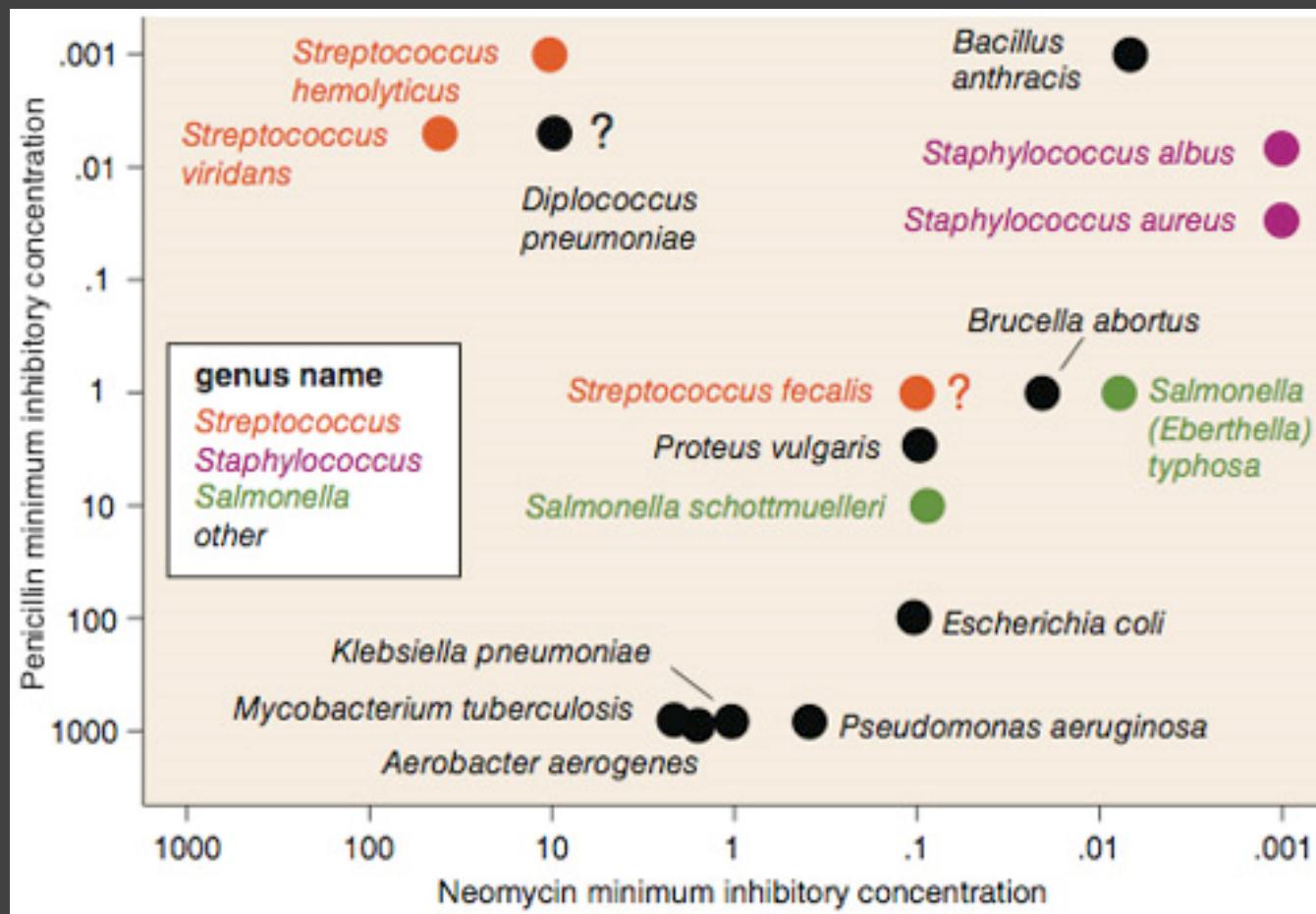
Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	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?



Mike Bostock
Stanford CS448B, Winter 2009

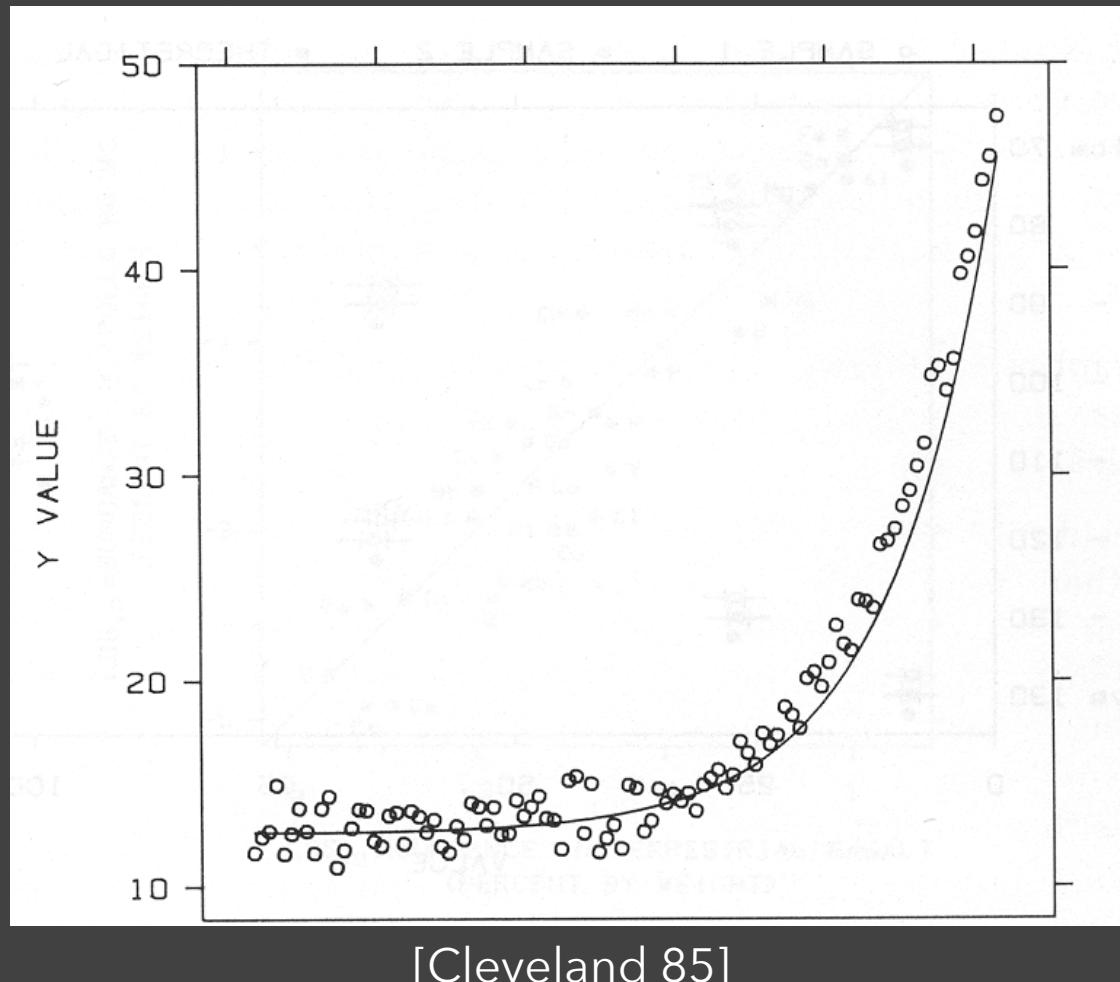


Do the bacteria group by resistance?
Do different drugs correlate?

Wainer & Lysen
American Scientist, 2009

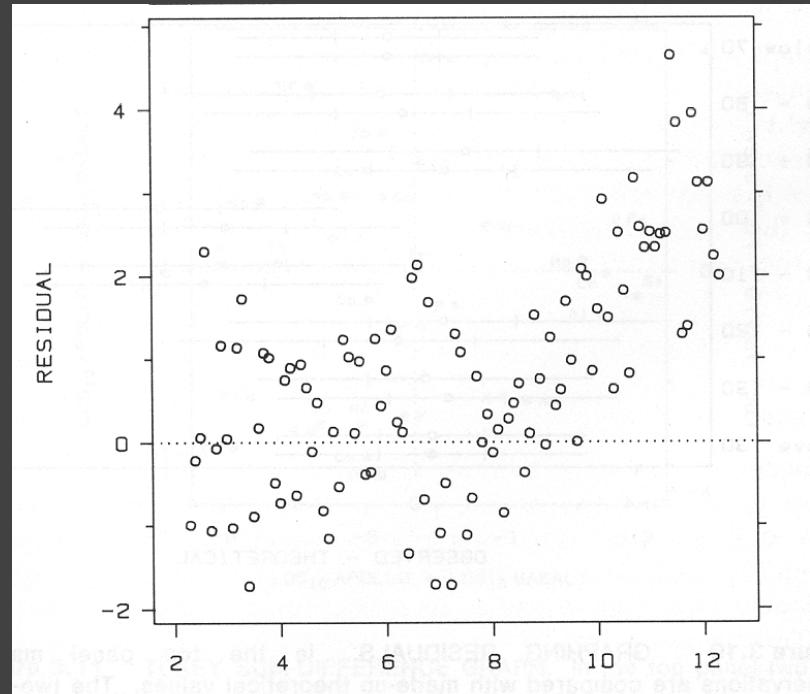
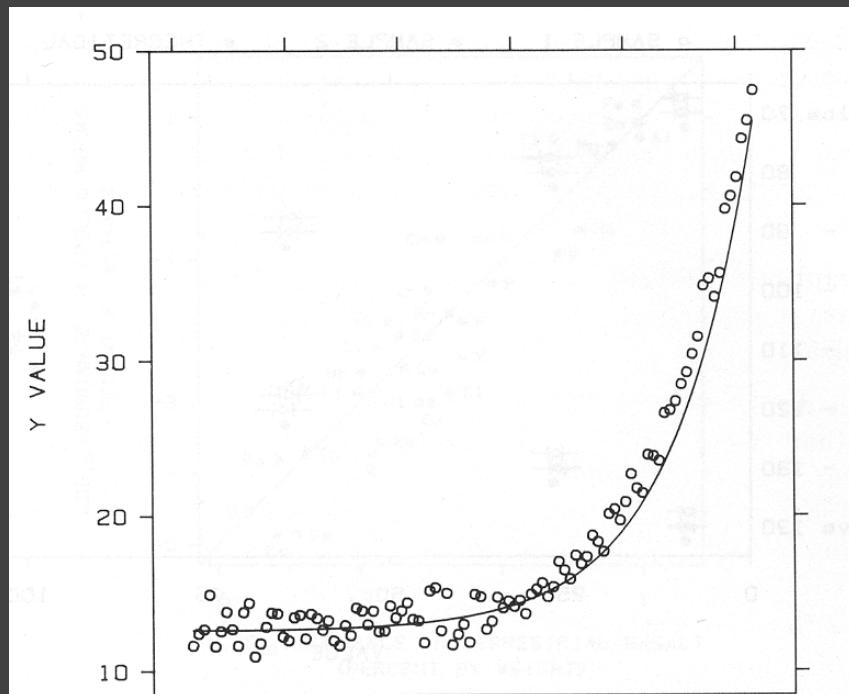
Transforming Data

How well does the curve fit the data?



Plot the Residuals

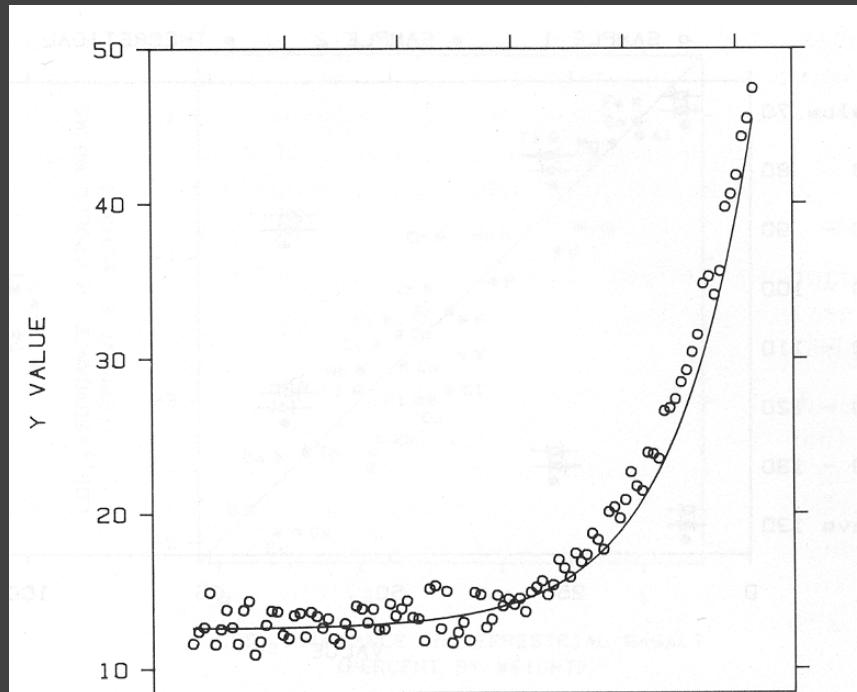
Plot vertical distance from best fit curve
Residual graph shows accuracy of fit



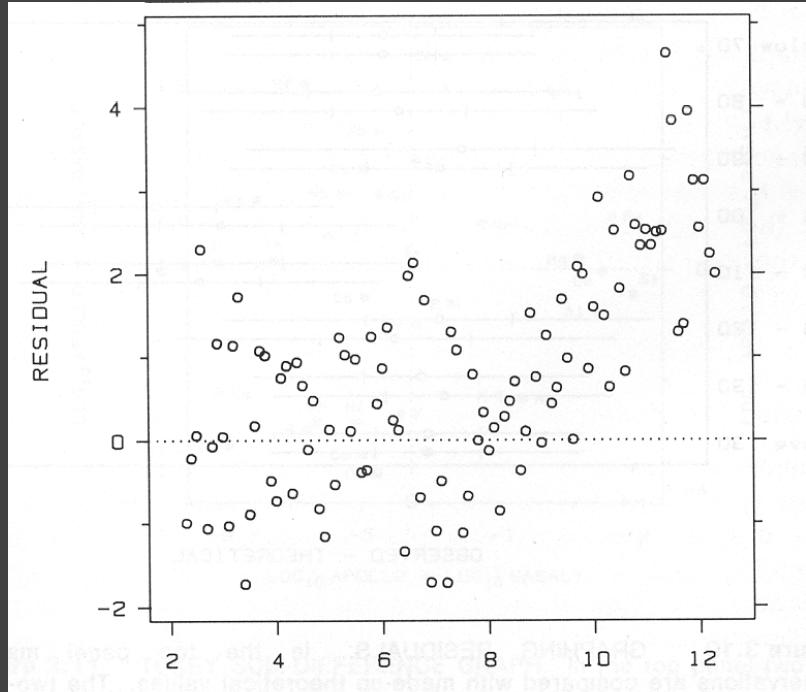
[Cleveland 85]

Multiple Plotting Options

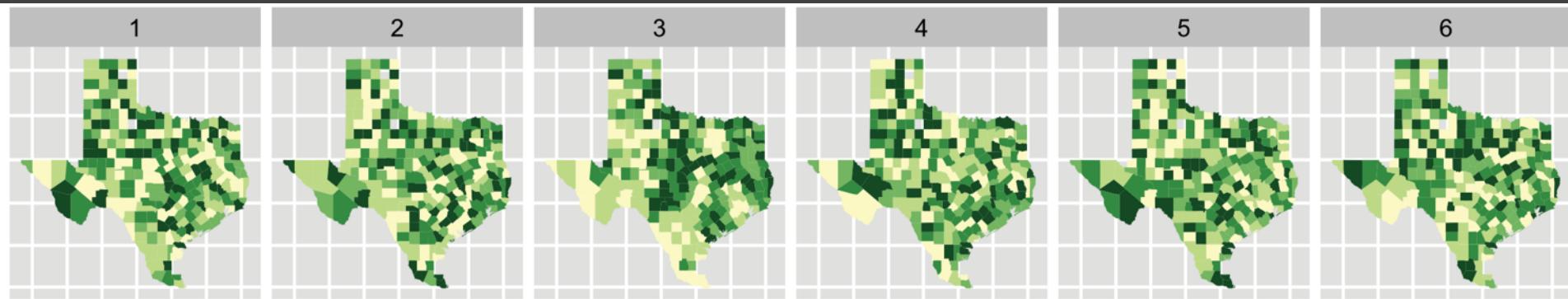
Plot model in data space



Plot data in model space



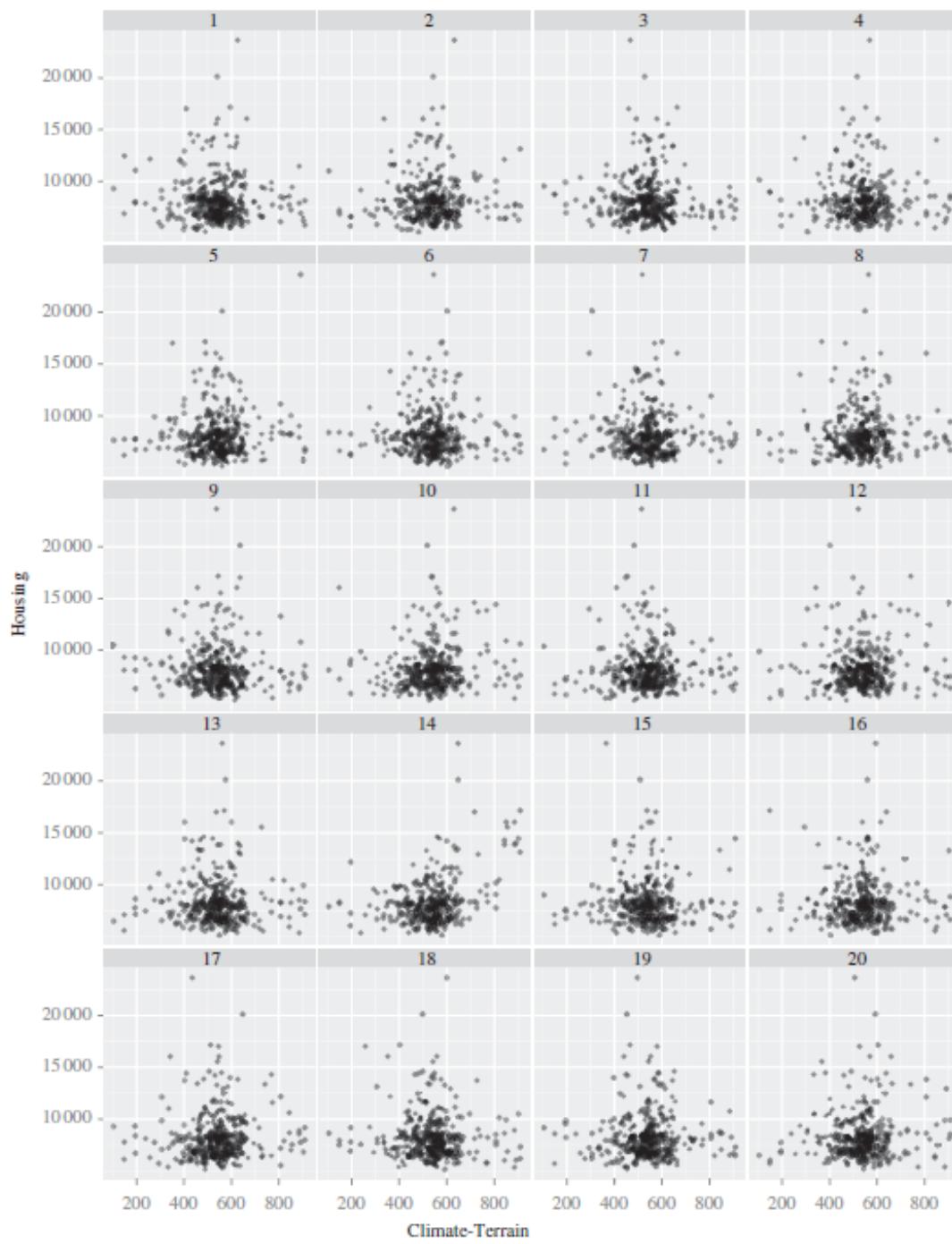
[Cleveland 85]

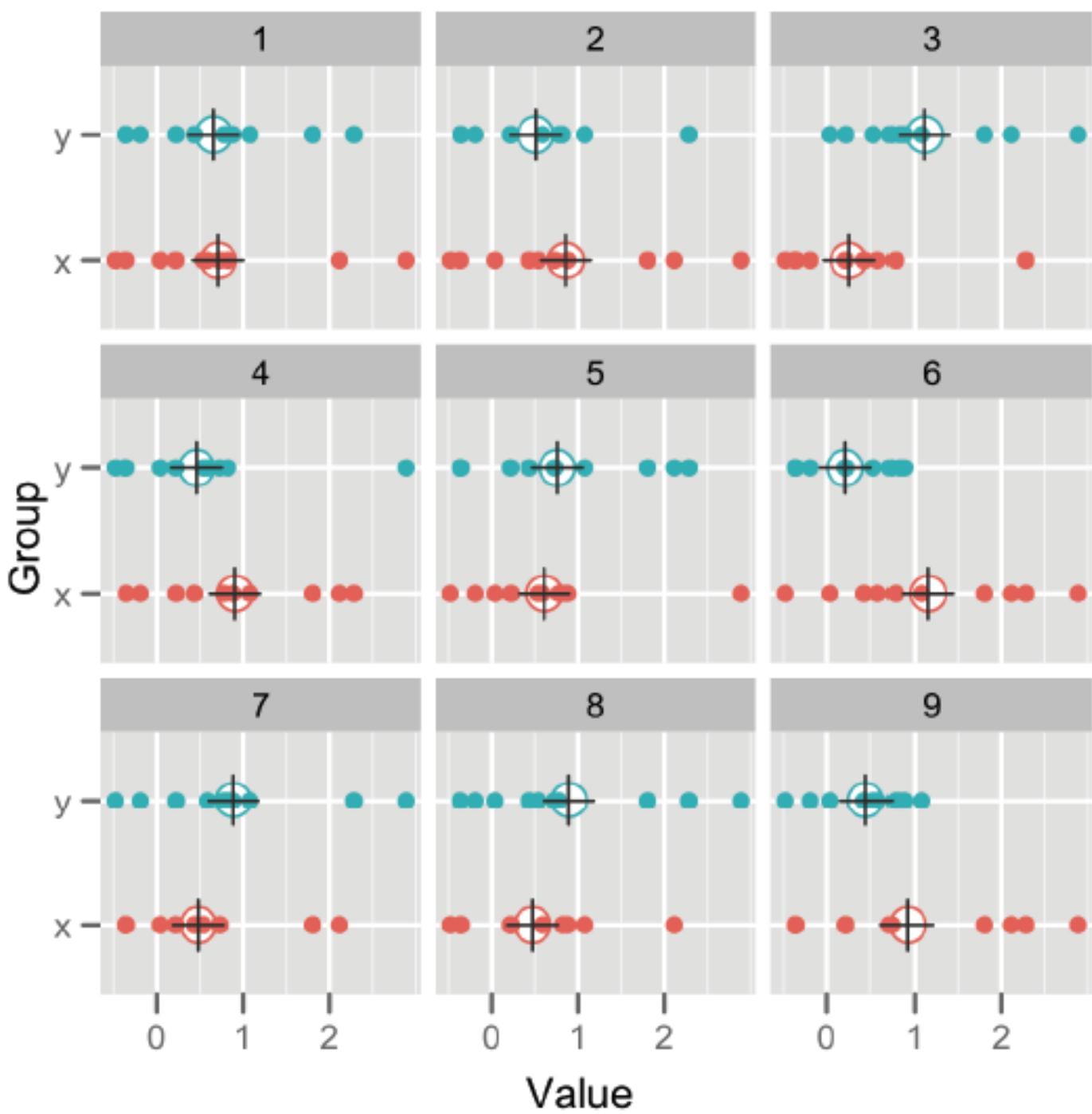


Choropleth maps of cancer deaths in Texas.

One plot shows a real data set. The others are simulated under the null hypothesis of spatial independence.

Can you spot the real data? If so, you have some evidence of spatial dependence in the data.





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

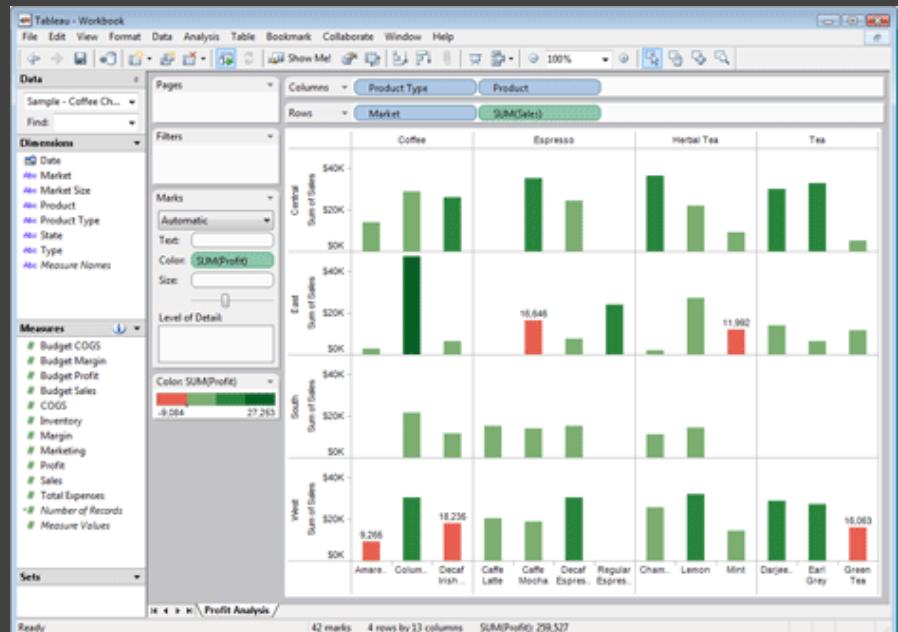
Interact with data

Refine your questions

Make a notebook

Keep record of your analysis

Prepare a final graphic and caption



Due by 5:00pm
Friday, April 17

Vis Tools Tutorial

Today, Tuesday April 14

3pm to 4:30pm in CSE 305

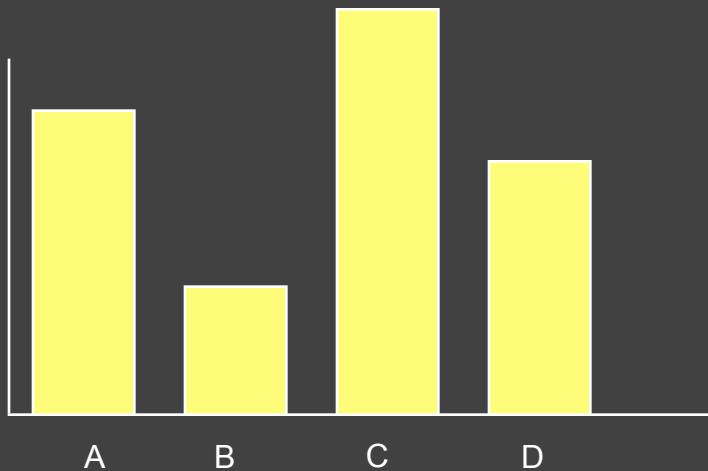
Become a **Tableau** power user

Learn **matplotlib**, valuable for iPython notebooks

See **new tools** coming out of CSE research

The Design Space of Visual Encodings

Univariate Data



	factors			
	A	B	C	
1				

variable

Univariate Data

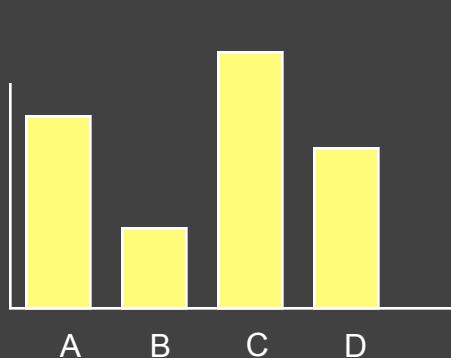
factors

1	A	B	C	

variable



Tukey box plot



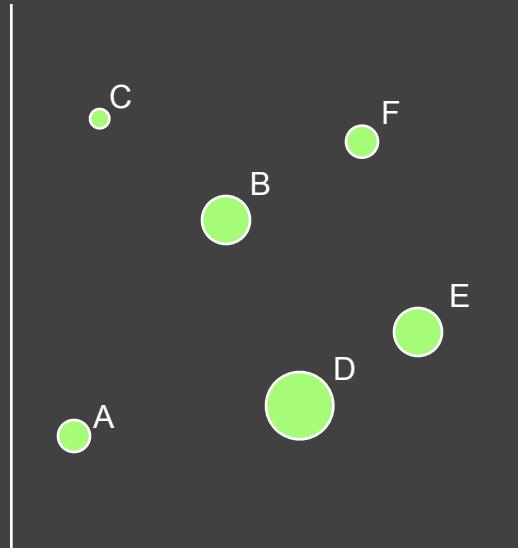
Bivariate Data

	A	B	C
1			
2			



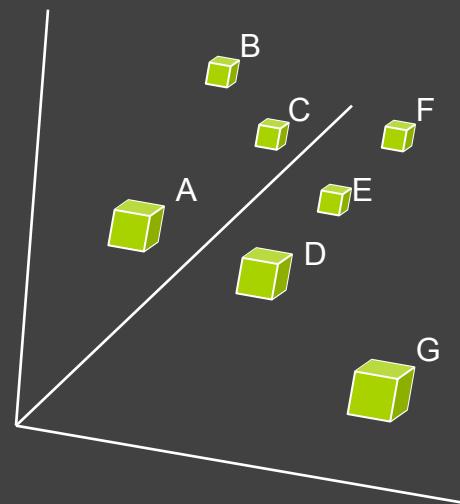
Scatter plot is common

Trivariate Data



	A	B	C
1			
2			
3			

3D scatter plot is possible

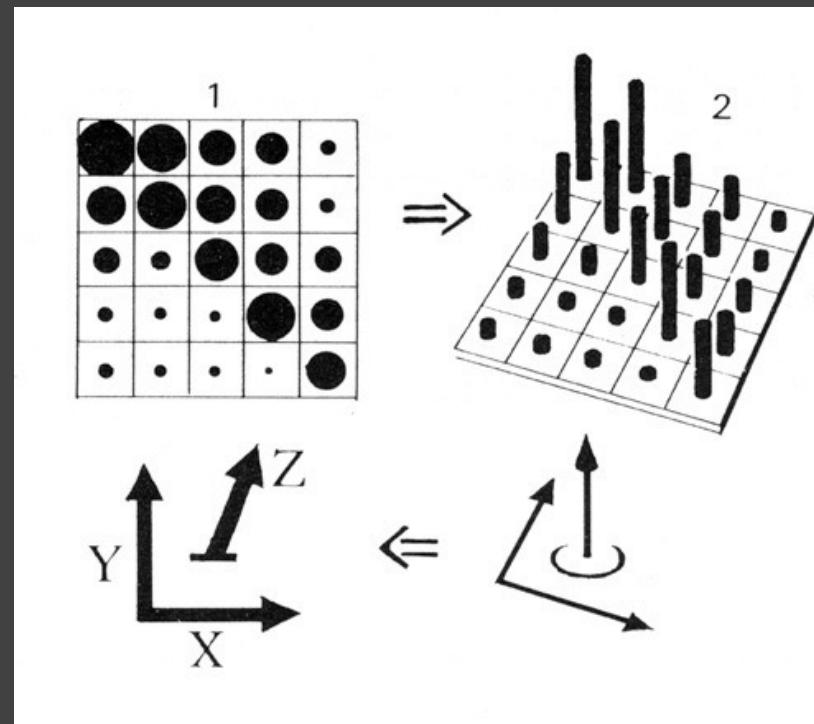


Three Variables

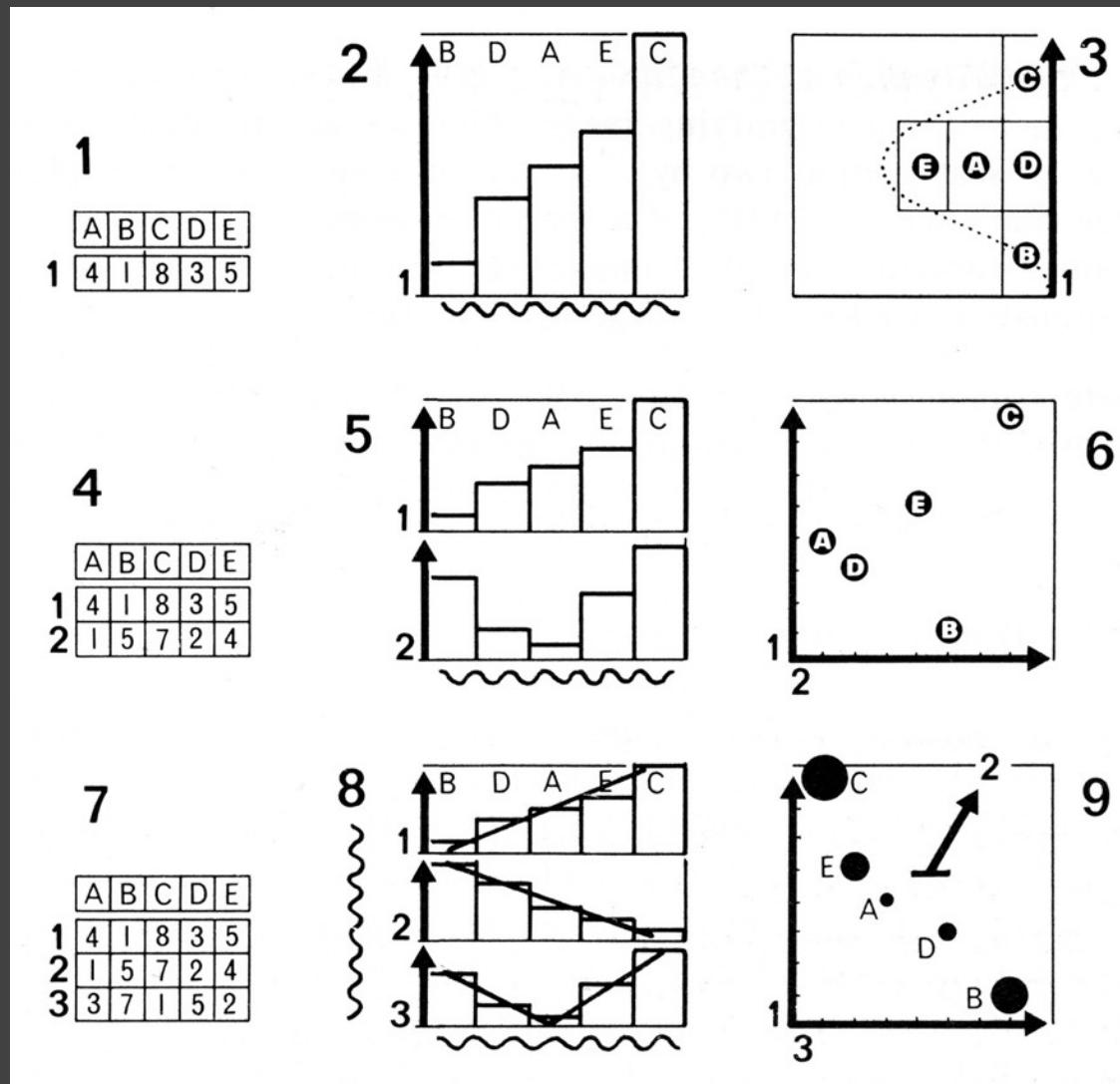
Two variables [x,y] can map to points

Scatterplots, maps, ...

Third variable [z] must use
Color, size, shape, ...



Large Design Space



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

Multidimensional Data

How many variables can
be depicted in an image?

	A	B	C
1			
2			
3			
4			
5			
6			
7			
8			

Multidimensional Data

How many variables can be depicted in an image?

"With up to three rows, a data table can be constructed directly as a single image ... However, an image has only three dimensions. And this barrier is impassible." - Bertin

	A	B	C
1			
2			
3			
4			
5			
6			
7			
8			

Multidimensional Data

Visual Encoding Variables

Position (X)

Position (Y)

Size

Value

Texture

Color

Orientation

Shape

~8 dimensions?

		LES VARIABLES DE L'IMAGE				
		POINTS	LIGNES	ZONES		
XY 2 DIMENSIONS DU PLAN	Z	x	x	x	12	12
	TAILLE	■	■	■	12	12
	VALEUR	■	■	■	12	12
		LES VARIABLES DE SÉPARATION DES IMAGES				
GRAIN		■■■	■■■	■■■	12	12
COULEUR		■	■	■	12	12
ORIENTATION		■	■	■	12	12
FORME		■	▲	●	12	12

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

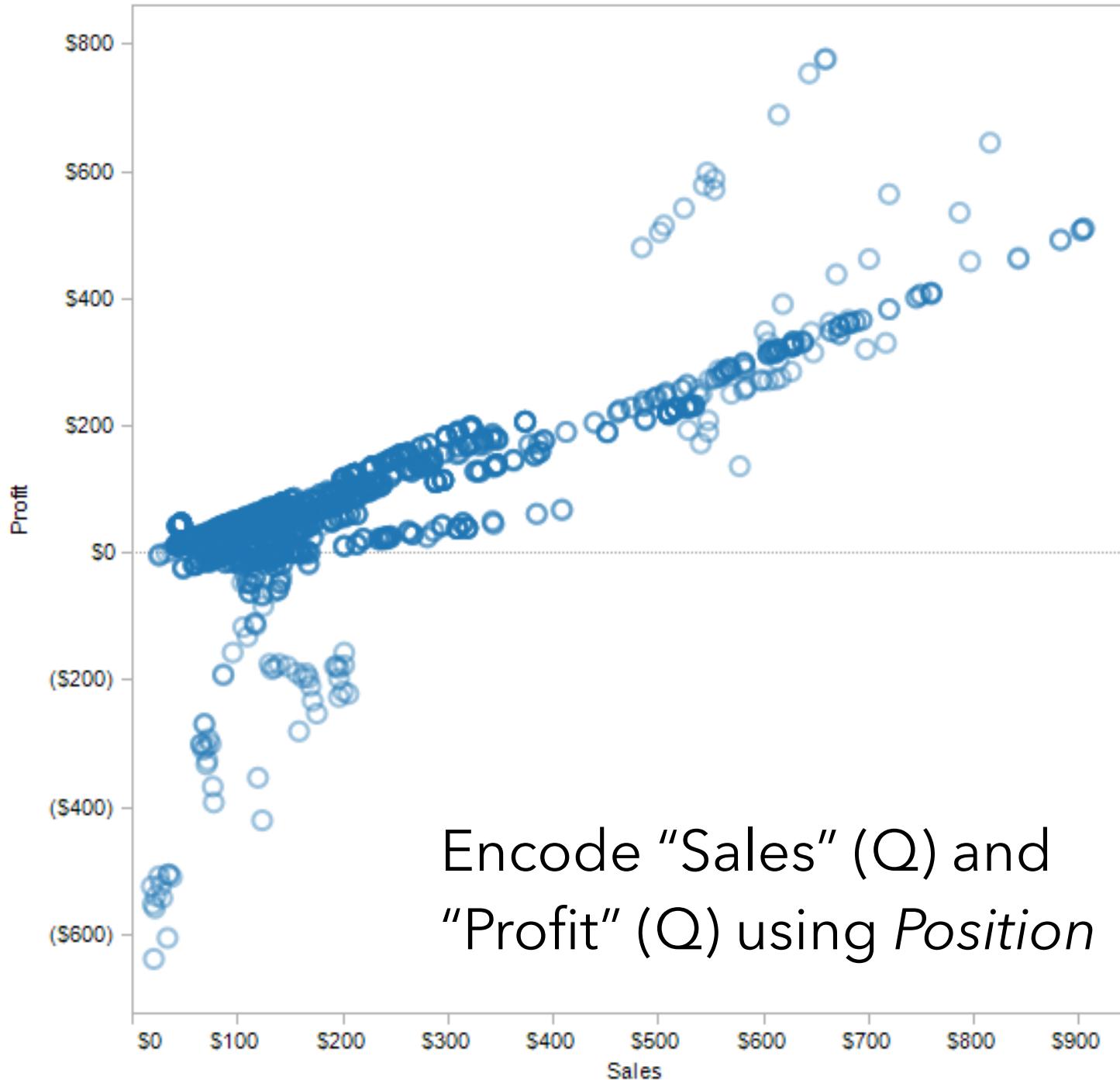
Color



Size



Level of Detail



Filters

YEAR(Date): 2010

Marks

x+ Automatic

Shape Label

Color ▾ Product Type

Size Level of Detail

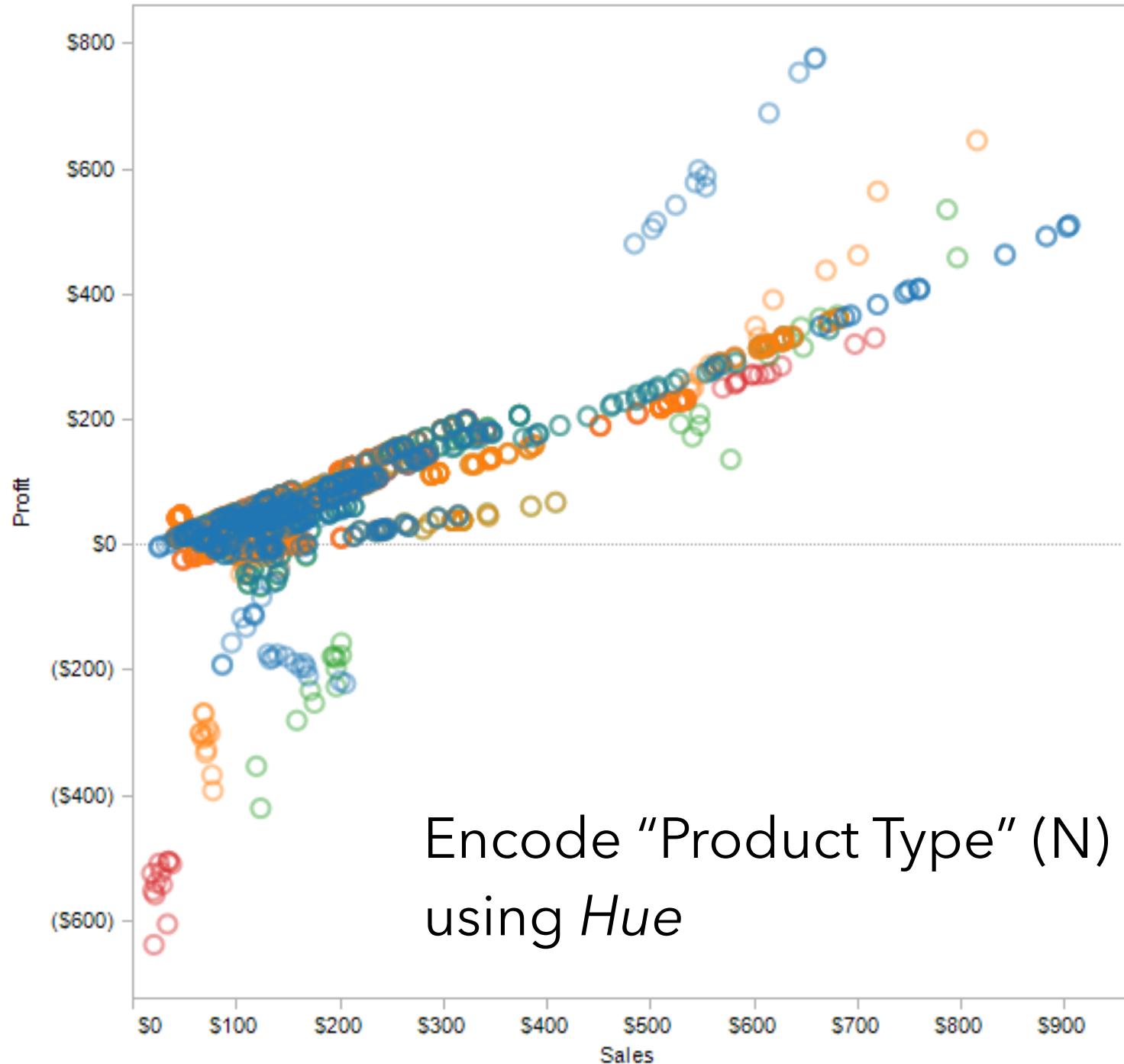
Product Type

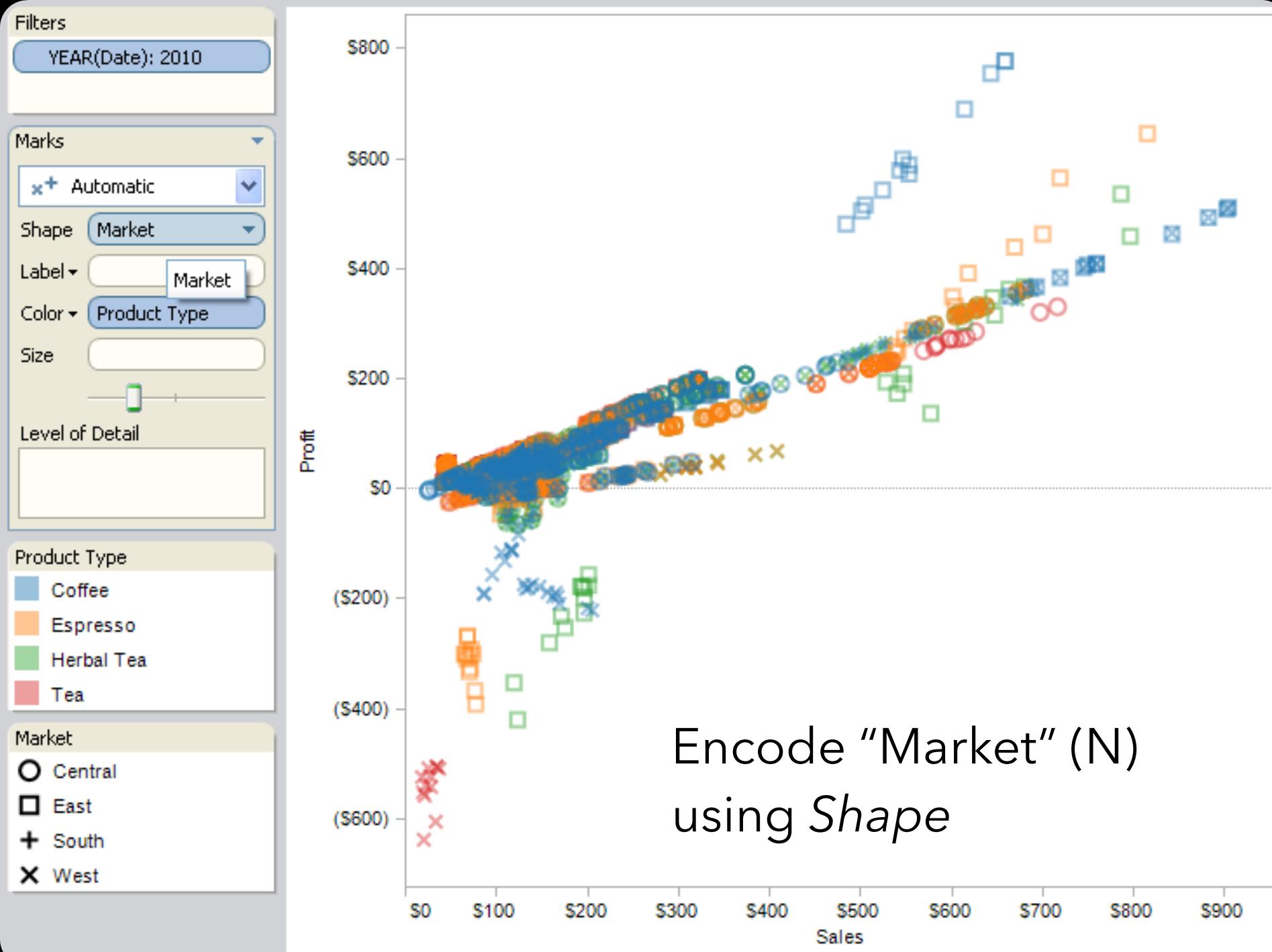
Coffee

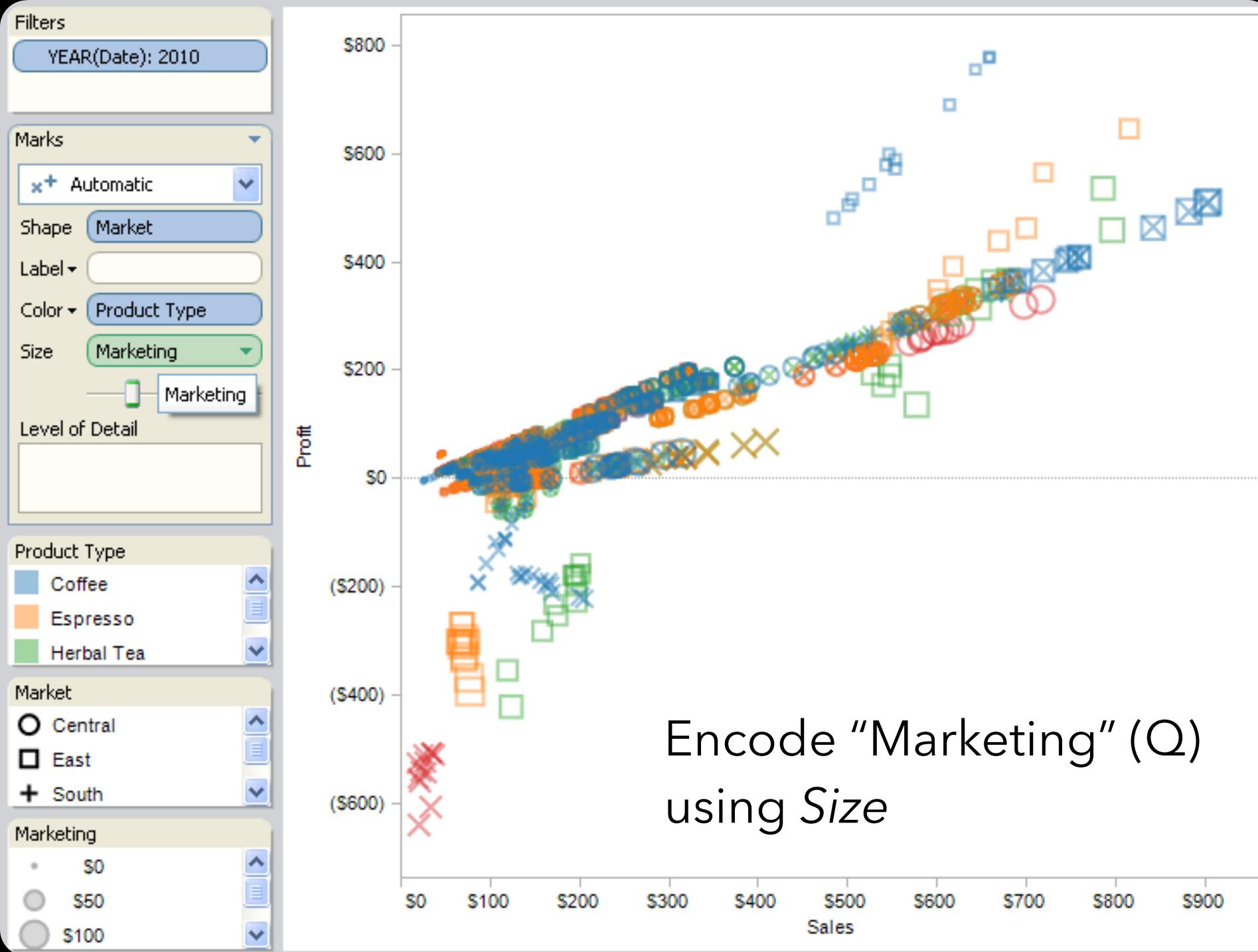
Espresso

Herbal Tea

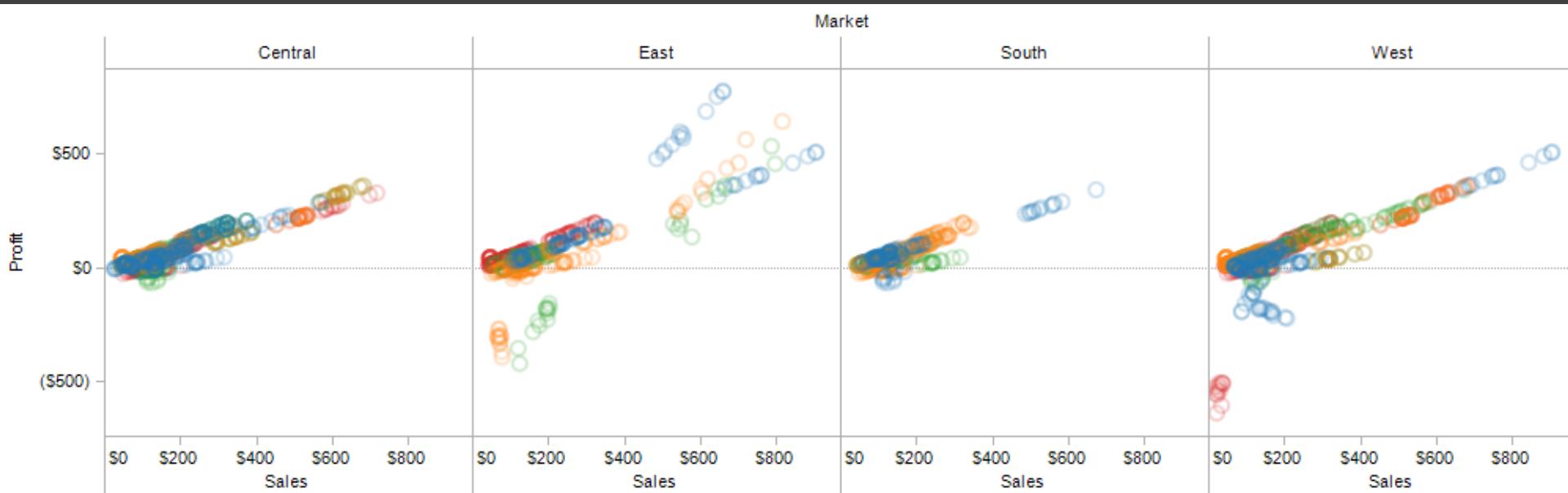
Tea







Trellis Plots



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

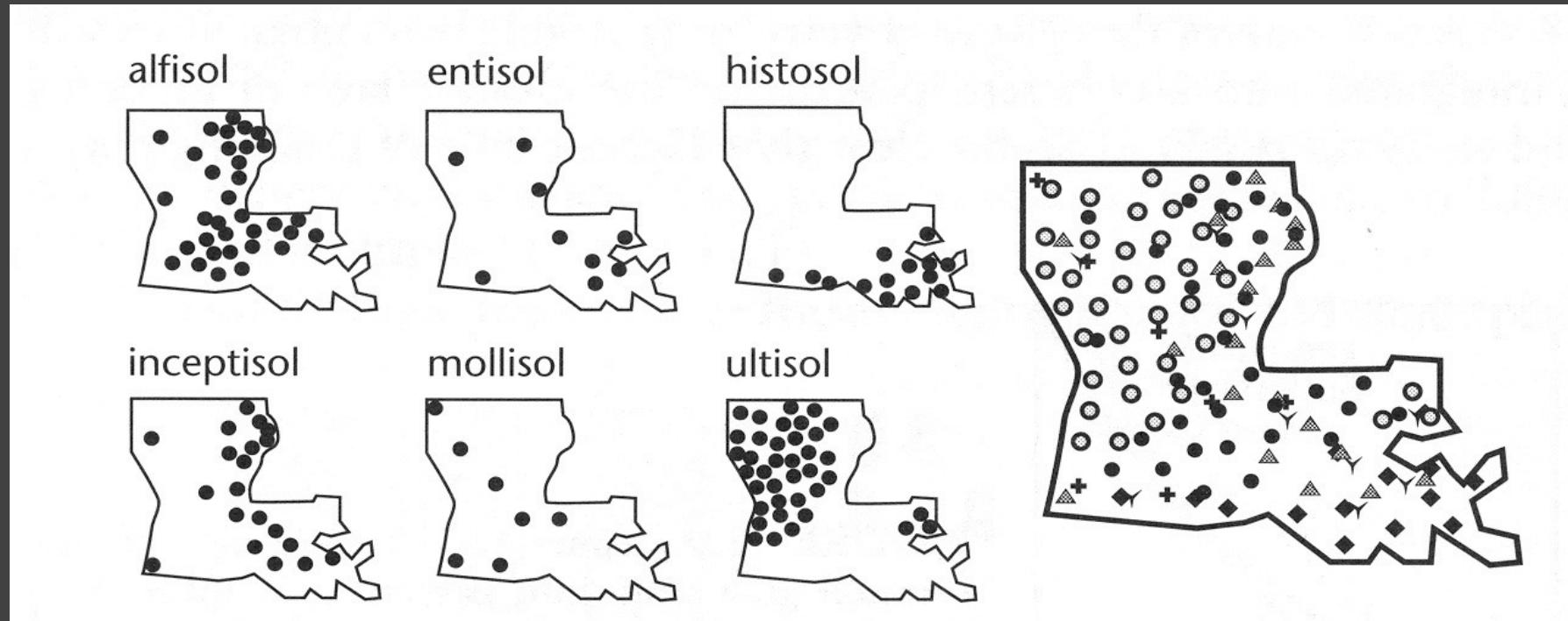
Typically nominal or ordinal variables are used as dimensions for subdivision.

Small Multiples



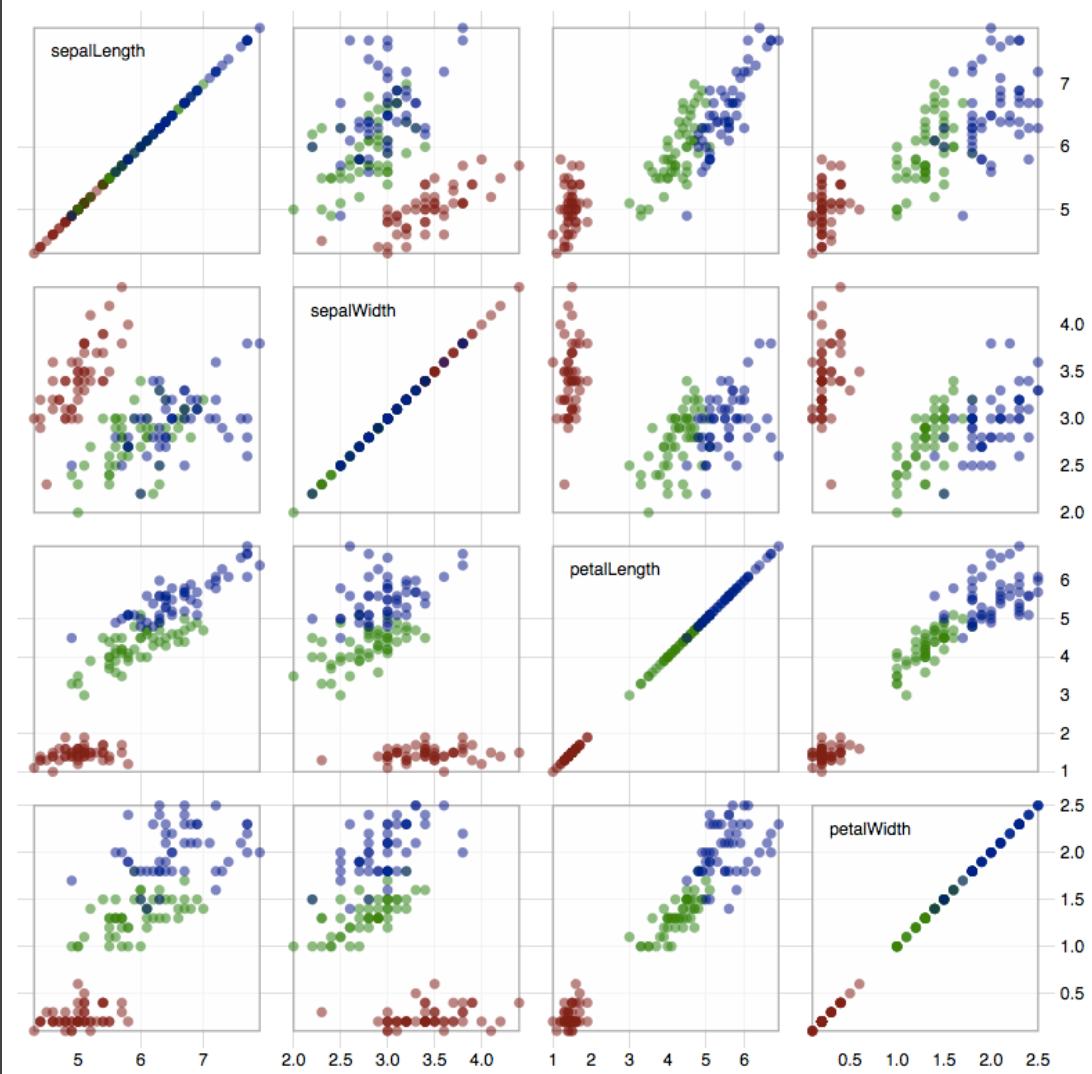
[MacEachren 95, Figure 2.11, p. 38]

Small Multiples

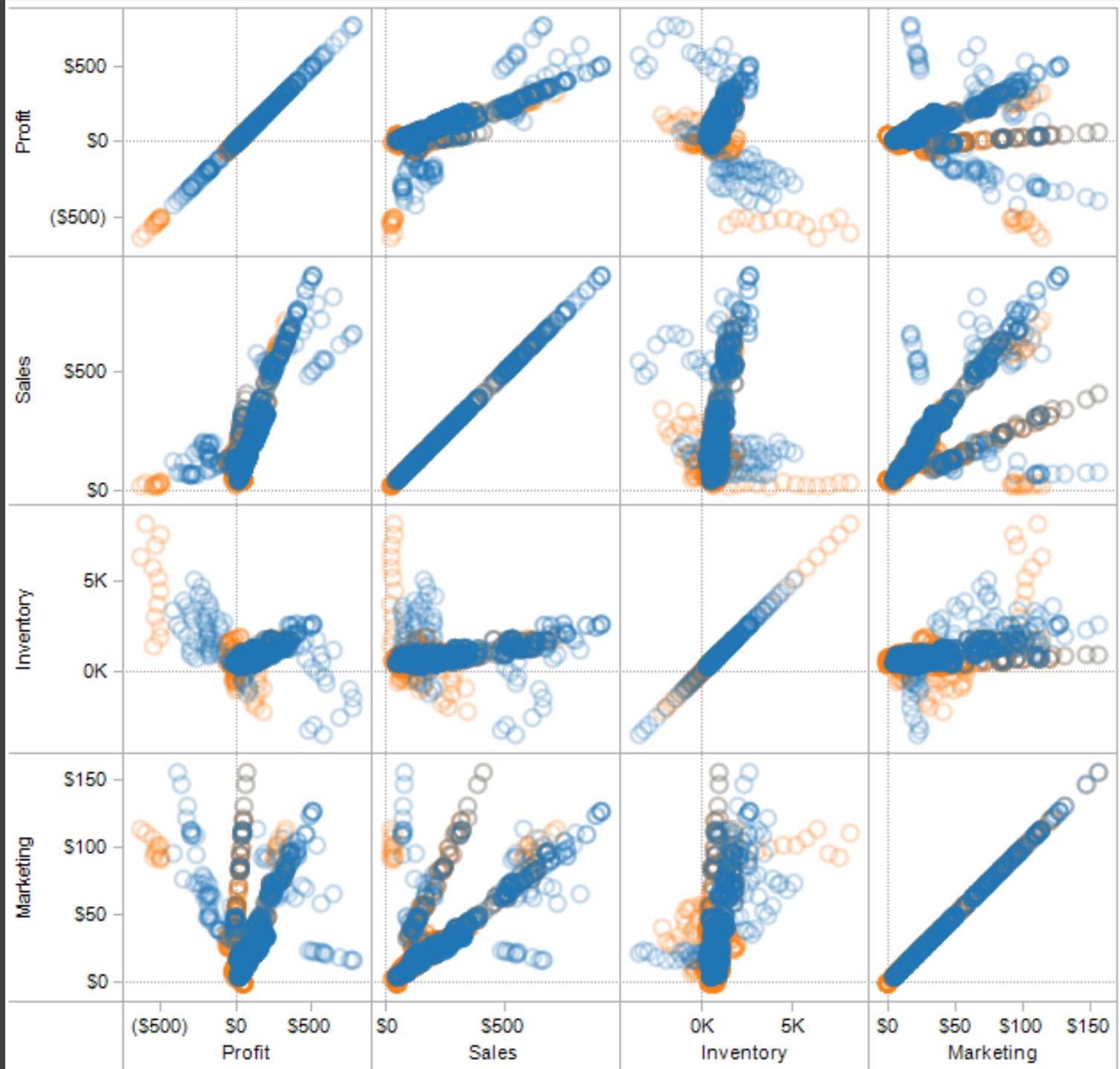


[MacEachren 95, Figure 2.11, p. 38]

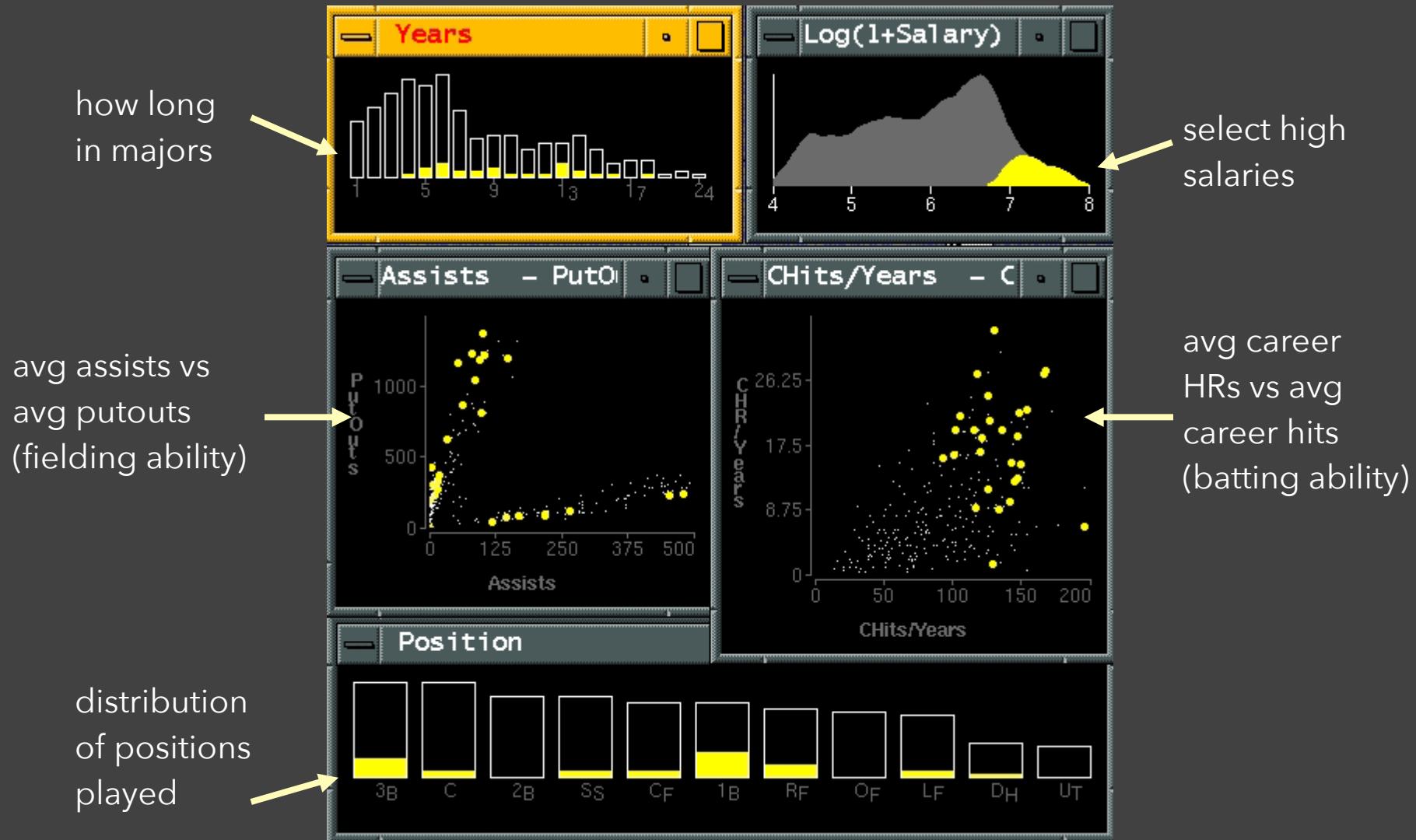
Scatterplot Matrix (SPLOM)



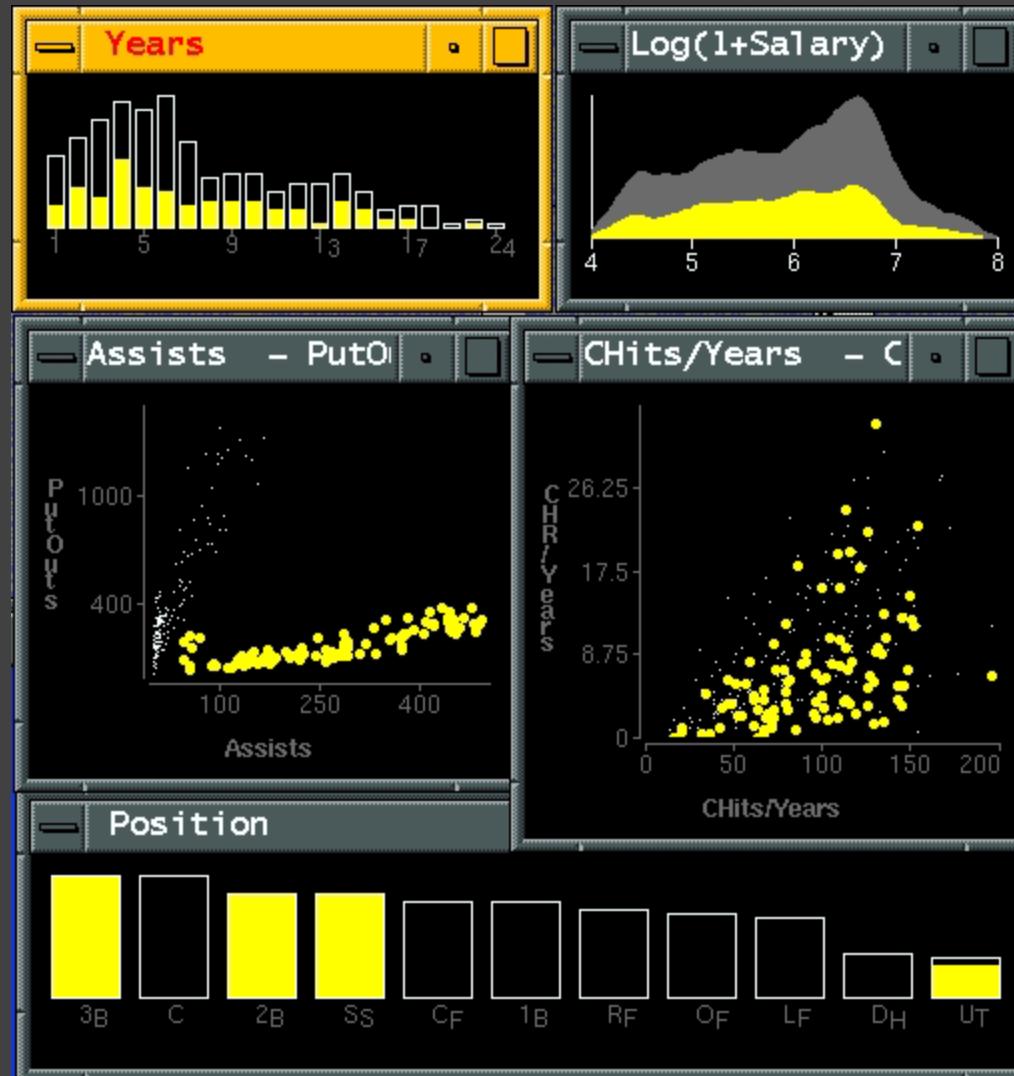
Scatter plots
for pairwise
comparison
of each data
dimension.



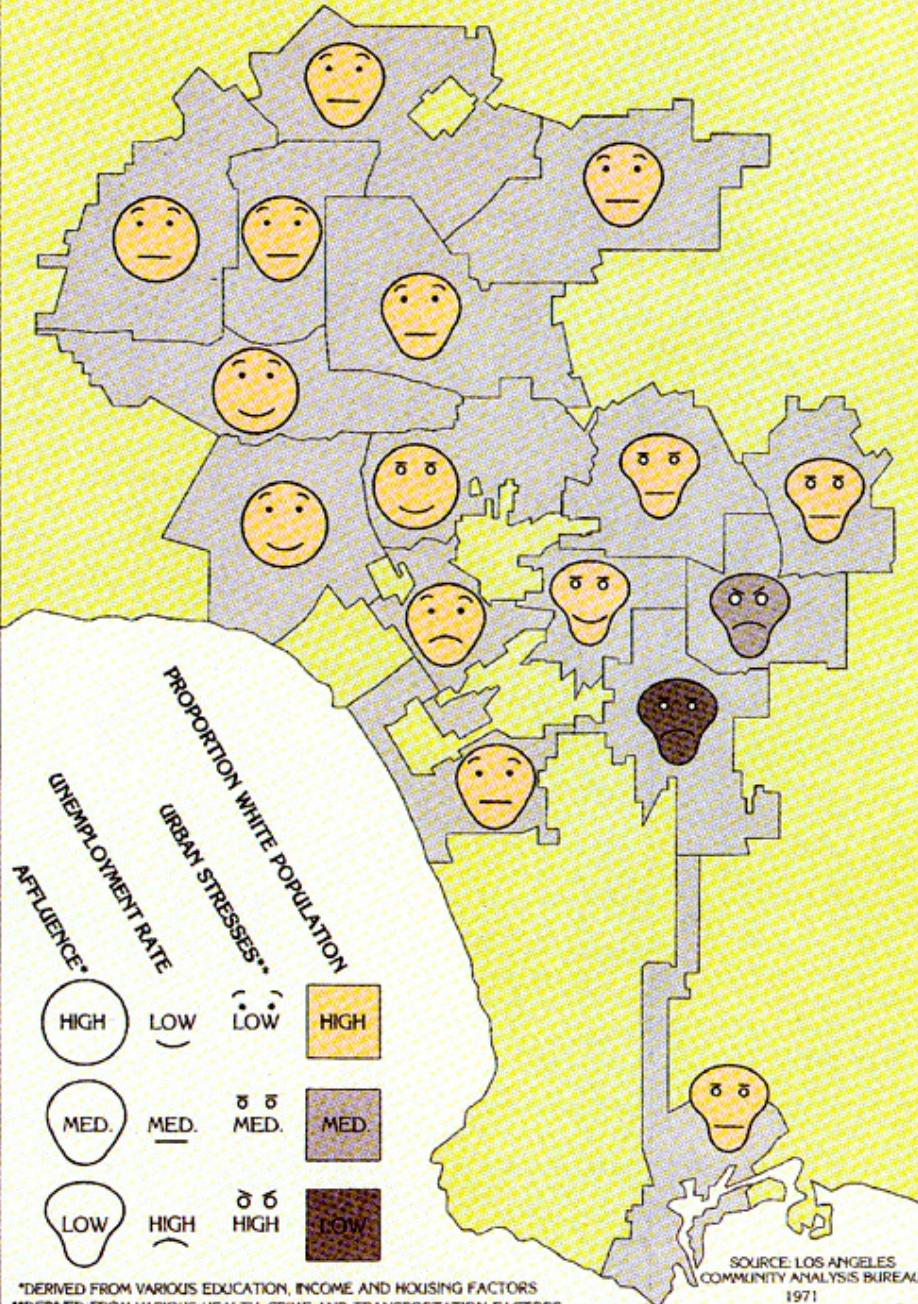
Multiple Coordinated Views



Linking Assists to Position



Life in Los Angeles



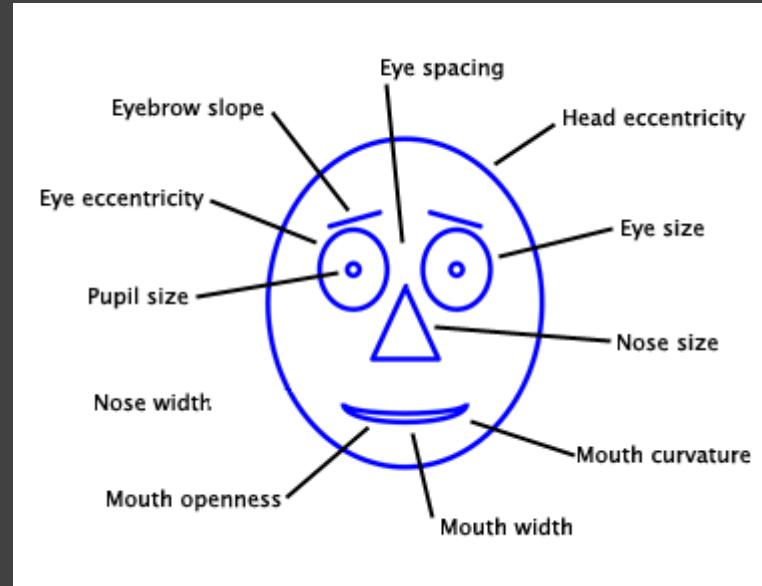
Chernoff Faces

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

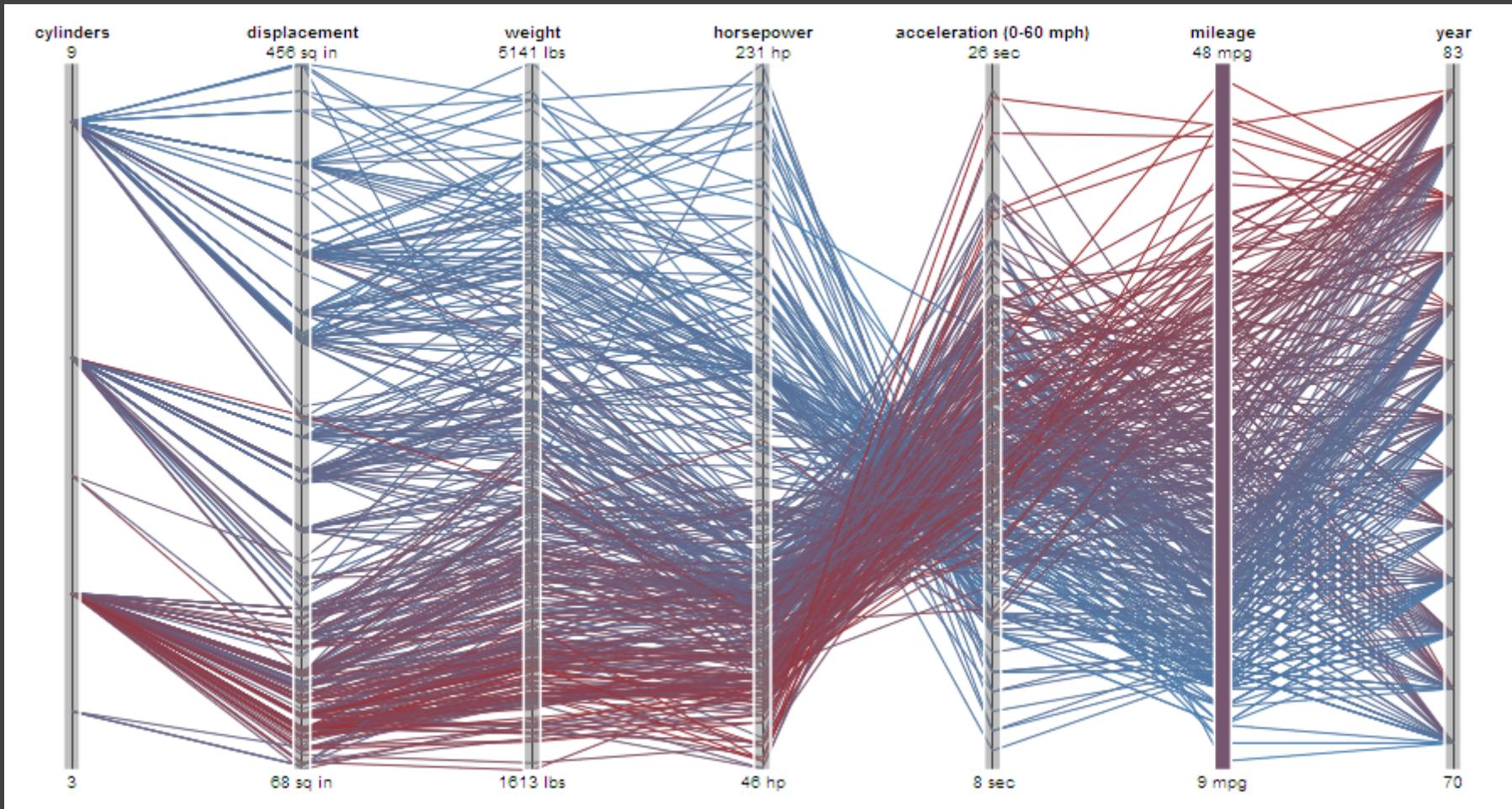
Reduce the problem space

Use interaction to generate *relevant* views

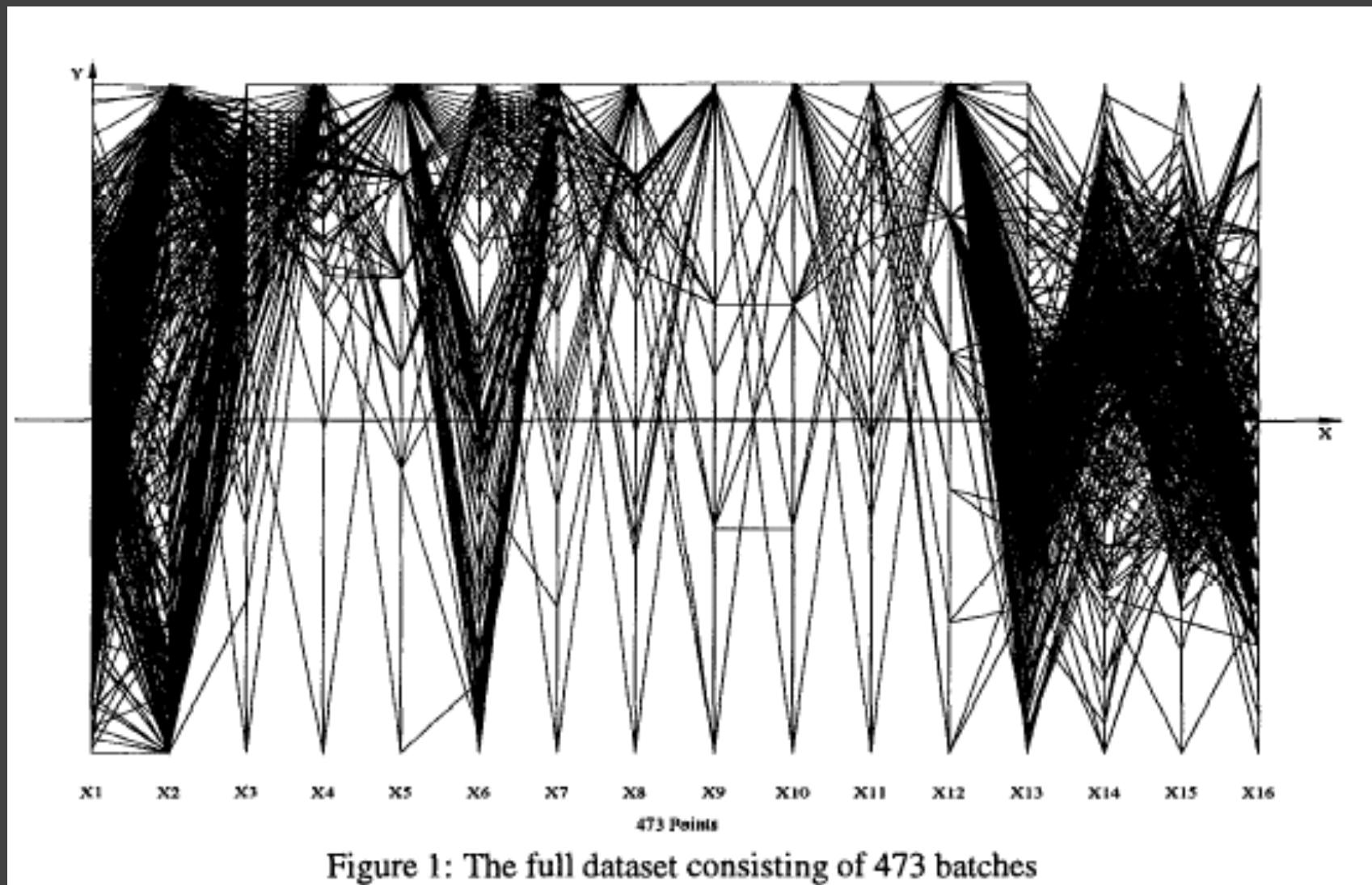
Rarely does a single visualization answer all questions. Instead, the ability to generate appropriate visualizations quickly is key.

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)

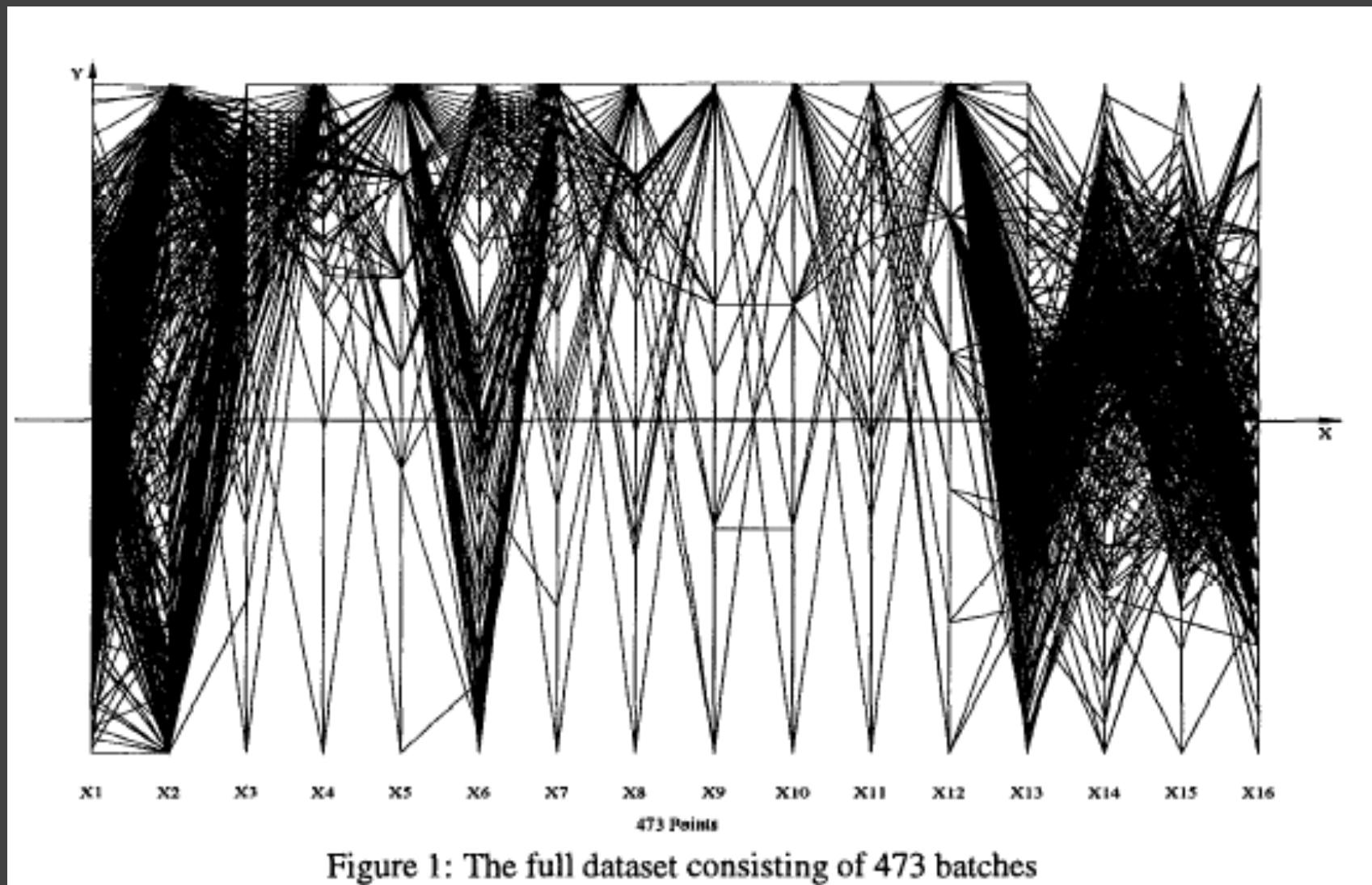
X3-12: 10 types of defects (0 defects shown at top)

X13-16: 4 physical parameters

Objective:

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

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 X_1 and X_2

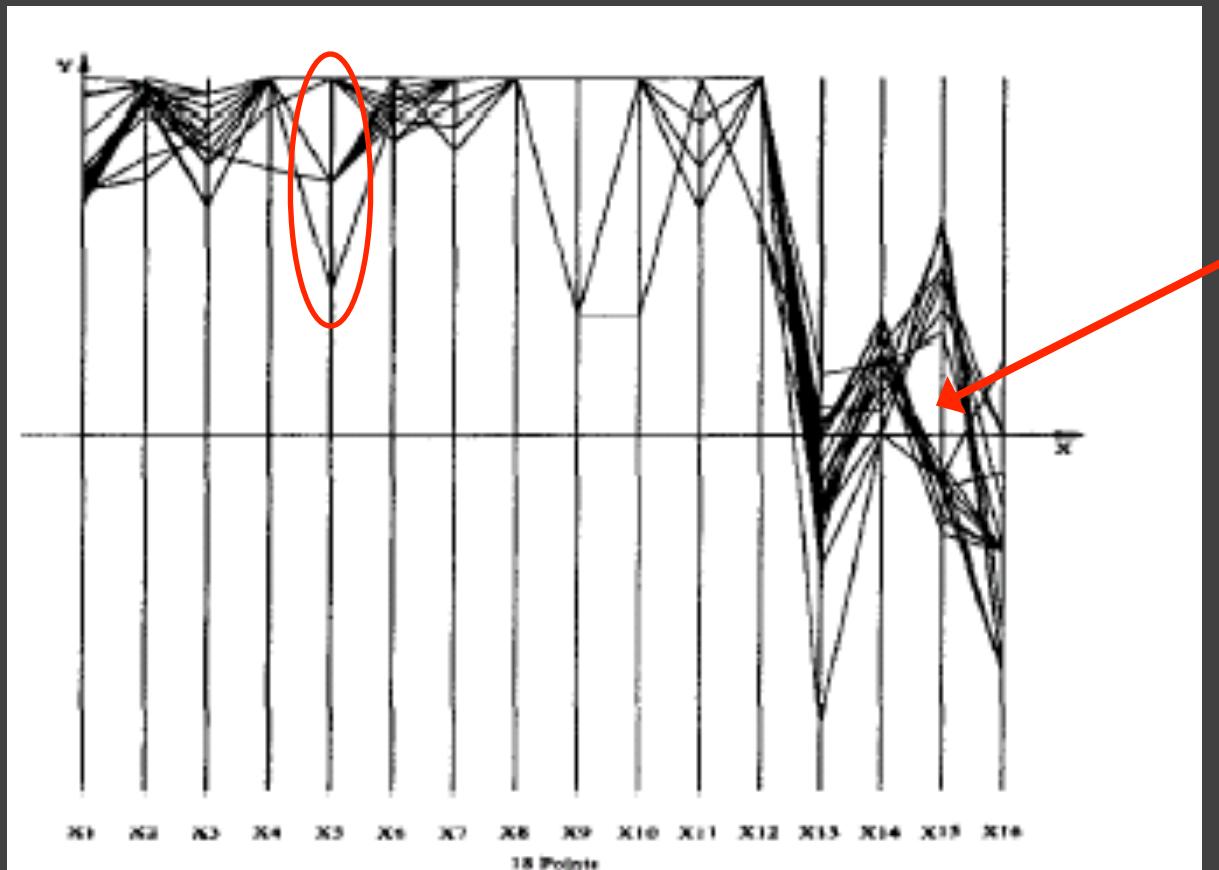
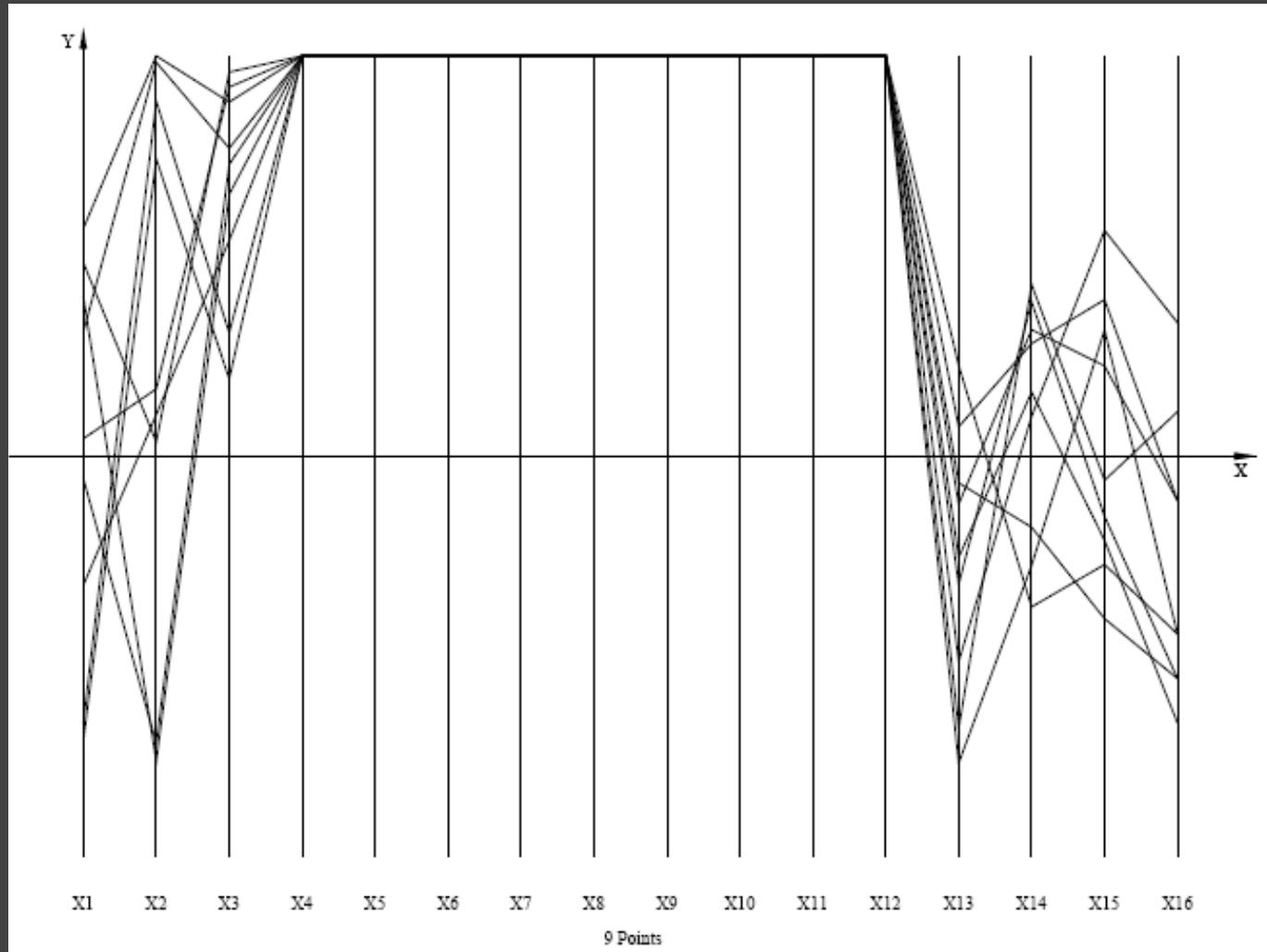


Figure 2: The batches high in Yield, X_1 , and Quality, X_2 .

Look for batches with *nearly* zero defects (9/10)

Most of these have low yields -> defects OK.



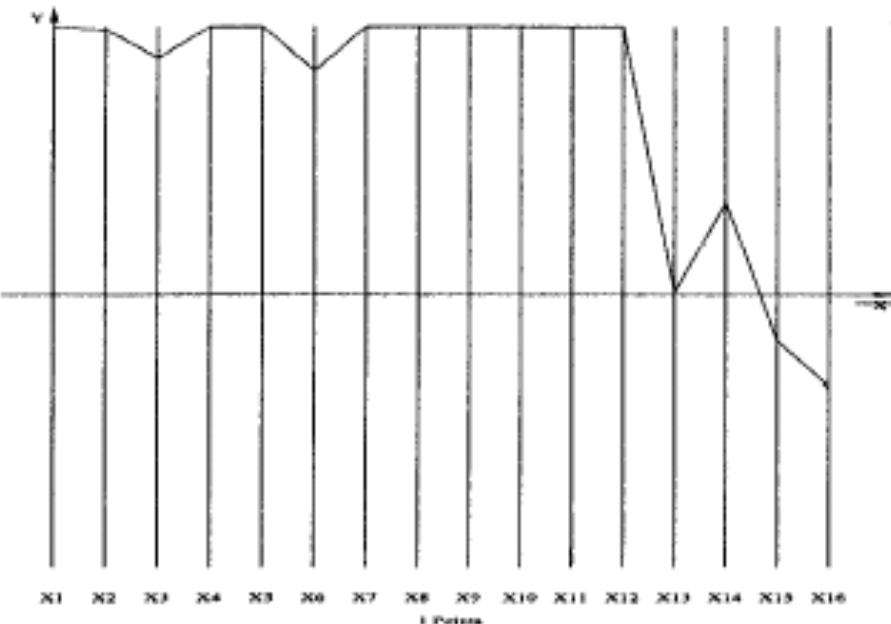


Figure 5: The best batch. Highest in Yield, X_1 , and very high in Quality, X_2 .

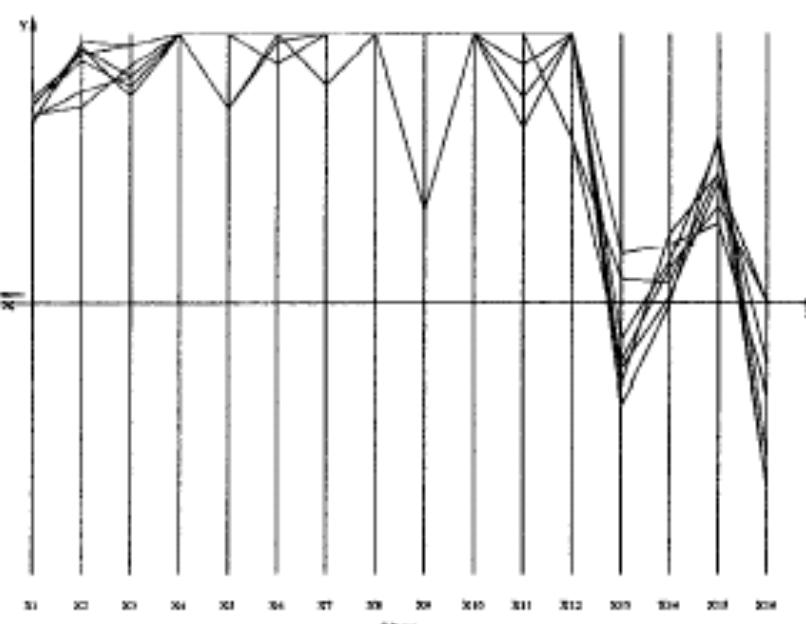
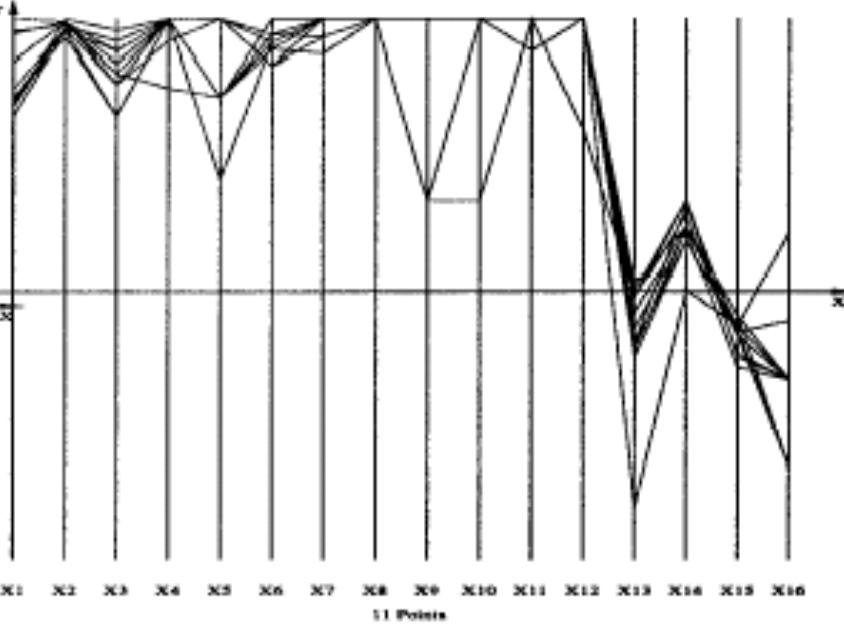
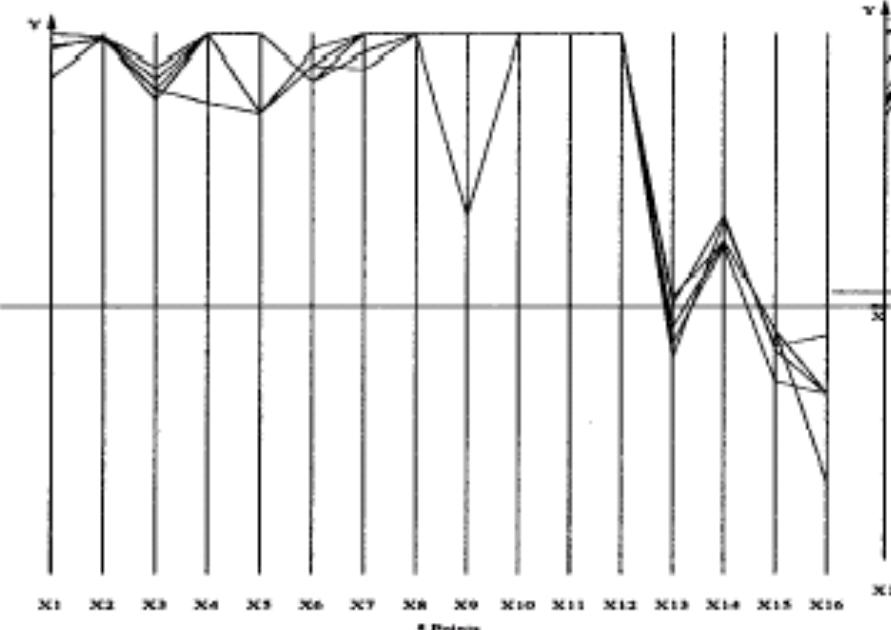
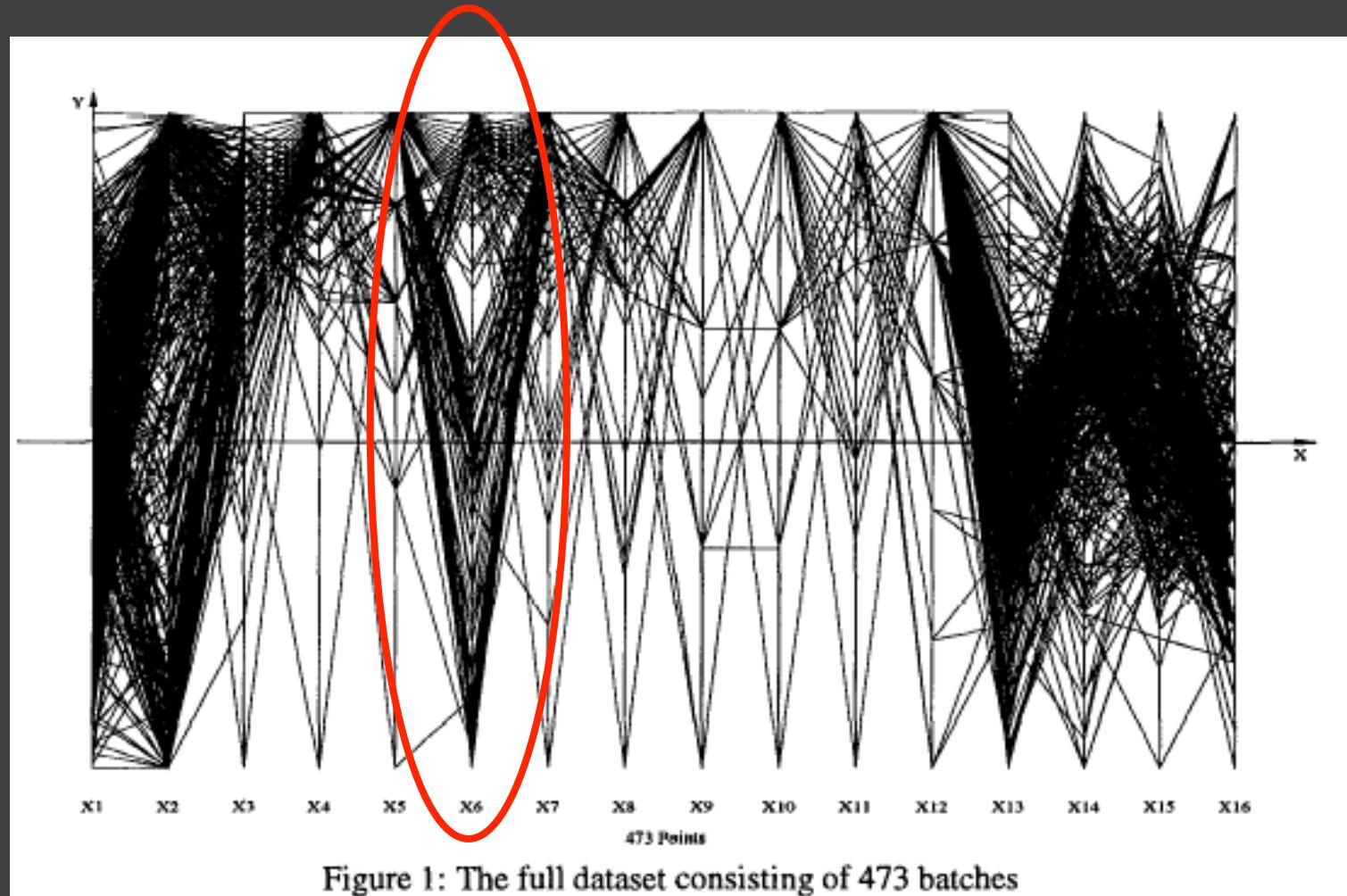


Figure 7: Upper range of split in X_{15}

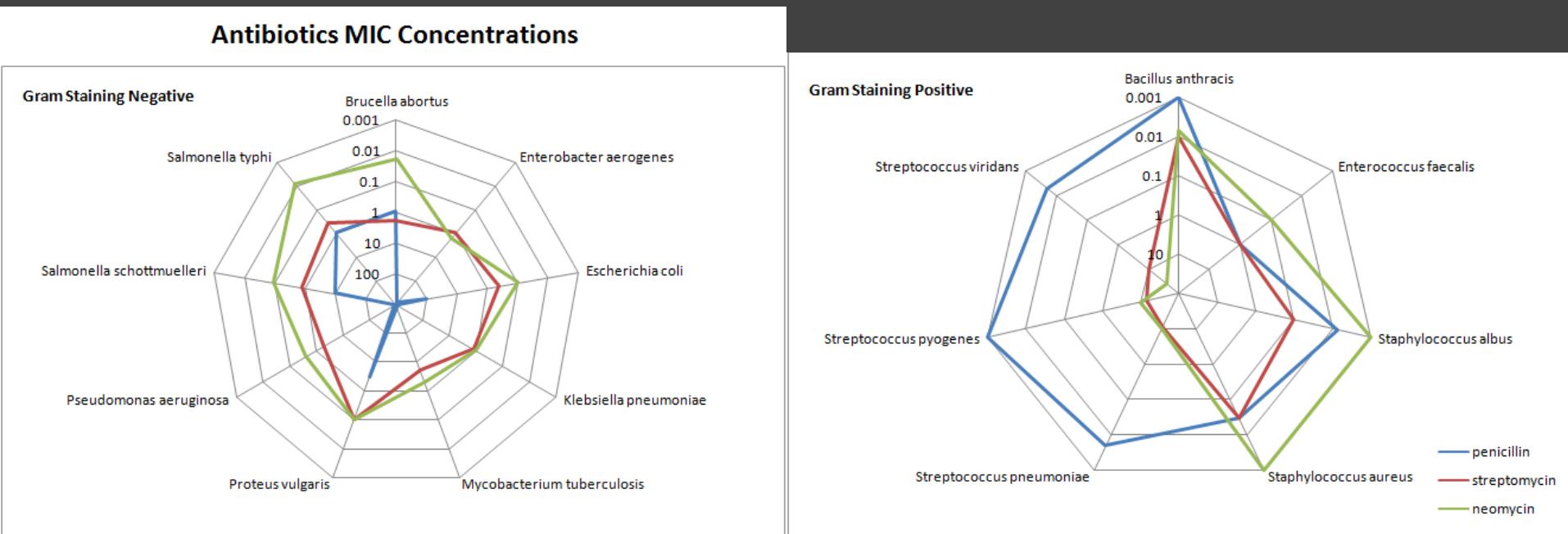


Notice that X6 behaves differently.

Allow 2 defects, including X6 -> best batches



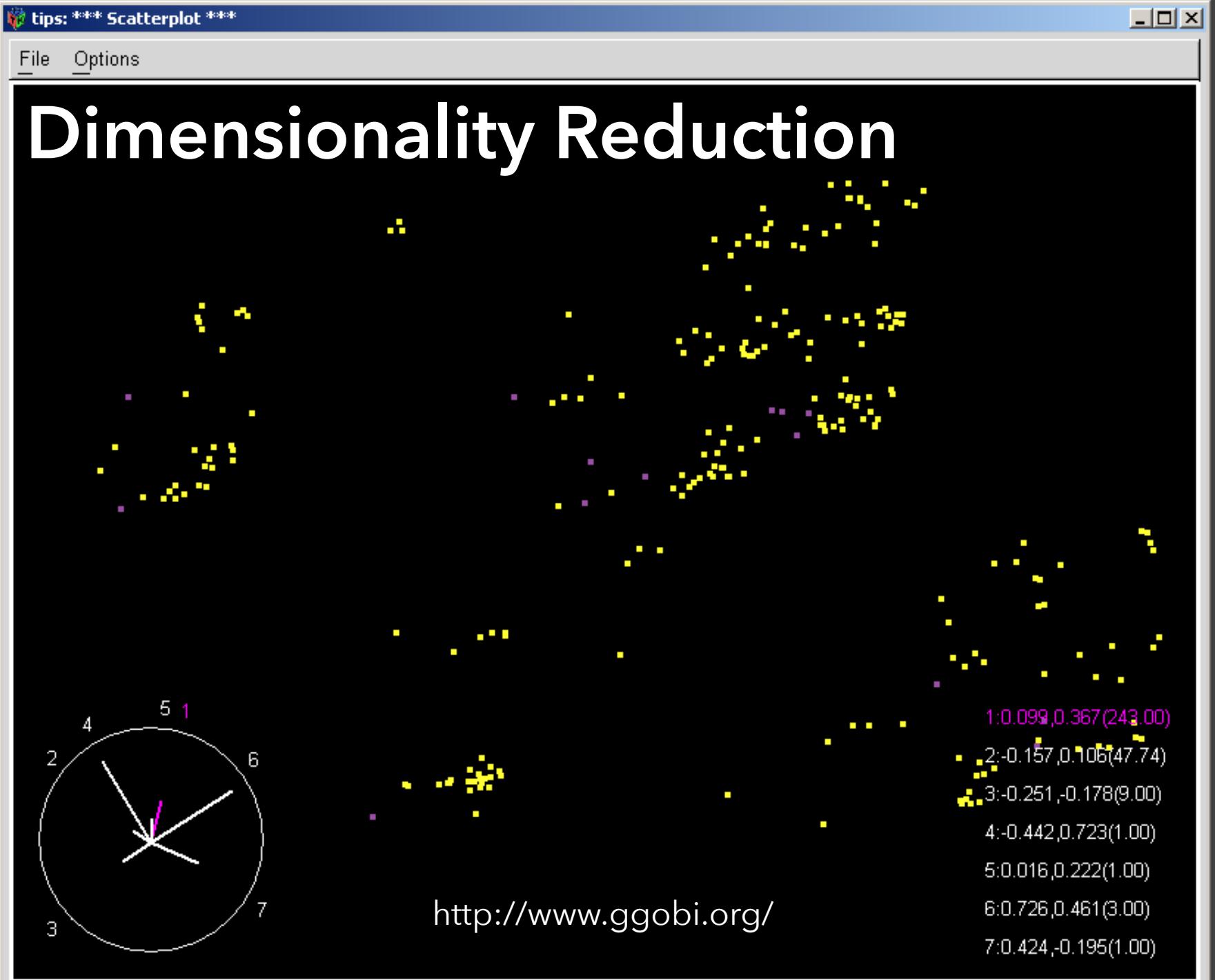
Radar Plot / Star Graph



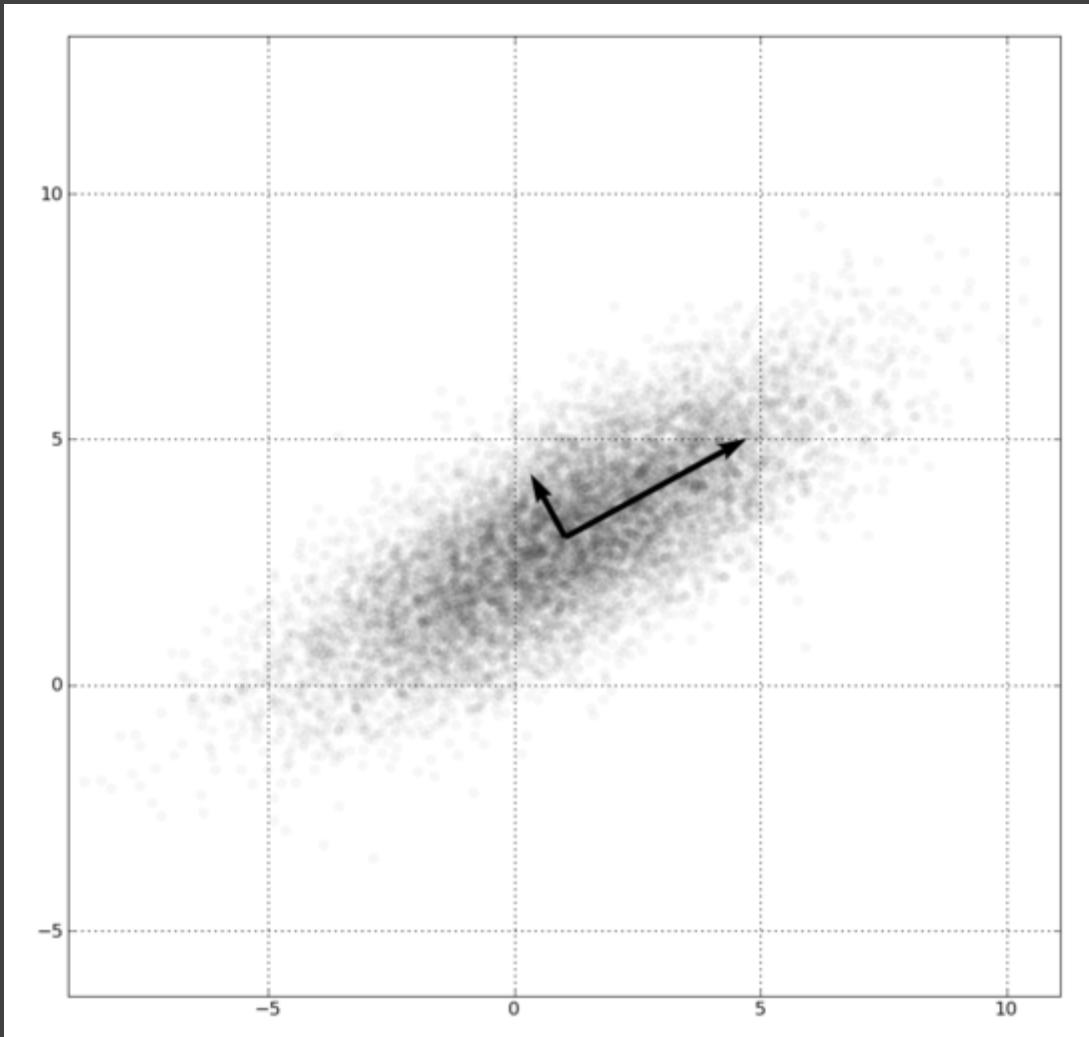
"Parallel" dimensions in polar coordinate space

Best if same units apply to each axis

Dimensionality Reduction

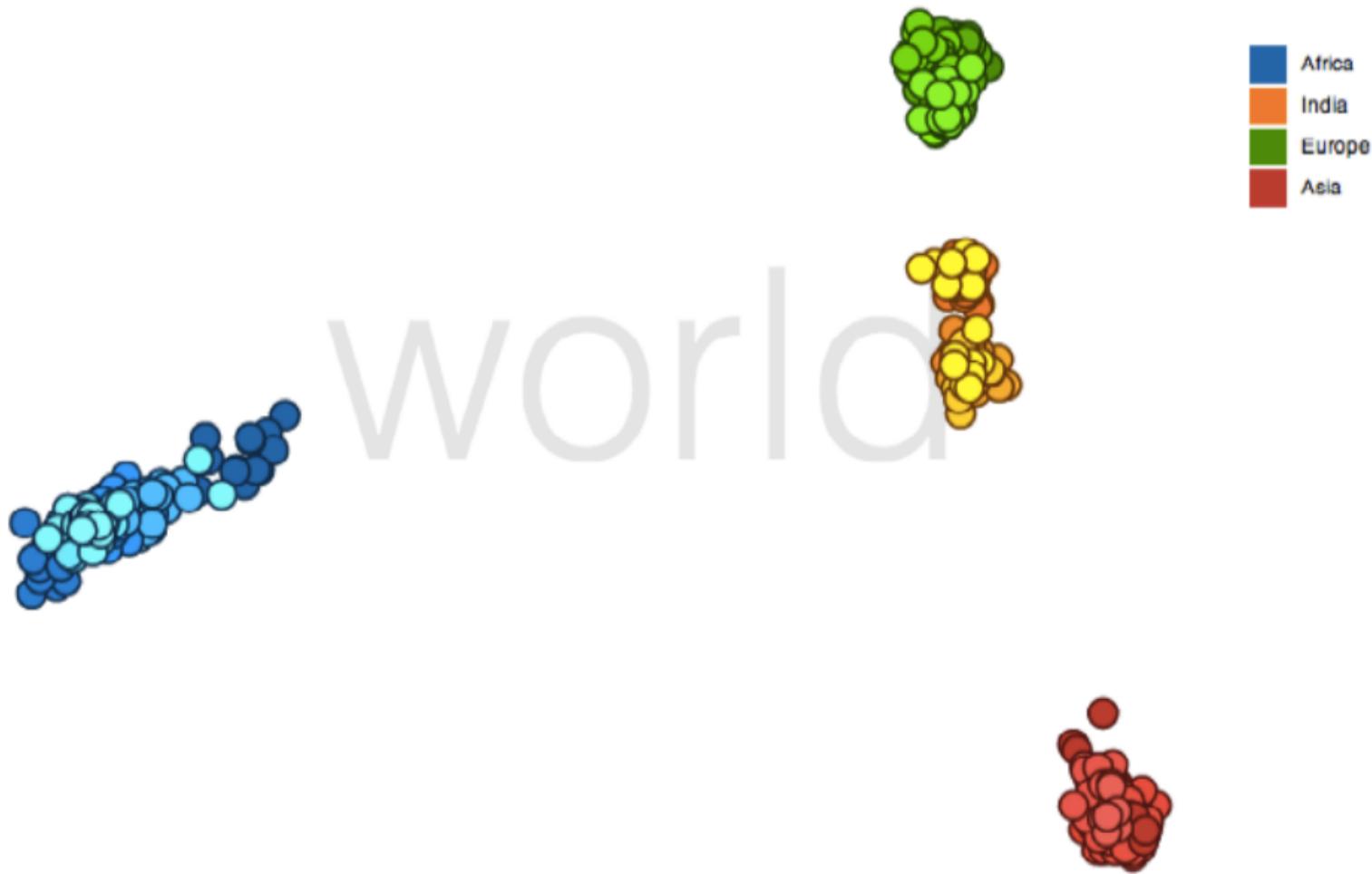


Principal Components Analysis



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

PCA on Genetic Sequences



Many Reduction Techniques!

Principal Components Analysis (PCA)

Multidimensional Scaling (MDS)

Locally Linear Embedding (LLE)

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

Isomap

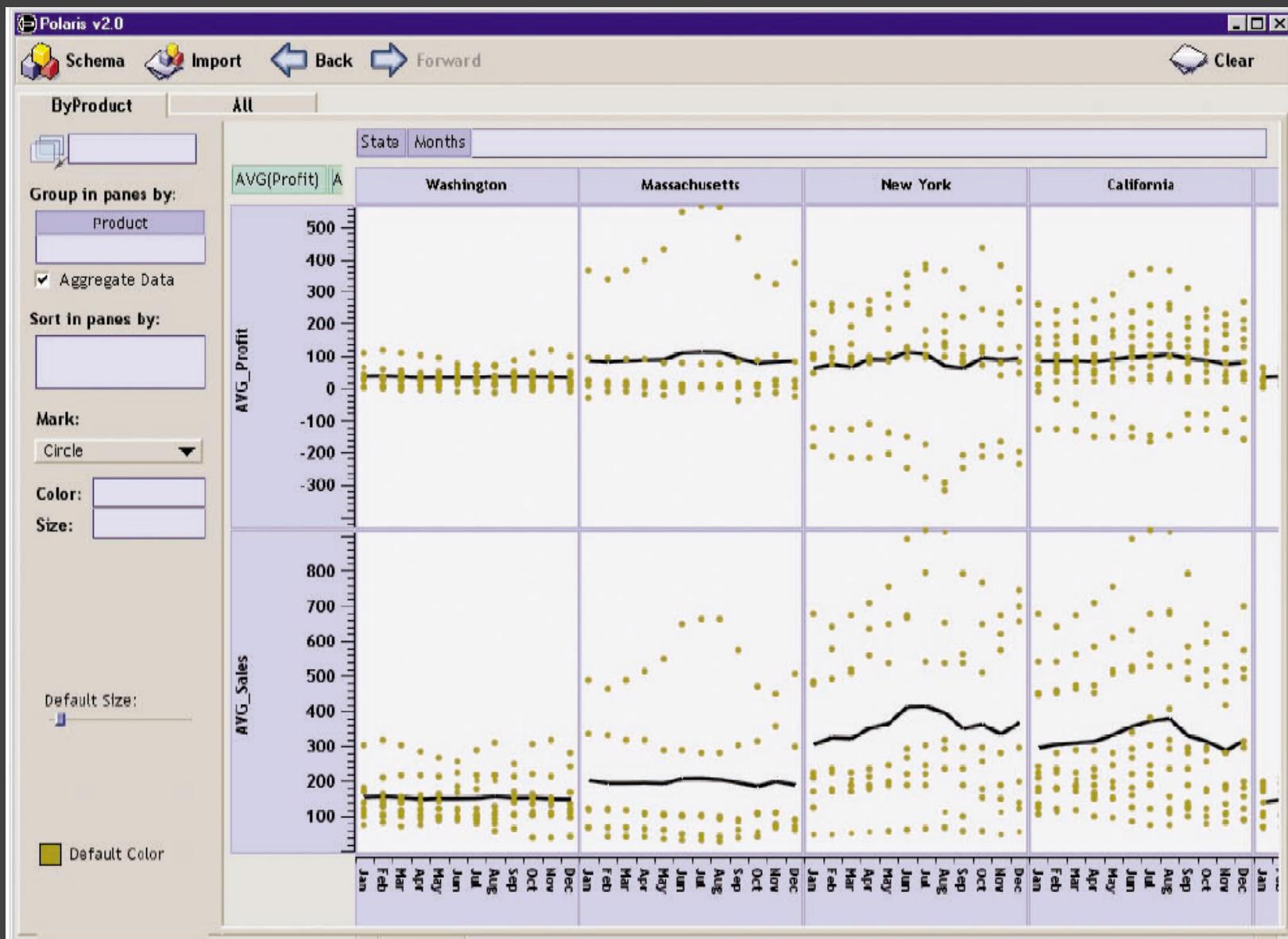
Auto-Encoder Neural Networks

Topological methods

...

Tableau / Polaris

Polaris [Stolte et al.]



Tableau

Encodings

Data Display

Data Model

The screenshot shows the Tableau interface with several panels:

- Schema:** Shows the connection to "congress.csv Connection".
- Dimensions:** Includes "Party", "Year", "Candidate", "Candidate ID", "General Elec Status", "Incumbent/Challenger/Open-Seal", "# Party", "Party Desig", "Primary Elec Status", "Runoff Elec Status", "Spec Elec Status", "State Code", "# Year", and "Measure Names".
- Measures:** Includes "# District", "# General Elec Pct", "# Total Receipts", and "# Measure Values".
- Groups:** An empty panel.
- Encodings:** A central panel where data is visualized. It shows a bar chart with three groups (1, 2, 3) over years 1996, 1998, 2000, and 2002. The Y-axis is "SUM(Total Receipts)" from 0M to 550M. The X-axis shows years grouped by group. The legend indicates blue for group 1, orange for group 2, and green for group 3. The chart shows high values for group 1 in 2000 and group 2 in 2000, with very low values for group 3.
- Columns:** Set to "Party" and "Year".
- Rows:** Set to "SUM(Total Receipts)".
- Filters:** An empty panel.
- Level of Detail:** An empty panel.
- Mark:** Set to "Automatic".
- Text:** An empty panel.
- Color:** Set to "Party".
- Size:** An empty panel.
- Legend:** Shows the color mapping for groups 1, 2, and 3.
- Size:** An empty panel.

The overall layout includes a top menu bar with File, Edit, View, Format, Data, Analysis, Table, Bookmark, Window, Help, and a toolbar with various icons. The bottom of the screen shows navigation buttons and a status bar indicating "Sheet 1 /".

Group	Year	Party	Sum(Total Receipts)
1	1996	Party 1	~350M
1	1998	Party 1	~360M
1	2000	Party 1	~530M
1	2002	Party 1	~480M
2	1996	Party 2	~430M
2	1998	Party 2	~410M
2	2000	Party 2	~520M
2	2002	Party 2	~490M
3	1996	Party 3	~10M
3	1998	Party 3	~10M
3	2000	Party 3	~10M
3	2002	Party 3	~10M

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 (/)

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

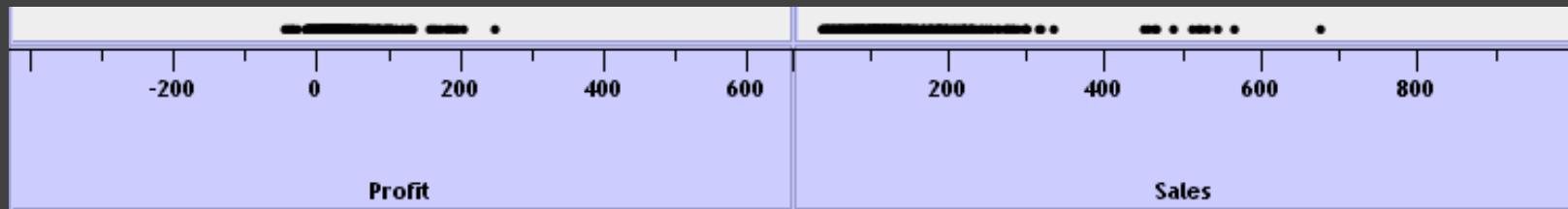
Ordered union of set interpretations

Quarter + Product Type

$$\begin{aligned} &= \{(Qtr1), (Qtr2), (Qtr3), (Qtr4)\} + \{(Coffee), (Espresso)\} \\ &= \{(Qtr1), (Qtr2), (Qtr3), (Qtr4), (Coffee), (Espresso)\} \end{aligned}$$

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

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



Cross (x) Operator

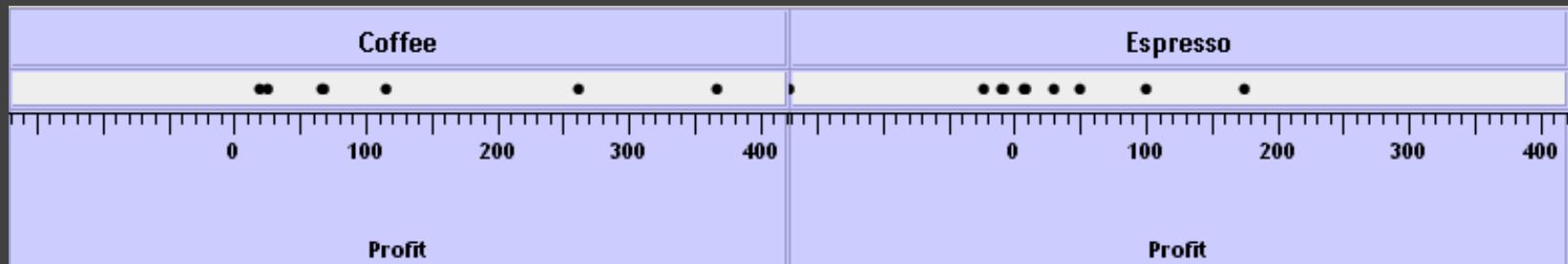
Cross-product of set interpretations

Quarter x Product Type =

$\{(Qtr1, \text{Coffee}), (Qtr1, \text{Tea}), (Qtr2, \text{Coffee}), (Qtr2, \text{Tea}), (Qtr3, \text{Coffee}), (Qtr3, \text{Tea}), (Qtr4, \text{Coffee}), (Qtr4, \text{Tea})\}$

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

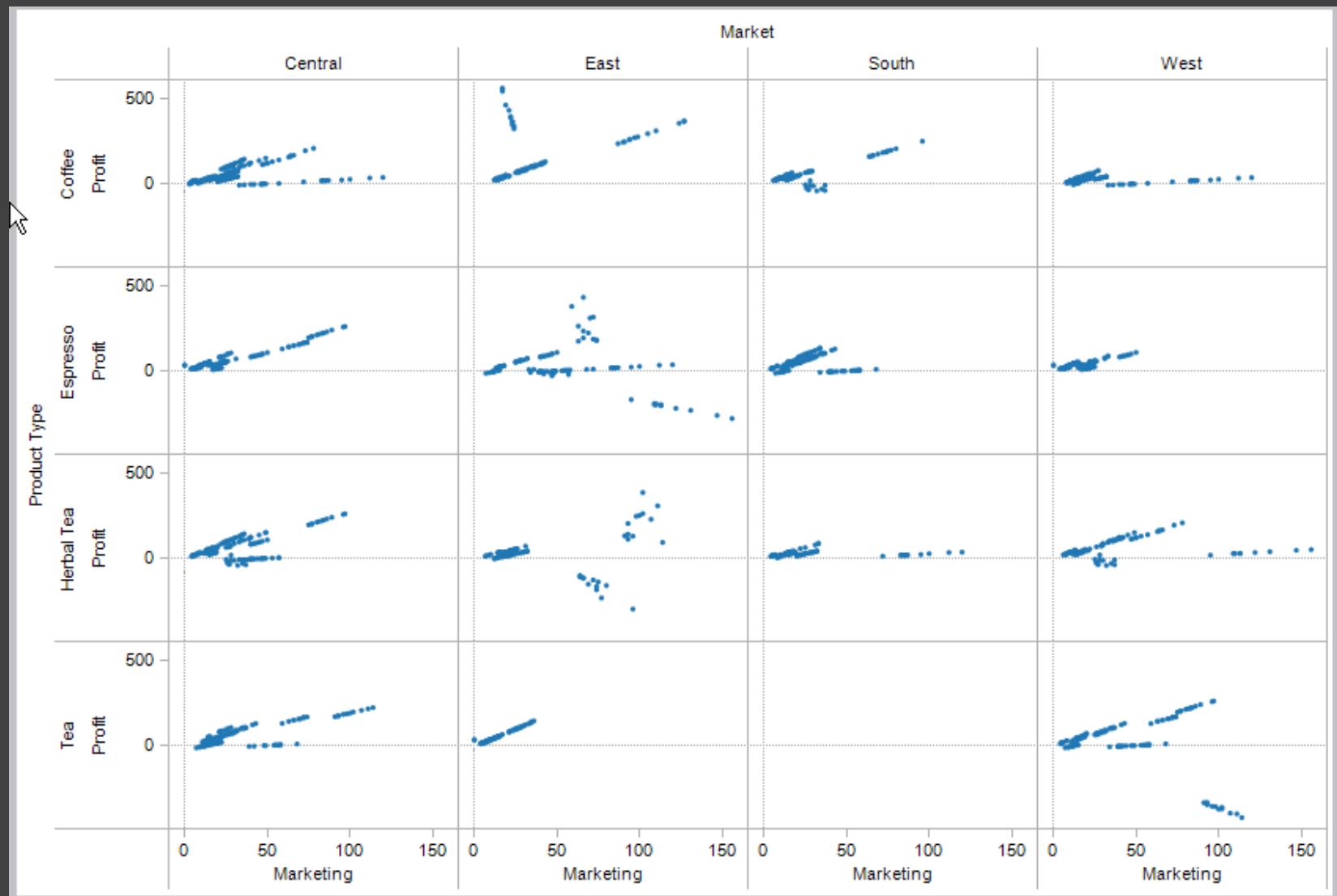
Note that this specifies operands *NOT* operators!

Operators are inferred by data type (O, Q)

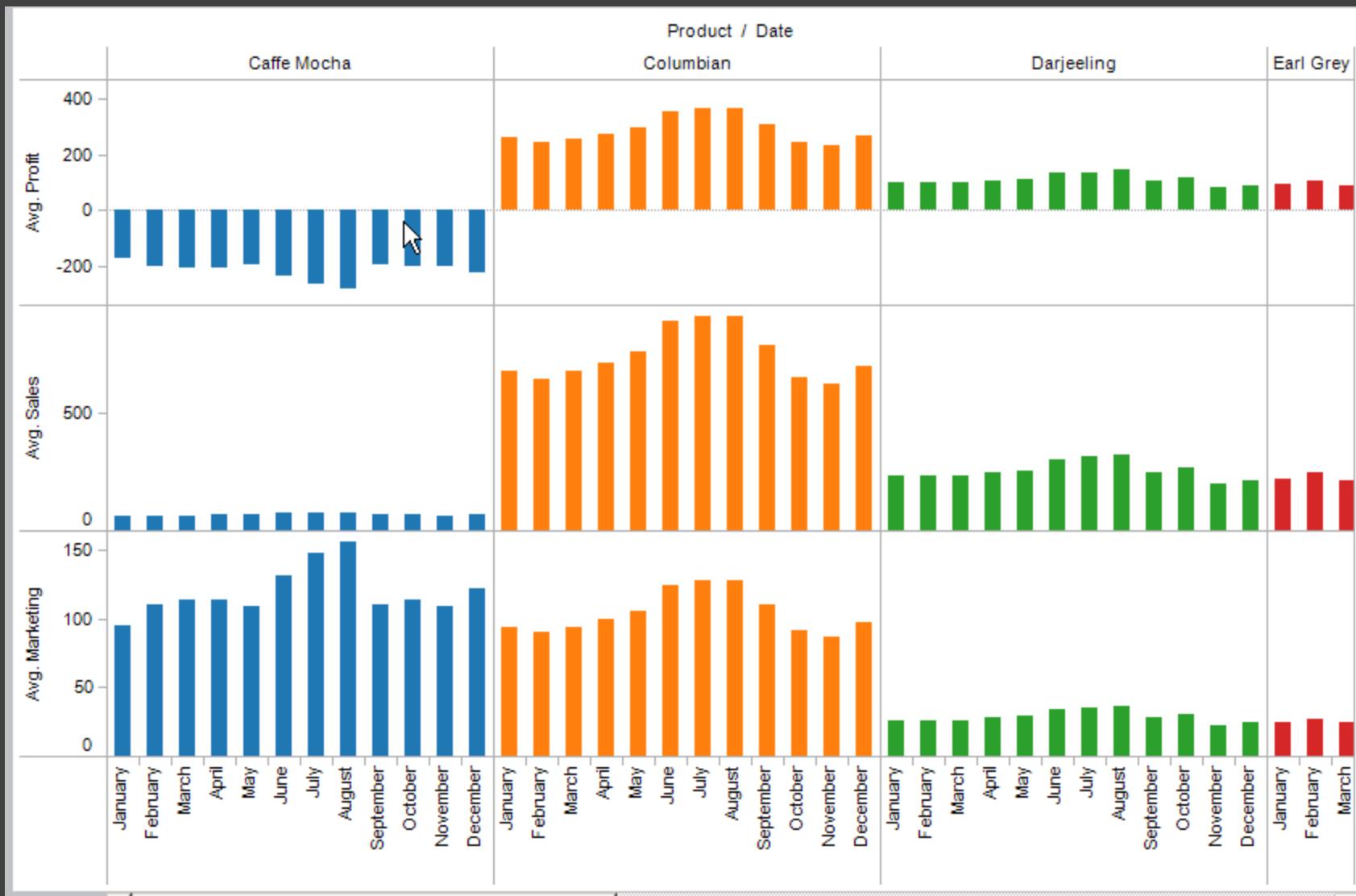
Ordinal-Ordinal

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	●	●	●	●

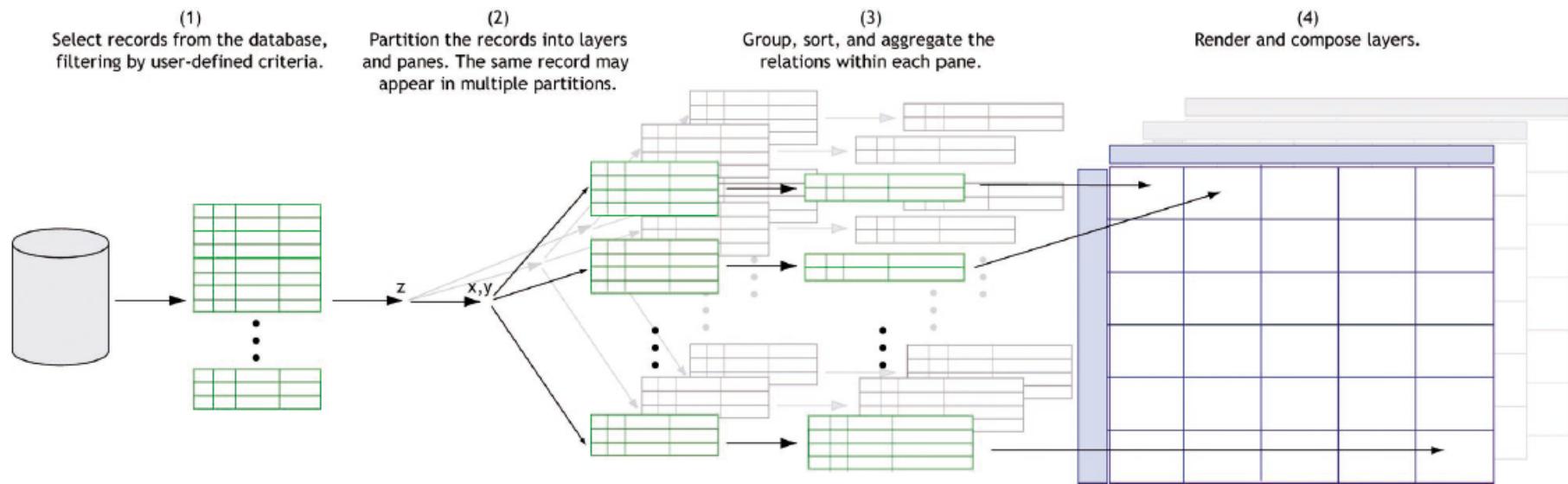
Quantitative-Quantitative



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.