

High Dimensional Data

Alark Joshi

High dimensional data

- Data with multiple dimensions, multiple variables or multiple attributes
- Cars dataset
 - Economy
 - Cylinders
 - Displacement
 - Power
 - Weight
 - Mph
 - Year

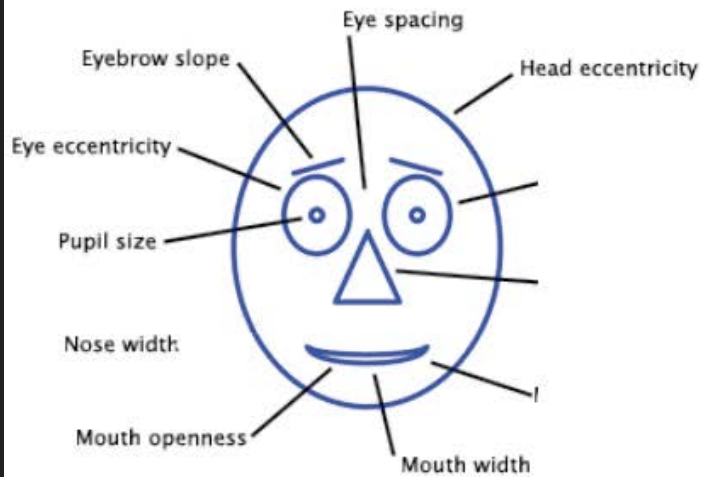
Scatterplots

- Great for visualizing 2D data
- Plot data attributes on x- and y-axis
- Scatterplot Matrix can be used to visualize multiple attributes

Scatterplot Matrix

- <http://mbostock.github.com/d3/ex/splom.html>

Chernoff Faces



Randomly selected parameters.
[Change face.](#)



Animation with random parameters.



All 0.



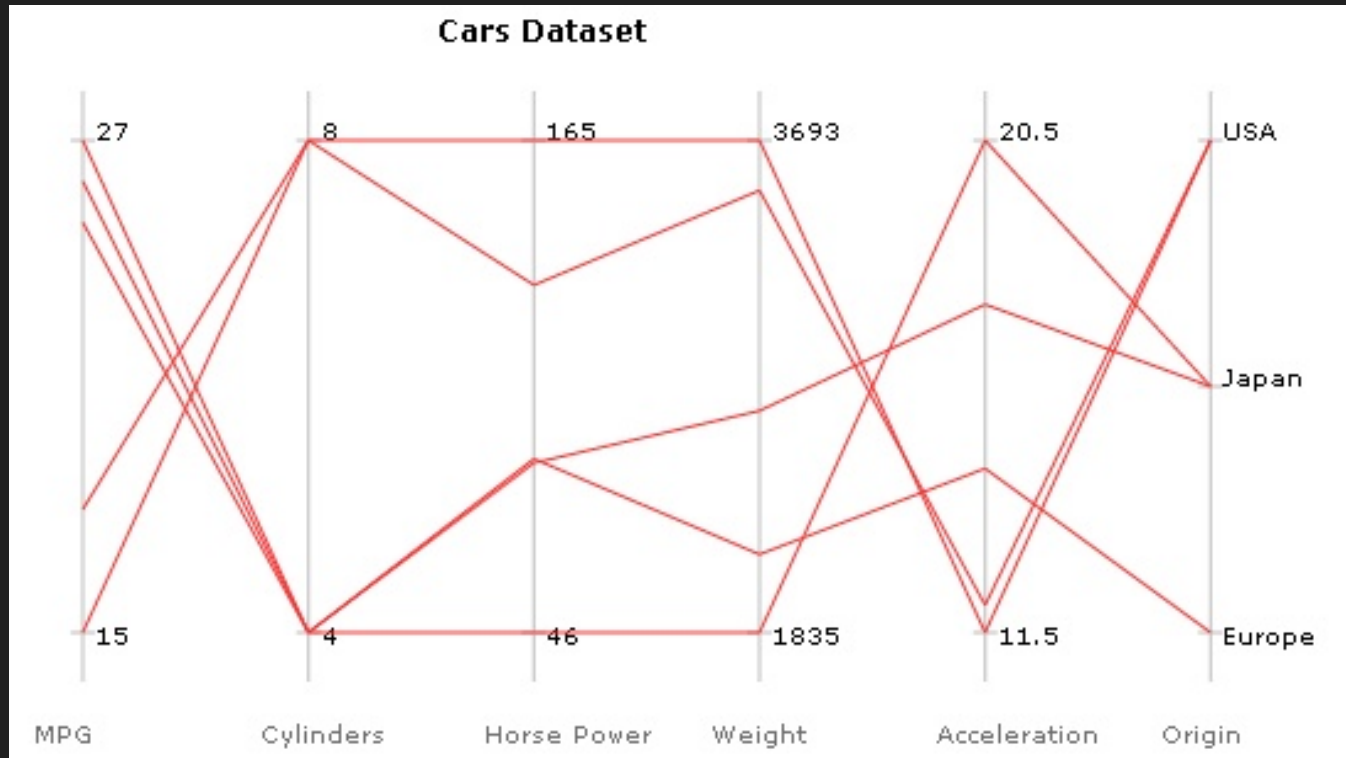
All 0.5.



All 1.

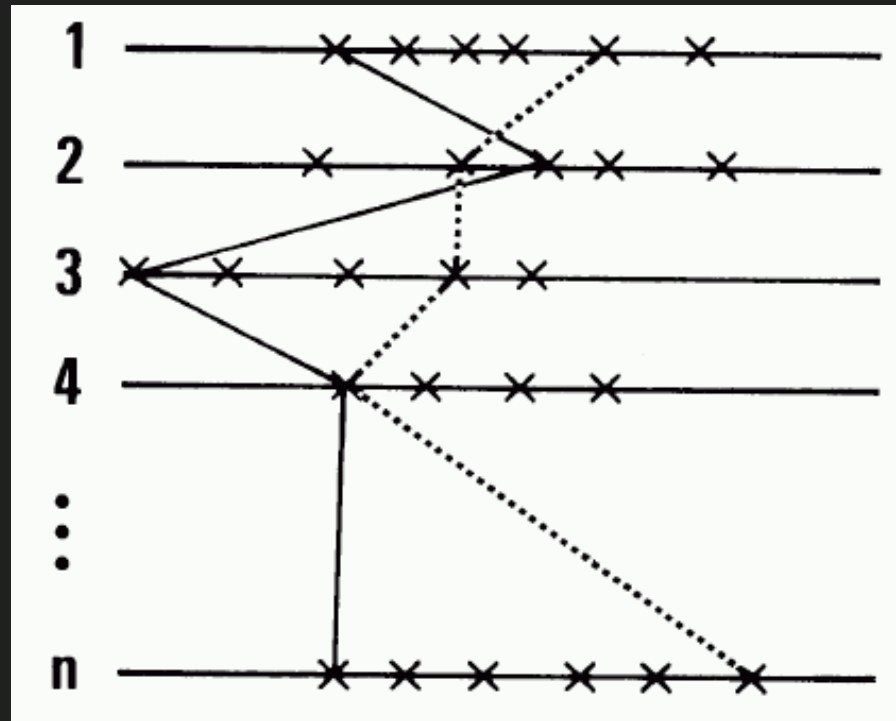
Parallel Coordinates

- Instead of having only 2 orthogonal axes (scatter plots), have parallel axes



Parallel Coordinates

- Connect variables for each data entity with a line



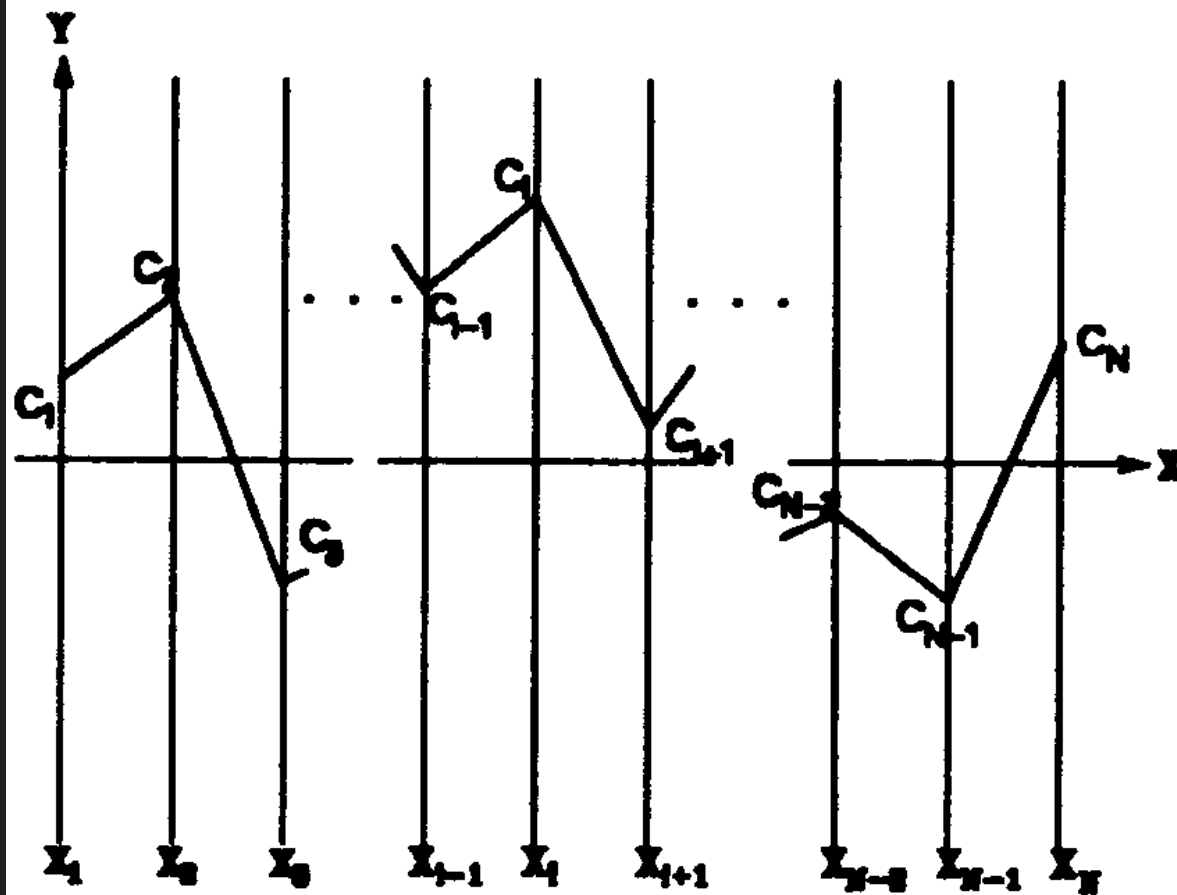


Figure 1: – Parallel axes for R^N .

The polygonal line shown represents the point $C = (c_1, \dots, c_{i-1}, c_i, c_{i+1}, \dots, c_N)$.

Five-dimensional hypersphere

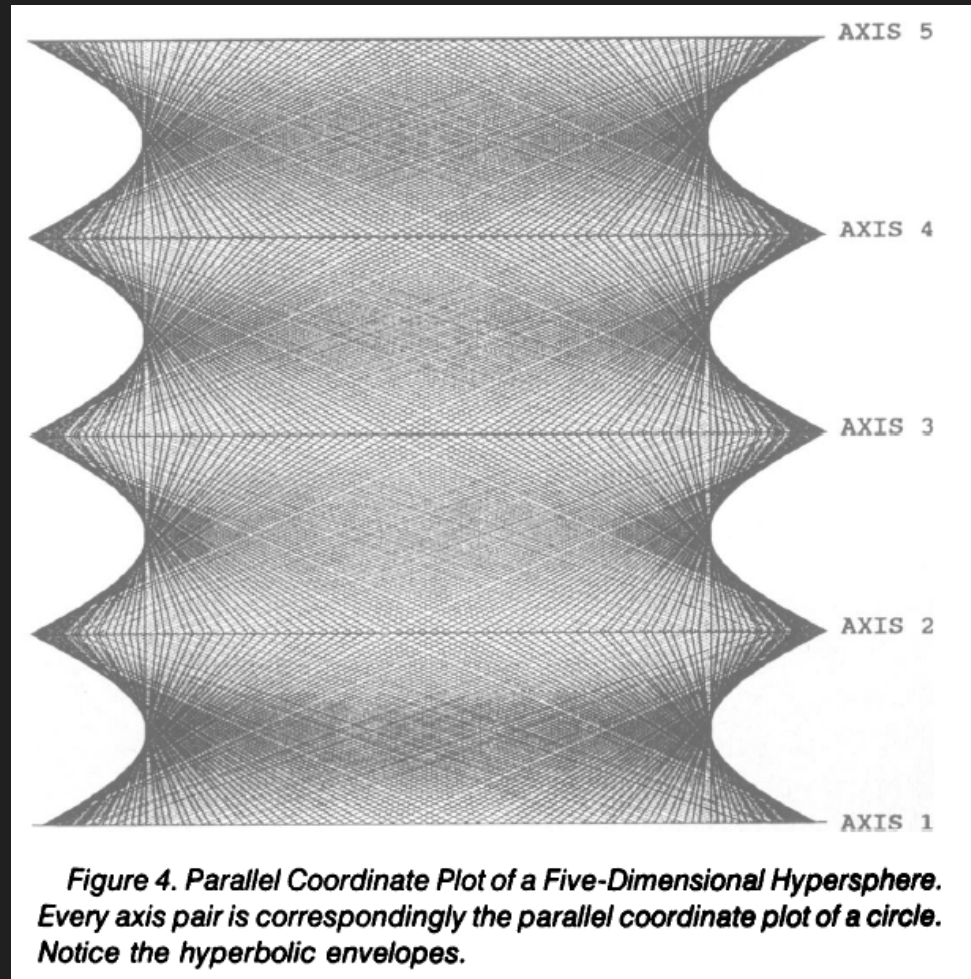
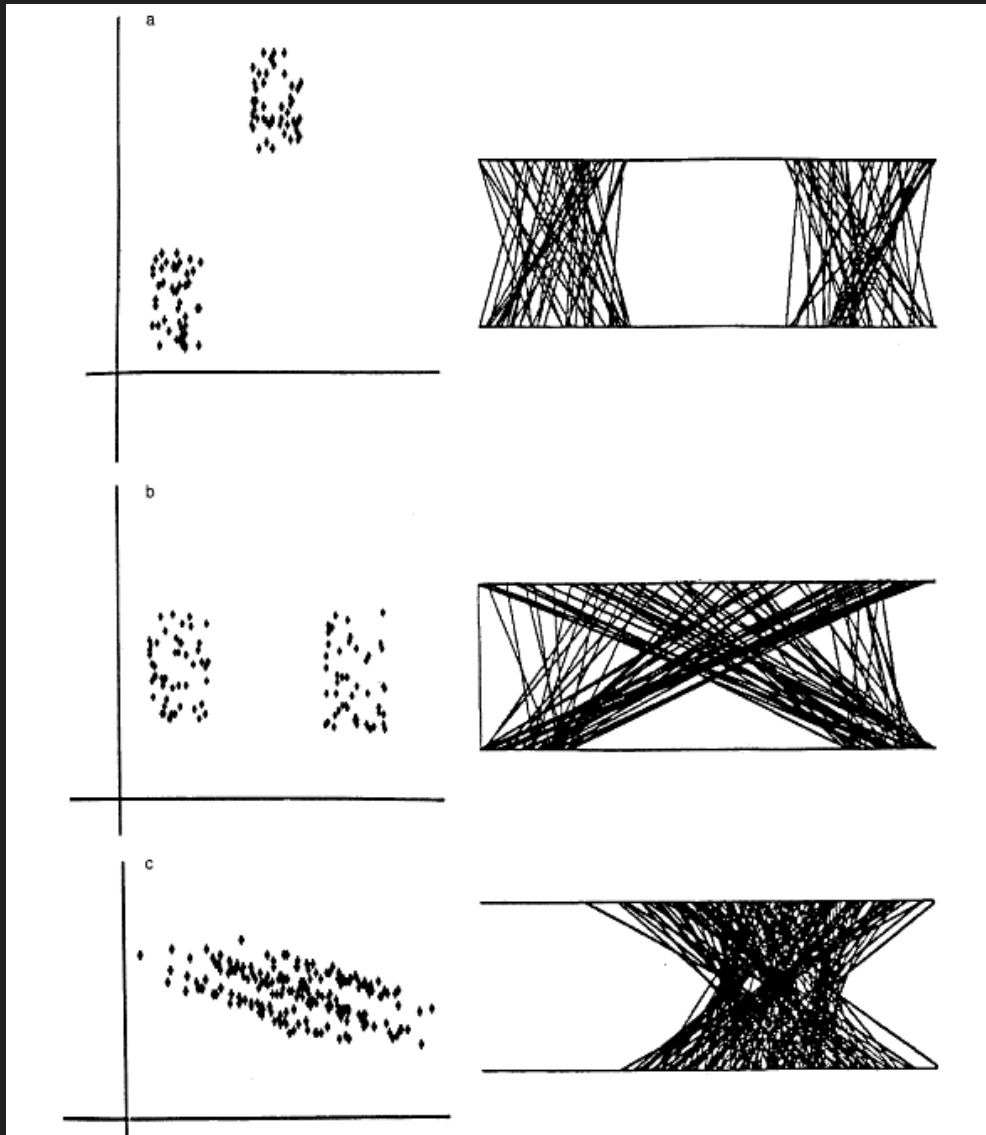


Image credits: Hyperdimensional Data Analysis Using Parallel Coordinates, Edward J. Wegman
Journal of the American Statistical Association , Vol. 85, No. 411 (Sep., 1990), pp. 664-675

Clustering in scatterplots vs PC



Clustering separated in x and y

Clustering separated in x but not in y

Clustering not separated in either projection

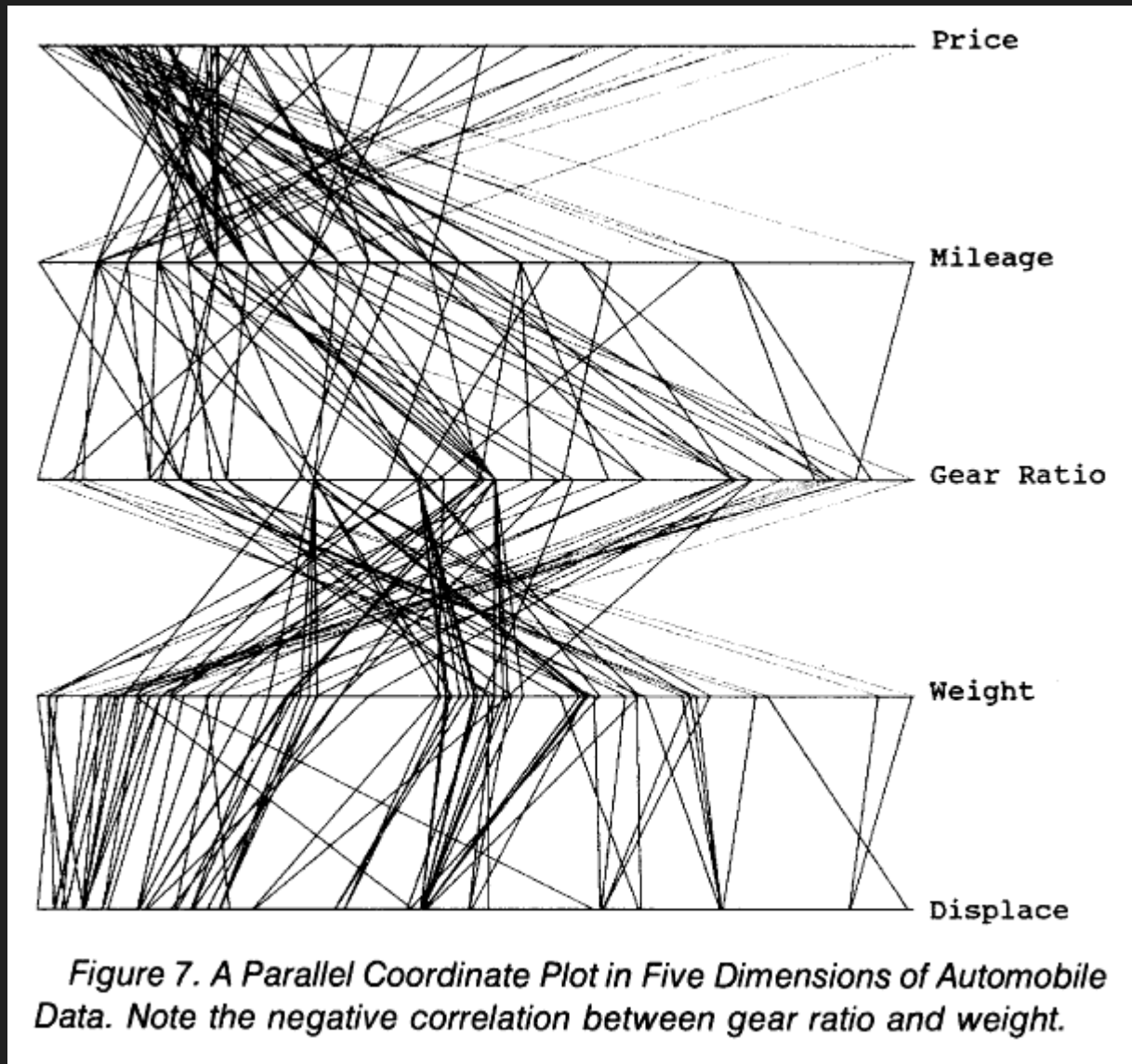


Image credits: Hyperdimensional Data Analysis Using Parallel Coordinates, Edward J. Wegman
Journal of the American Statistical Association , Vol. 85, No. 411 (Sep., 1990), pp. 664-675

PC Plot showing American Cars

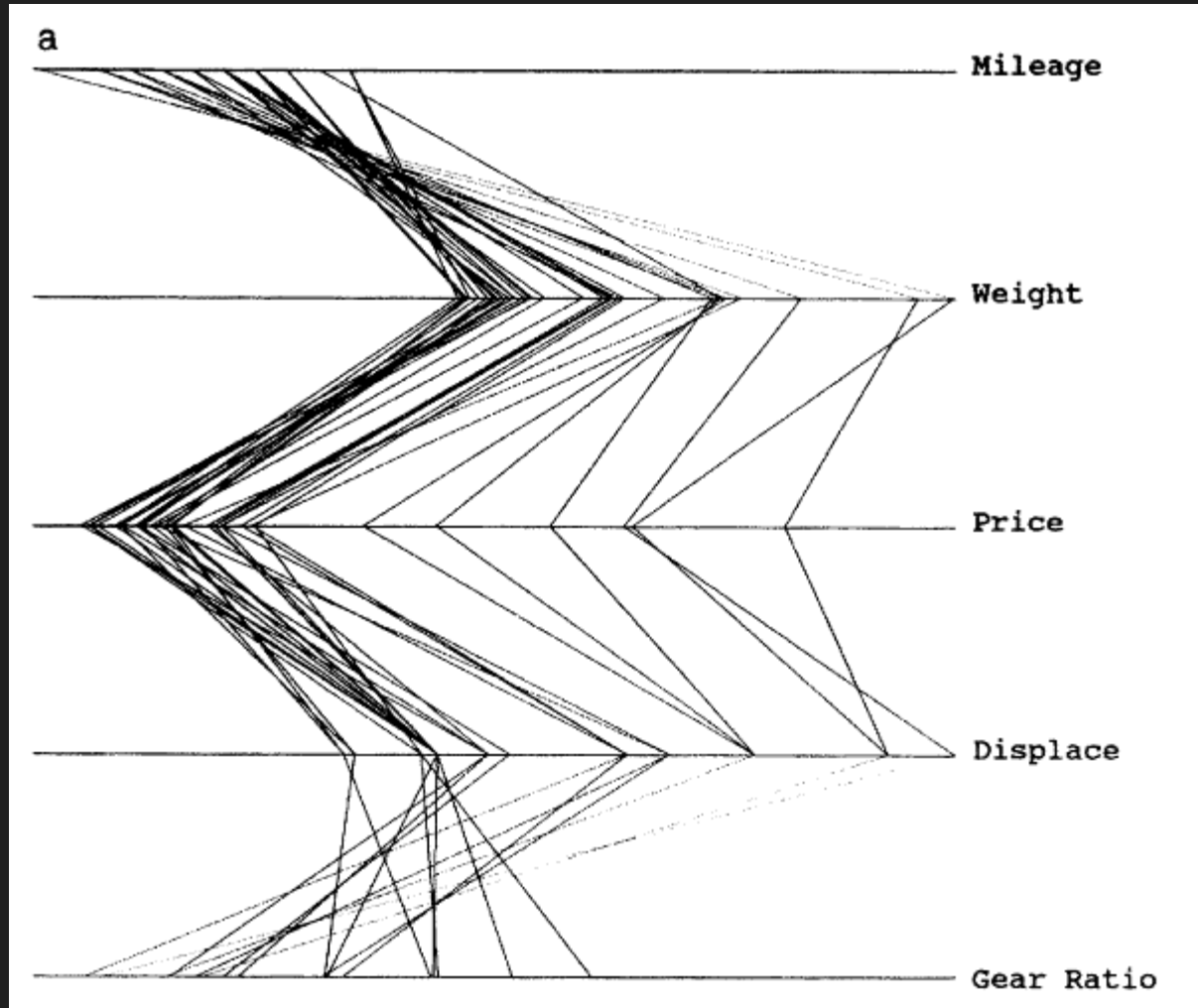


Image credits: Hyperdimensional Data Analysis Using Parallel Coordinates, Edward J. Wegman
Journal of the American Statistical Association , Vol. 85, No. 411 (Sep., 1990), pp. 664-675

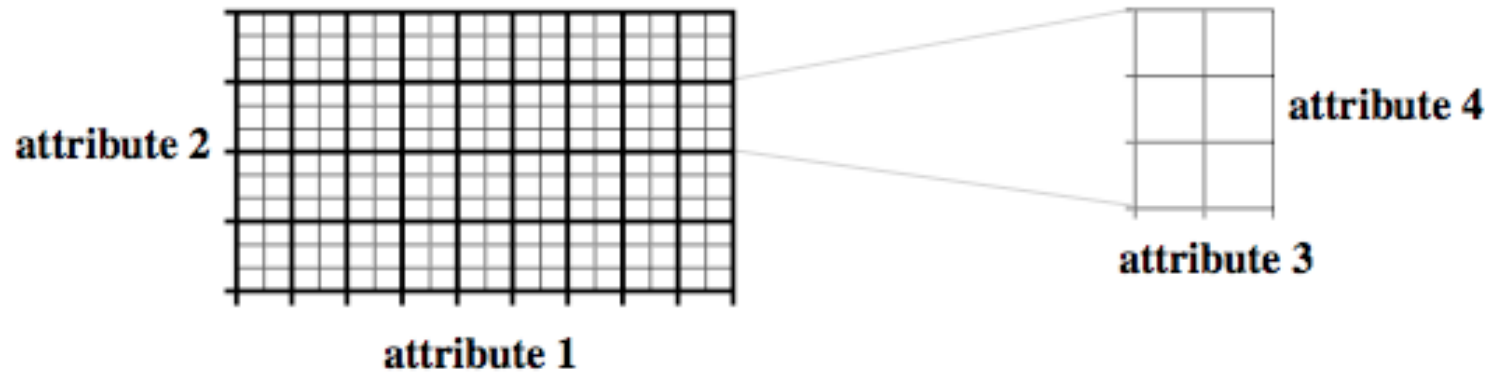
Demo

- Parallel coordinates in D3
 - <http://bl.ocks.org/1341281>

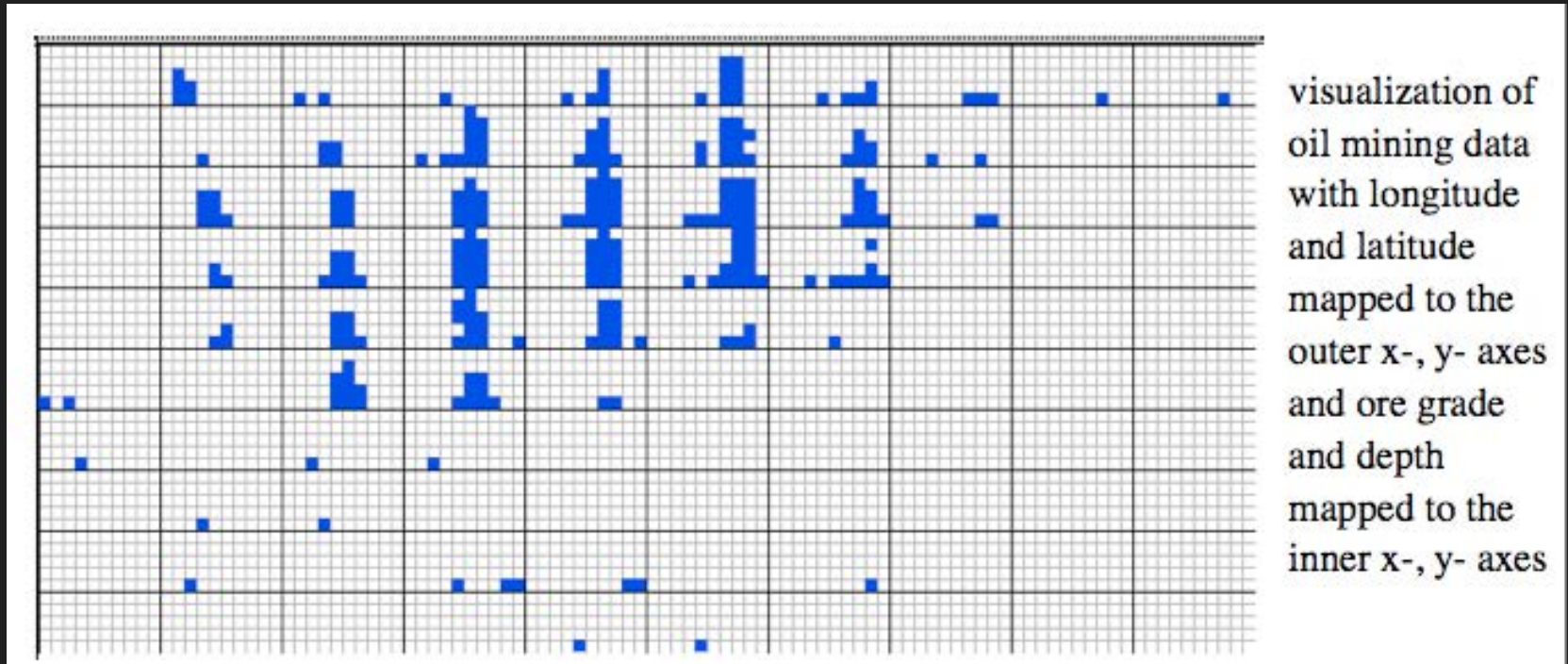
PC: Axis Ordering

- Geometric interpretations
 - Hyperplane, hypersphere
 - Points do have an intrinsic order
- Nominal data
 - No intrinsic order
 - Indeterminate/arbitrary order
 - Weakness of many techniques
 - Downside: human-powered search
 - Upside: Powerful interaction technique
- In most implementations, a user can interactively swap axes

Dimensionality Stacking



Dimensionality Stacking





Alphabetical



Median Value

Pixel-oriented techniques

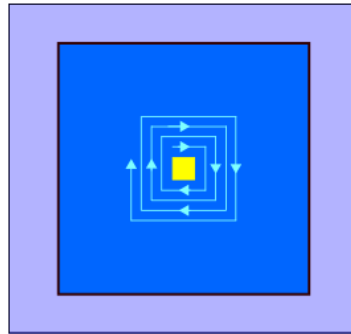


Figure 1:
Spiral Shaped Arrangement
of one Dimension

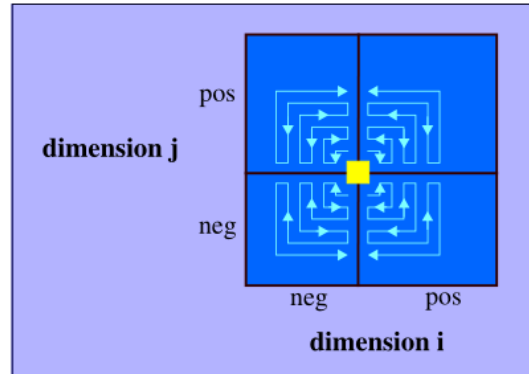


Figure 3: 2D-Arrangement of one Dimension

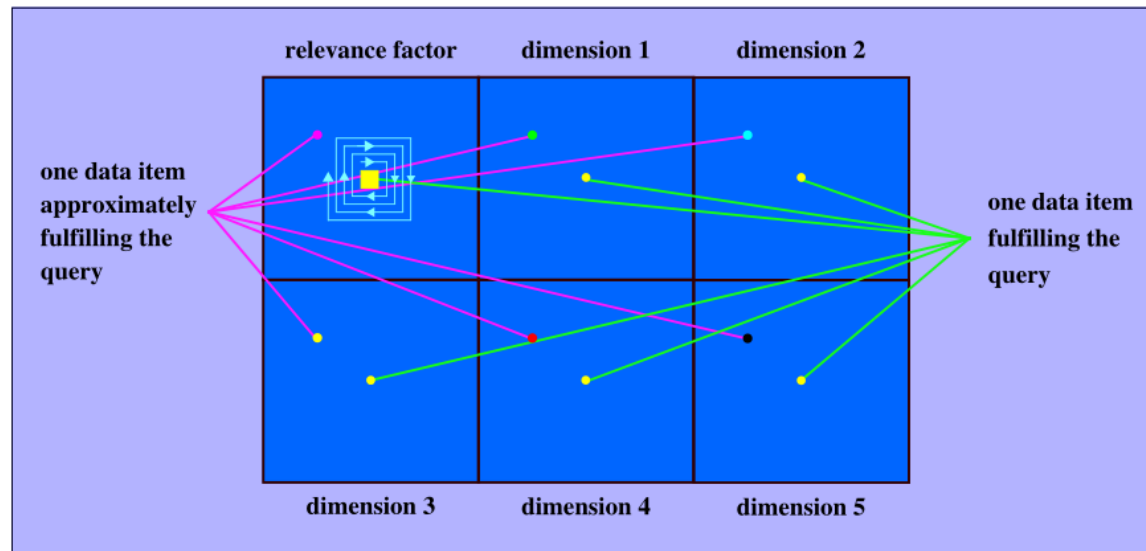
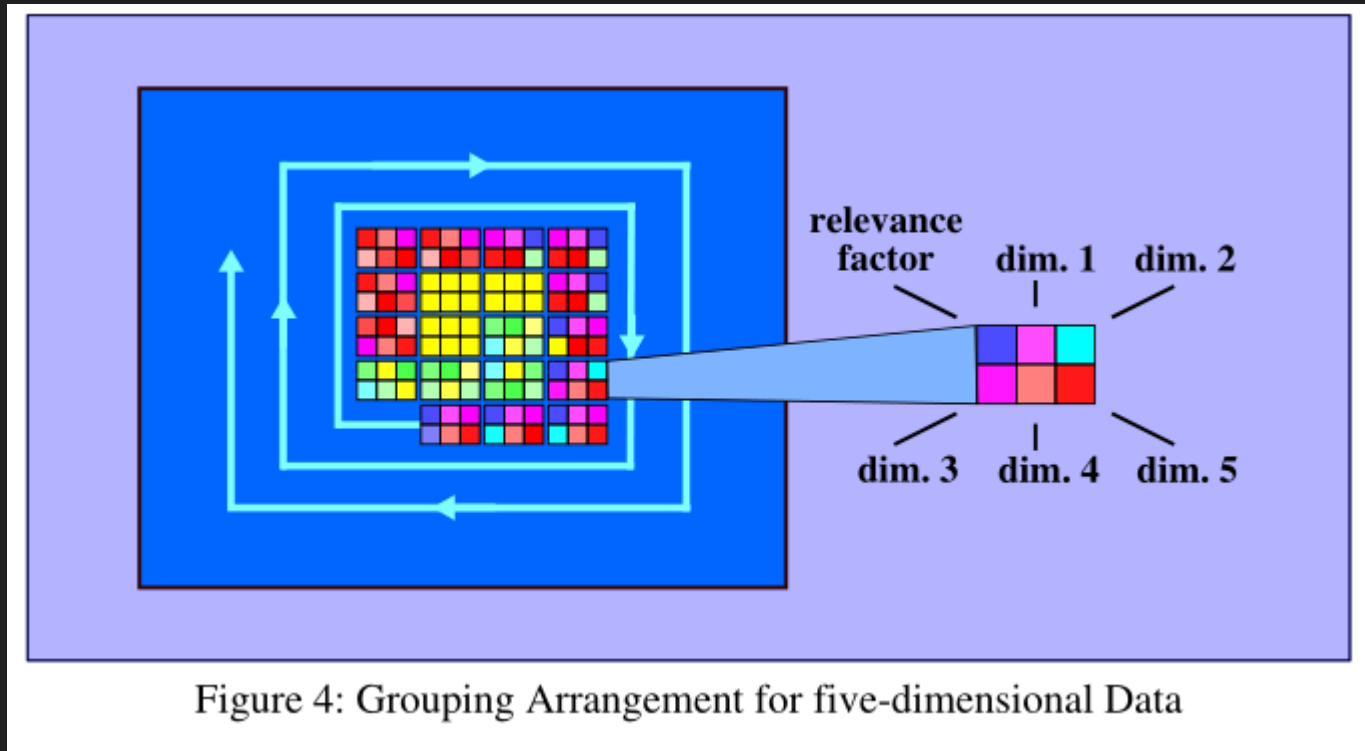
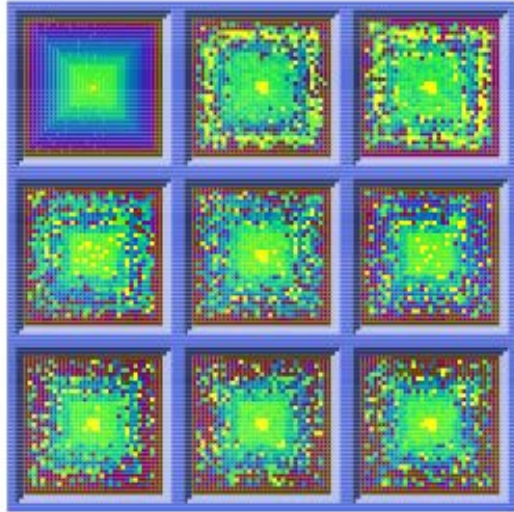


Figure 2: Arrangement of Windows for Displaying five-dimensional Data

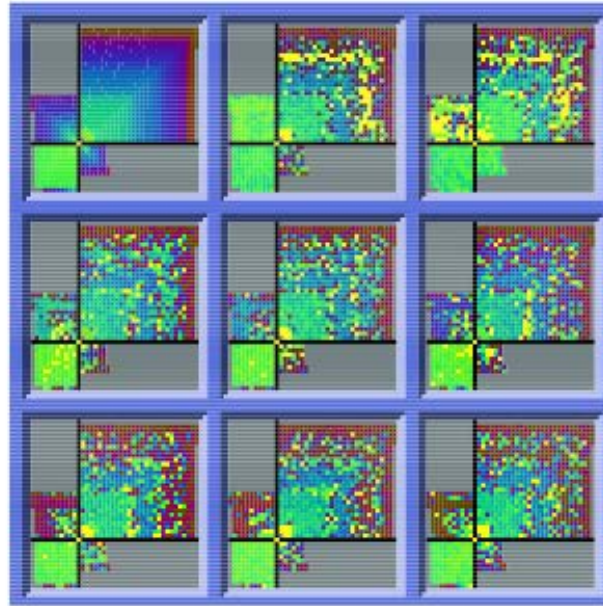
Pixel-oriented techniques



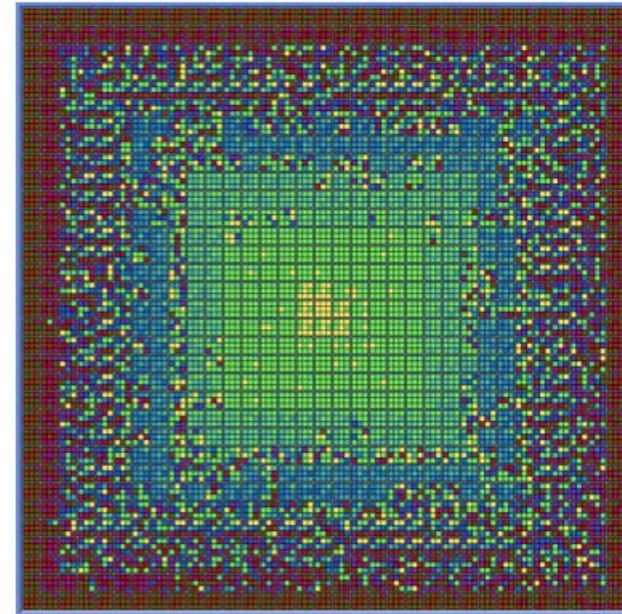
Visualizing 8-dimensional data



a. Basic Visualization Technique



b. 2D-Arrangement



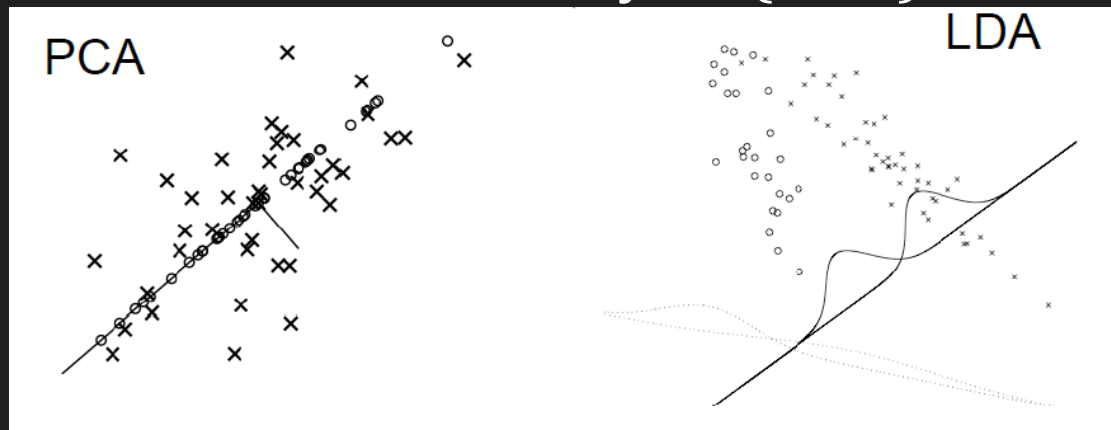
c. Grouping Arrangement

Dimensionality Reduction

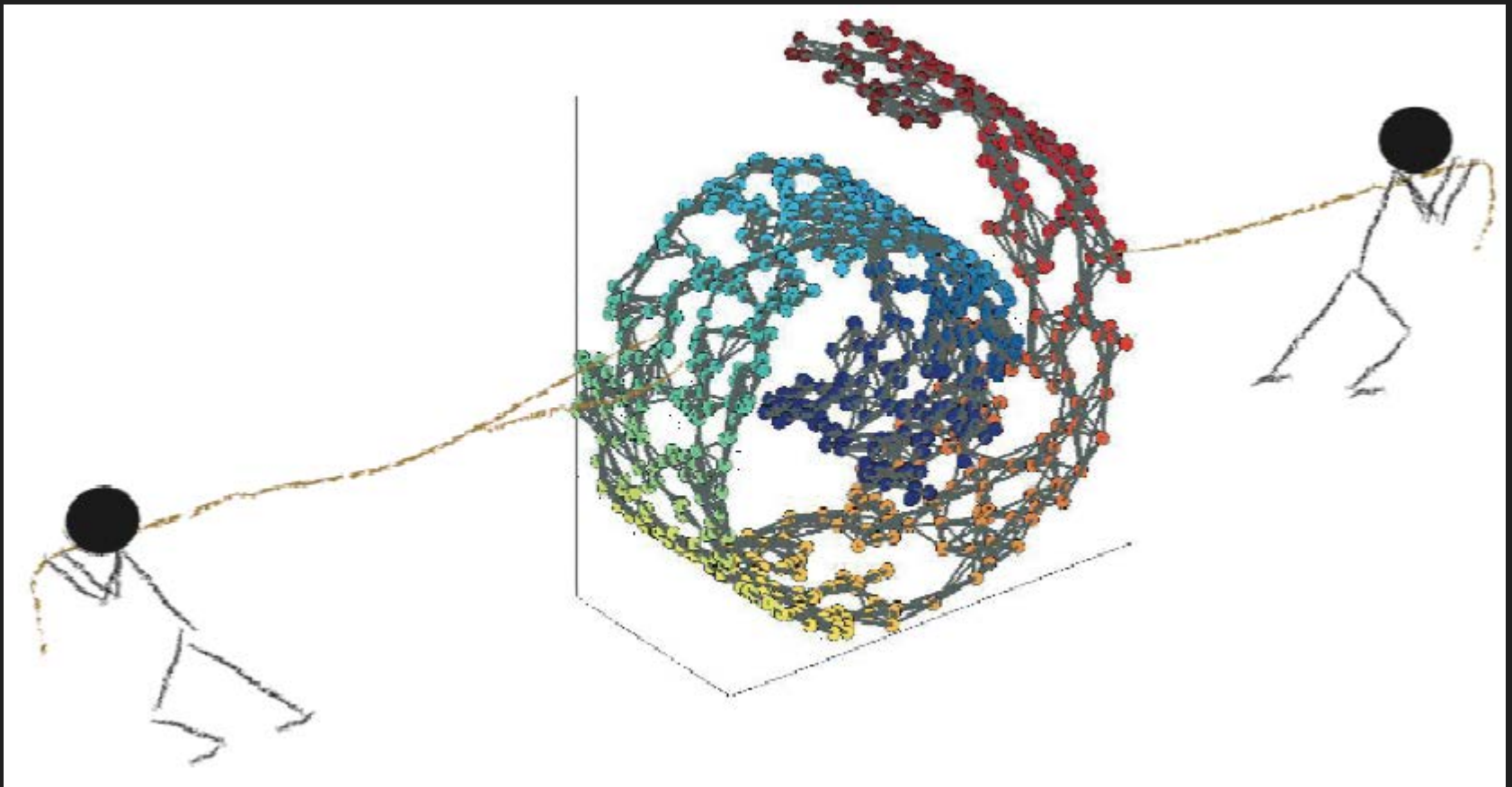
- Mapping multidimensional space into space of fewer dimensions
 - Typically 2D for clarify
 - 1D/3D possible
 - Preserve and communicate variance in data as much as possible
 - Show underlying structure of data
- Linear vs non-linear approaches

Linear Dimensionality Reduction

- Based on linear projections
- Given dimensions has a strong meaning
- Preserve the linearity in the layout
- Examples:
 - Principal Component Analysis (PCA)
 - Independent Component Analysis (ICA)
 - Linear Discriminant Analysis (LDA), ...



Problems for Linear Approaches

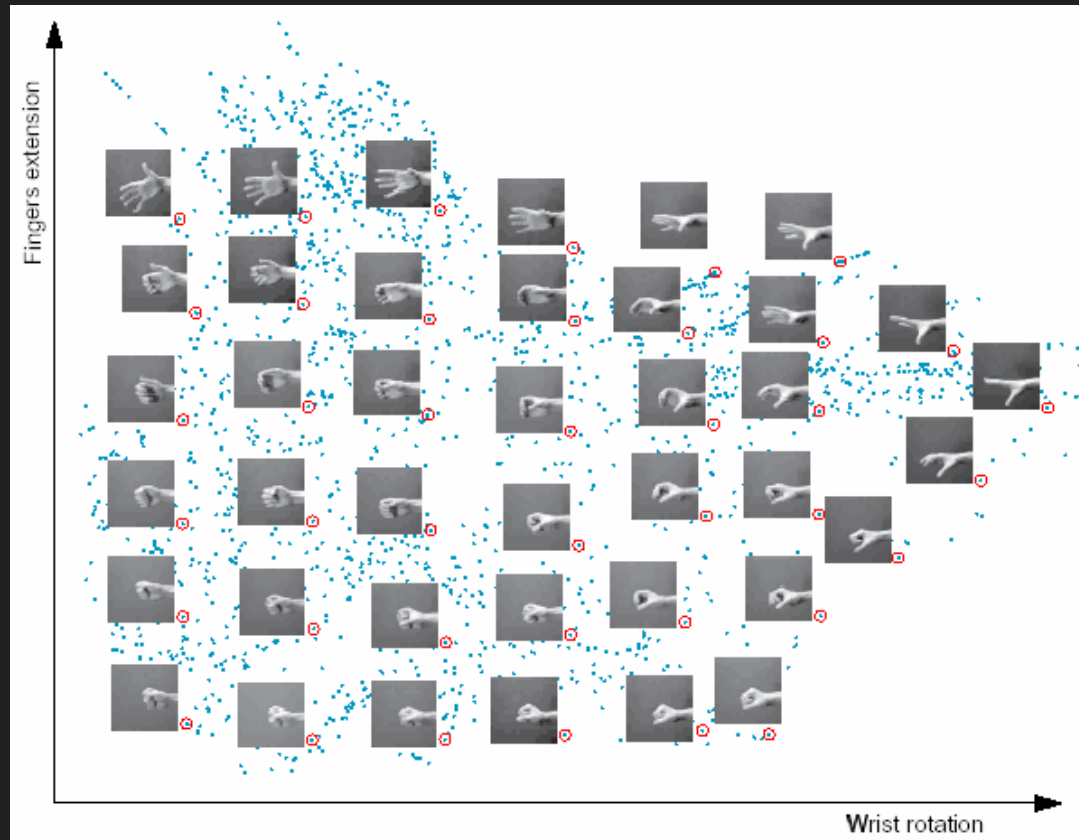


Non-linear Dimensionality Reduction

- Does not assume any inherent meaning to given dimensions
- Minimize differences between interpoint distances in high and low dimensions
- Examples:
 - Multidimensional scaling (MDS)
 - Isomap
 - Local linear embedding (LLE)

Isomap

- 4096 D to 2D
- 2D: wrist rotation, fingers extension



Goals

- Preserve and communicate as much variance as possible
- Find and display clusters
 - Compare/evaluate with previous clustering algorithms
- Understand structure
 - Absolution position is not reliable
 - Fine grained structure not reliable

Hierarchical Parallel Coordinates

- YH Fua, MO Ward, and IA Rundensteiner (1999), Hierarchical Parallel Coordinates for Exploration of Large Datasets, Proceedings of IEEE Visualization '99, pp. 43-50.
- Interactive visualization of large multivariate data sets
- Proposed a number of novel extensions to the parallel coordinates display technique
- Presentation by Danny

Dimension Ordering

- Determining dimension ordering important
 - Heuristic
 - Divide and conquer
 - Iterative hierarchical clustering
 - Representative dimensions

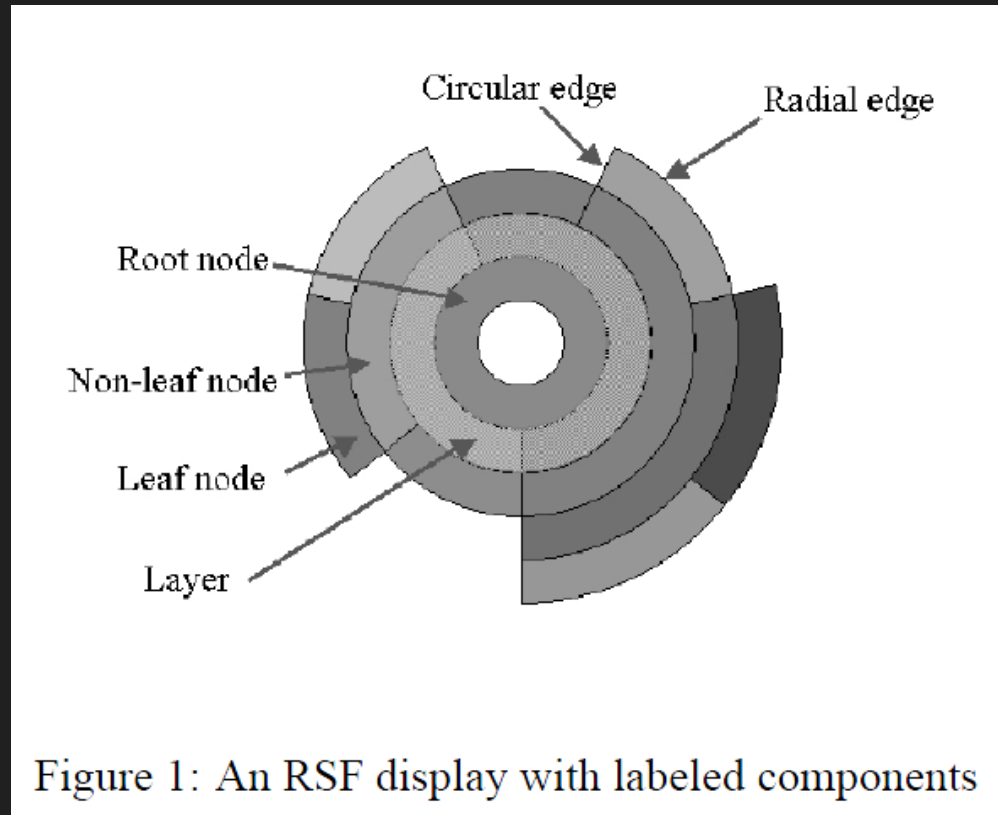
Dimension Ordering

- Choices
 - Similarity metrics
 - Importance metrics (variance, etc.)
 - Ordering algorithms
 - Optimal
 - Random swap
 - Simple depth-first traversal

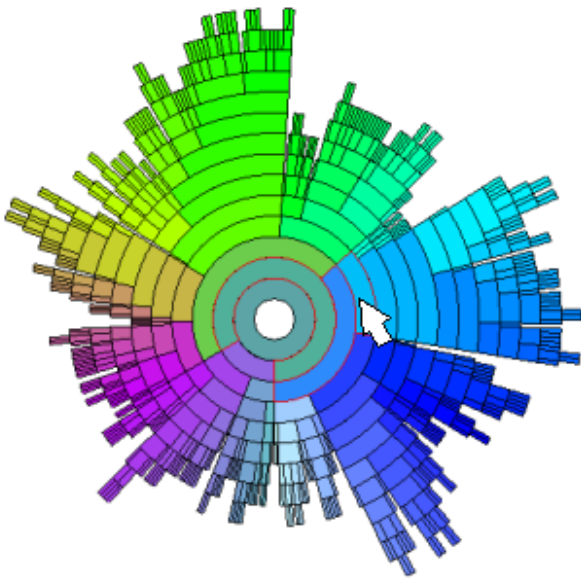
Dimension Filtering

- Interaction
 - Structure-based brushing
 - Focus + context
 - Manual interaction through UI components

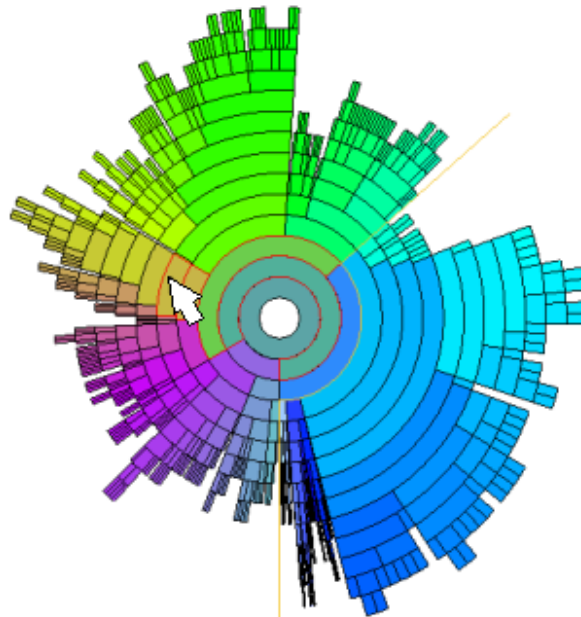
InterRing – Hierarchical Data Navigation



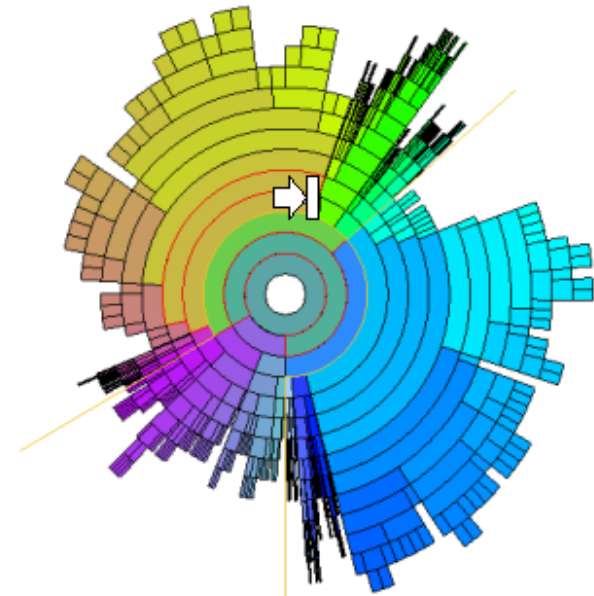
InterRing - MultiFocus Distortion



(a)

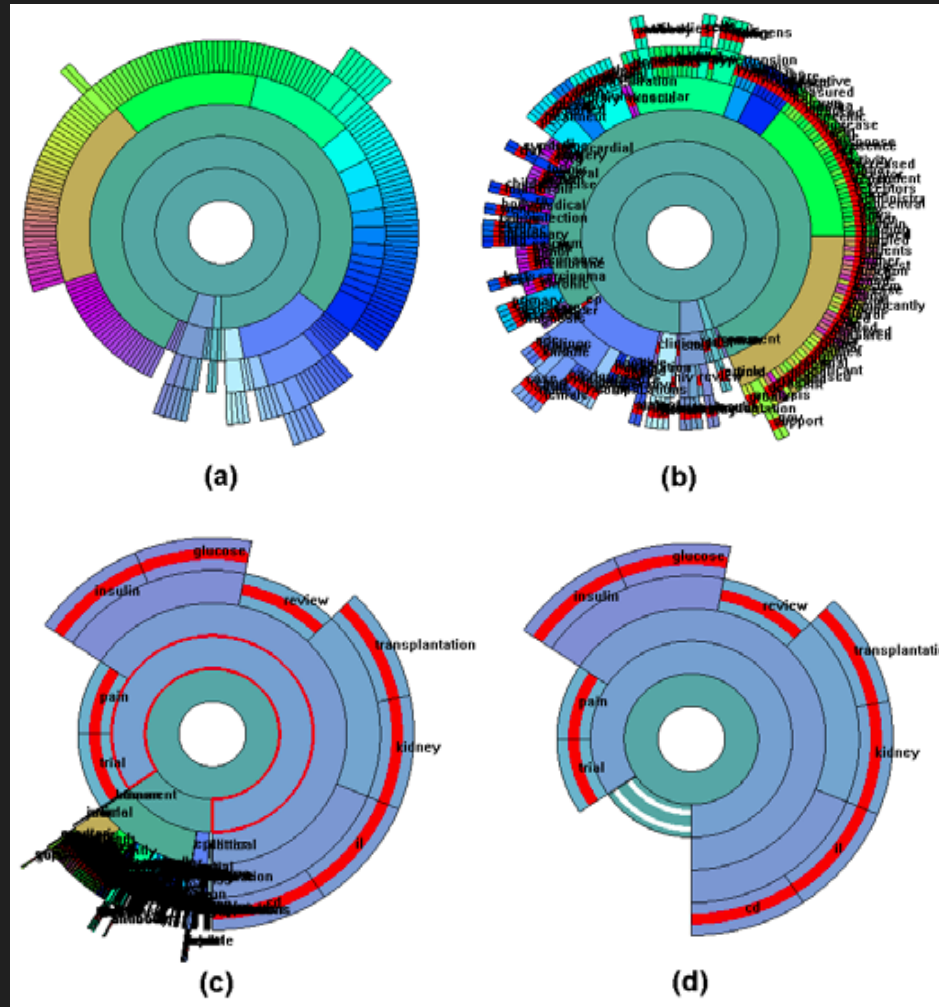


(b)



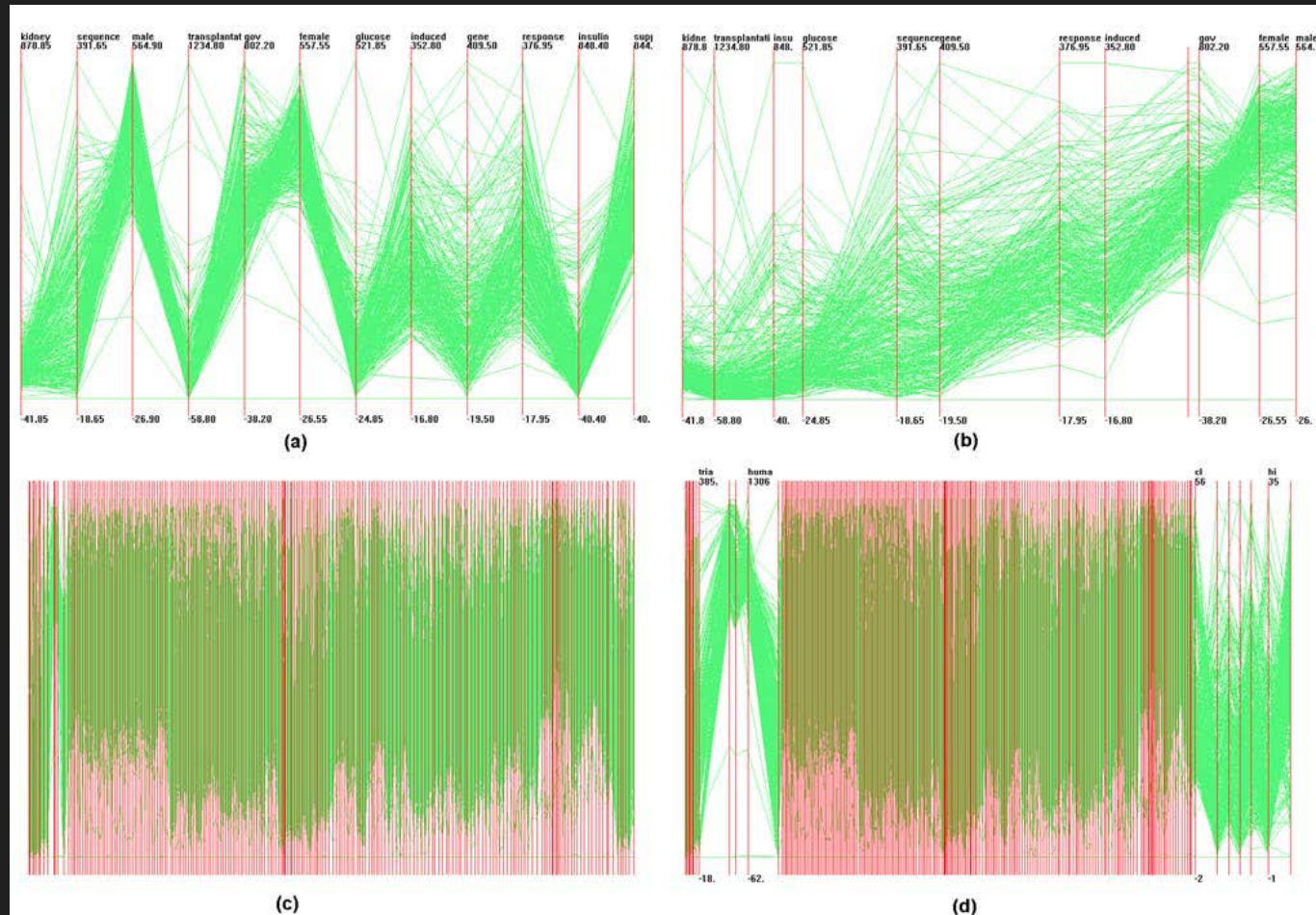
(c)

Filtering Interfaces - InterRing



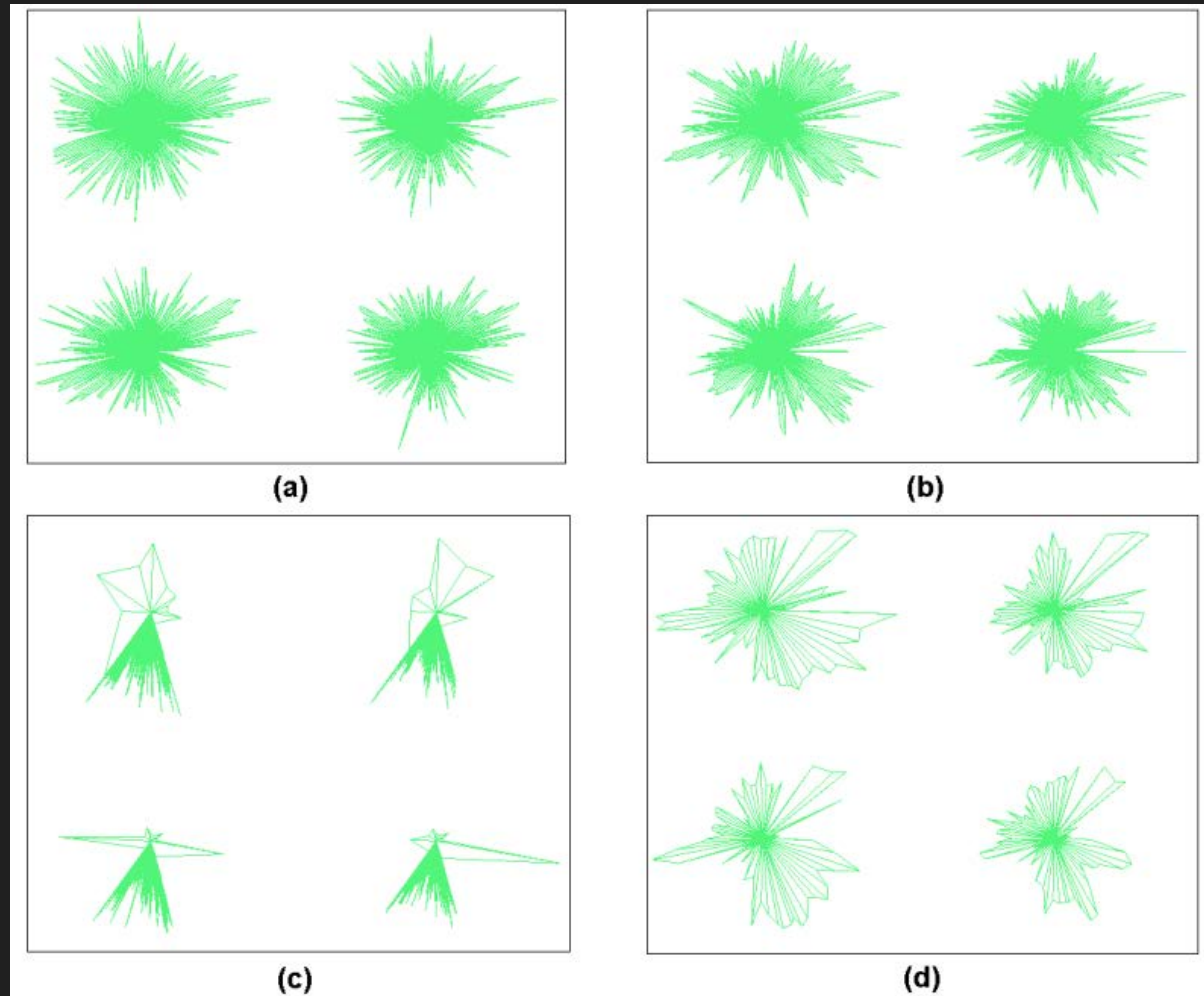
Raw, order, distort and rollup (filter)

Filtering Interfaces – Parallel Coordinates



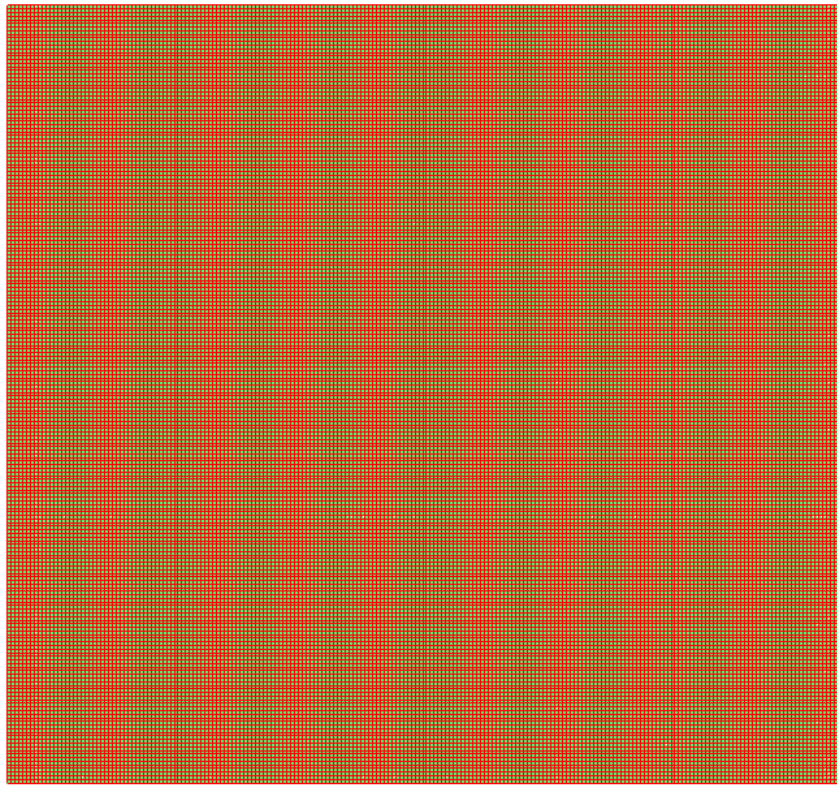
Raw, order/space, zoom and filter

Filtering Interfaces - InterRing

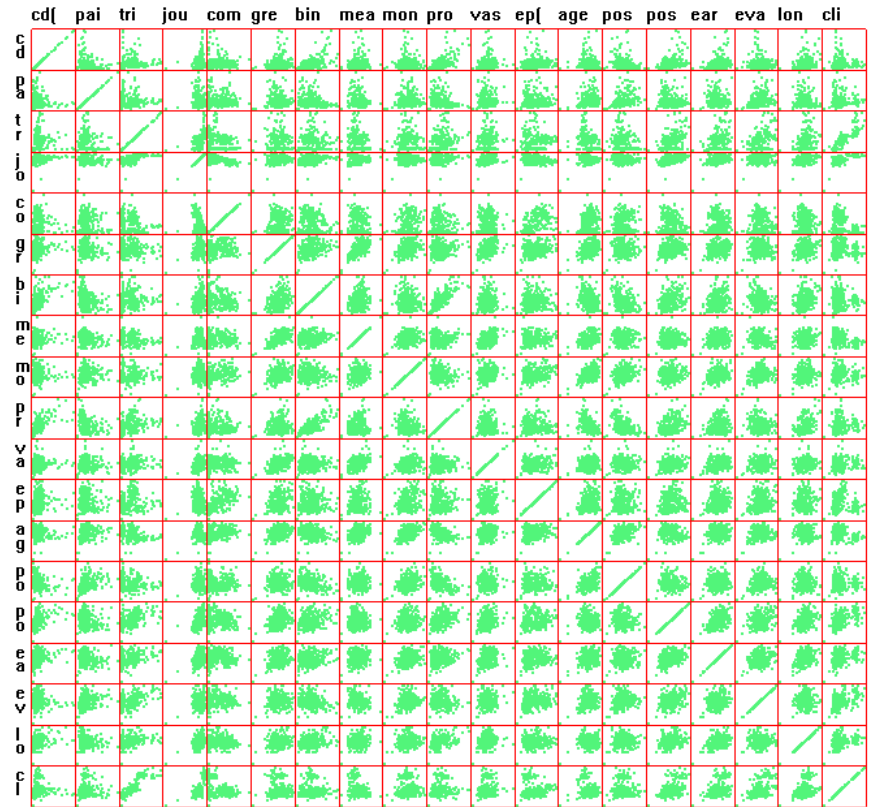


Raw, order/space, distort and filter

Filtering Interfaces - InterRing



(a)



(b)

Raw and filter

Polaris

- Multiscale Visualization Using Data Cubes, Chris Stolte, Diane Tang and Pat Hanrahan, Proc. InfoVis 2002.
- Stolte, C., Tang, D., and Hanrahan, P., Polaris: a system for query, analysis, and visualization of multidimensional databases, Commun. ACM 51, 11 (Nov. 2008), 75-84.

Large, Multi-Dimensional Databases

- Data acquisition not a problem anymore
- Extracting useful meaning from the data is a challenge
- “Path of exploration is unpredictable”
- Analysts want to be able to change the type of data and the visualization technique to examine the data
- Need to be able to visualize large subsets of data

Polaris

- An interactive exploration system that facilitates exploration of large, multi-dimensional relational databases
- Treat each attribute as a data cube (n-dimensional databases = n data-cubes)
- Polaris can facilitate multi-dimensional data exploration through a table-based display

Database Schema:

The user drags fields from the database schema to shelves to define the visual specification.

Layer Tabs:

Each layer has its own tab; different transformations and mappings can be specified for each layer.

Axis Shelves:

The fields placed here determine the structure of the table and the types of graphs in each table pane.

Context Menu:

The context menu provides access to the data transformation and interaction capabilities of Polaris such as sorting, filtering, and aggregation.

Layer Shelf:

The fields placed here determine how records are partitioned into layers.

Grouping and Sorting Shelves:

The fields placed here determine how records are grouped and sorted within the table panes.

Mark Pulldown:

Relations in each pane are mapped to marks of the selected type.

Retinal Property Shelves:

The fields placed here determine how data is encoded in the retinal properties of the marks.

Legends:

Legends enable the user to see and modify the mappings from data to retinal properties.

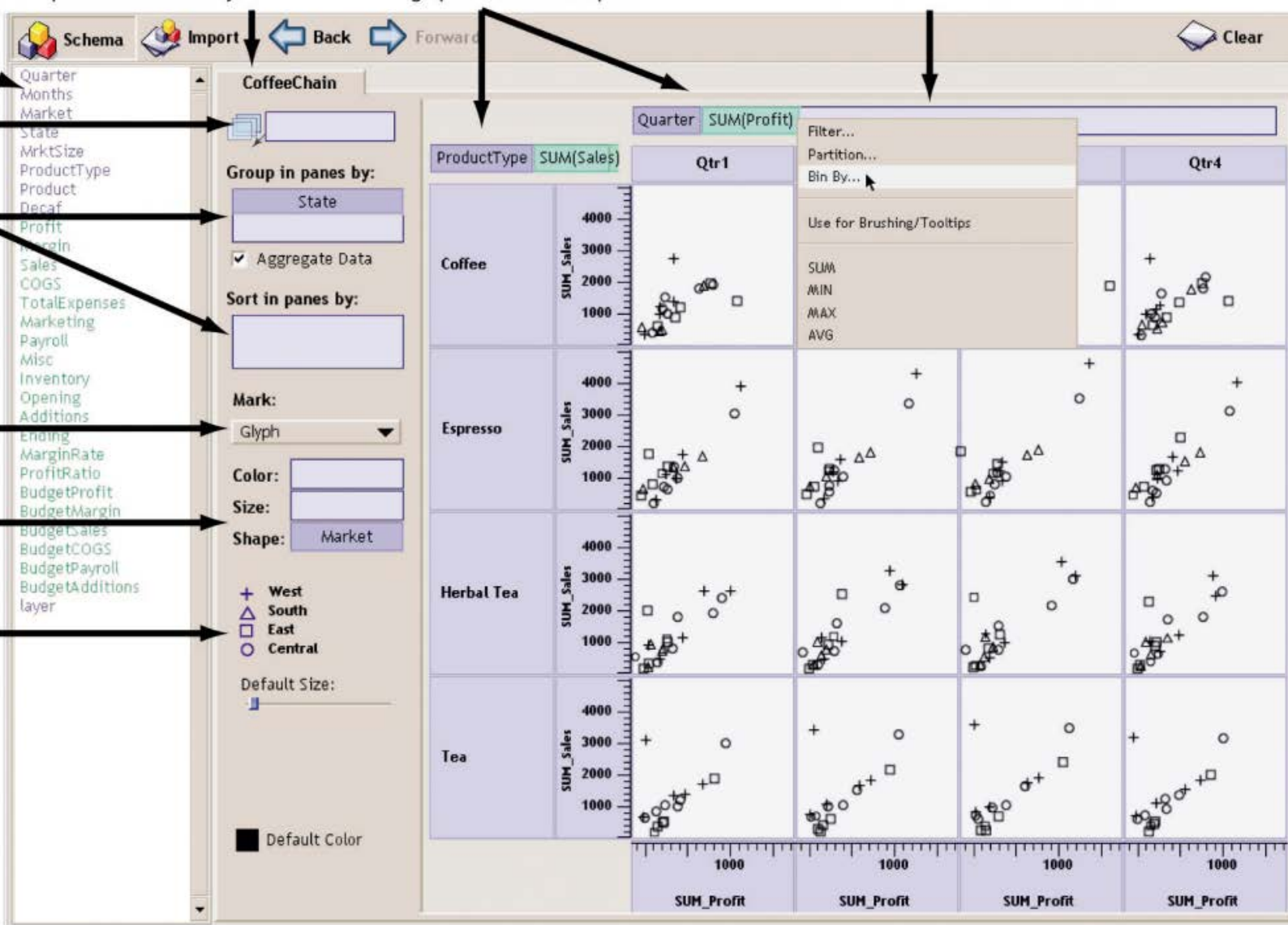


Table Algebra

- Define a formal mechanism to specify table configurations
- Consists of three separate expressions
 - Two expressions define the x and y axes of the table
 - Third expression defines the z-axis (partitions the display into layers)

The screenshot displays the Table Algebra interface with three main sections: Database Schema, Layer Tabs, and Axis Shelves. Arrows point from descriptive text labels to specific parts of the interface.

Database Schema: The user drags fields from the database schema to shelves to define the visual specification.

Layer Tabs: Each layer has its own tab; different transformations and mappings can be specified for each layer.

Axis Shelves: The fields placed here determine the structure of the table and the types of graphs in each table pane.

Layer Shelf: The fields placed here determine how records are partitioned into layers.

Grouping and Sorting Shelves: The fields placed here determine how records are grouped and sorted within the table panes.

The interface shows a 'CoffeeChain' layer tab. The 'Layer Shelf' contains 'Quarter', 'Months', 'Market', 'State', 'MrktSize', 'ProductType', 'Product', and 'Decaf'. The 'Grouping and Sorting Shelves' are set to 'State' and 'Aggregate Data'. The 'Axis Shelves' show 'ProductType' and 'SUM(Sales)' on the x-axis, and 'SUM(Profit)' on the y-axis. A scatter plot titled 'Coffee' is displayed, showing 'SUM_Sales' on the y-axis (ranging from 1000 to 4000) and 'Qtr1' on the x-axis. The plot shows data points for different products, with a '+' symbol indicating a specific data point.

Operators

$$A \times B = \{a_1, \dots, a_n\} \times \{b_1, \dots, b_m\}$$

$$= \{a_1 b_1, \dots, a_1 b_m, a_2 b_1, \dots, a_2 b_m, \dots, a_n b_1, \dots, a_n b_m\}$$

- Cross (x) operator: Cartesian product

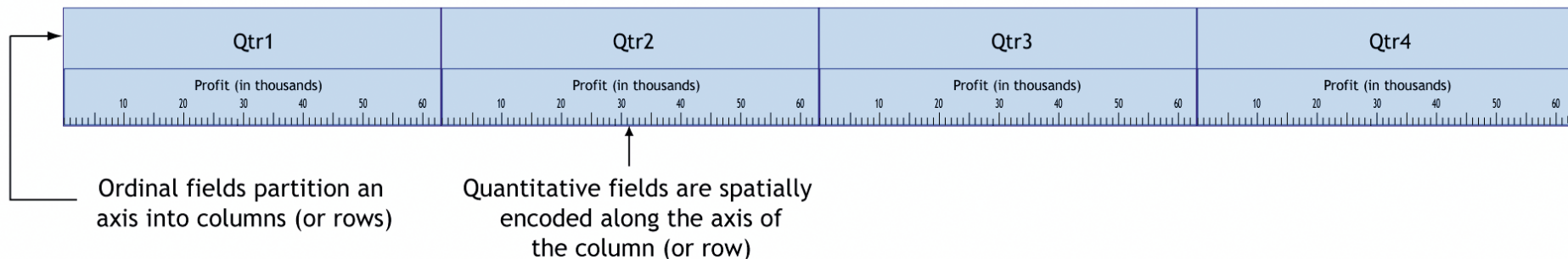
$O = \text{Quarter} = \{\text{Qtr1}, \text{Qtr2}, \text{Qtr3}, \text{Qtr4}\} = \text{Qtr1} + \text{Qtr2} + \text{Qtr3} + \text{Qtr4}$:

Qtr1	Qtr2	Qtr3	Qtr4
------	------	------	------

$O \times O = \text{Quarter} \times \text{Product} = \{(\text{Qtr1}, \text{Coffee}), (\text{Qtr1}, \text{Espresso}), (\text{Qtr1}, \text{Herbal Tea}), (\text{Qtr1}, \text{Tea}), (\text{Qtr2}, \text{Coffee}) \dots (\text{Qtr4}, \text{Tea})\}$:

Qtr1				Qtr2				Qtr3				Qtr4			
Coffee	Espresso	Herbal Tea	Tea	Coffee	Espresso	Herbal Tea	Tea	Coffee	Espresso	Herbal Tea	Tea	Coffee	Espresso	Herbal Tea	Tea

$O \times Q = \text{Quarter} \times \text{Profit} = \{(\text{Qtr1}, \text{Profit}), (\text{Qtr2}, \text{Profit}), (\text{Qtr3}, \text{Profit}), (\text{Qtr4}, \text{Profit})\}$:



$$A/B = \{a_i b_j \mid \exists r \in R \text{ st } A(r) = a_i \text{ \& } B(r) = b_i\}.$$

Operators

- Nest (/) operator: $A/B = B$ within A

$O = \text{Quarter} = \{\text{Qtr1}, \text{Qtr2}, \text{Qtr3}, \text{Qtr4}\} = \text{Qtr1} + \text{Qtr2} + \text{Qtr3} + \text{Qtr4}$:

Qtr1	Qtr2	Qtr3	Qtr4
------	------	------	------

$O/O = \text{Quarter} / \text{Month} = \{(\text{Qtr1}, \text{Jan}), (\text{Qtr1}, \text{Feb}), (\text{Qtr1}, \text{Mar}), (\text{Qtr2}, \text{Apr}), (\text{Qtr2}, \text{May}) \dots (\text{Qtr4}, \text{Dec})\}$:

Qtr1			Qtr2			Qtr3			Qtr4		
Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec



The set entry (Qtr4,Nov)
corresponds to this column

$$A + B = \{a_1, \dots, a_n\} + \{b_1, \dots, b_m\} \\ = \{a_1, \dots, a_n, b_1, \dots, b_m\}$$

Operators

- Concatenation (+) operator

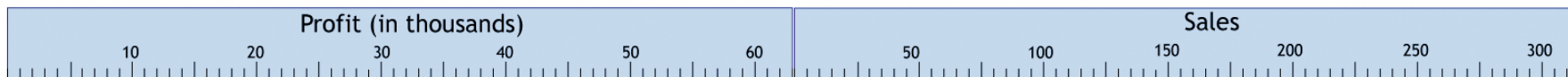
$O = \text{Quarter} = \{\text{Qtr1}, \text{Qtr2}, \text{Qtr3}, \text{Qtr4}\} = \text{Qtr1} + \text{Qtr2} + \text{Qtr3} + \text{Qtr4}$:

Qtr1	Qtr2	Qtr3	Qtr4
------	------	------	------

$O + O = \text{Quarter} + \text{Product} = \{\text{Qtr1}, \text{Qtr2}, \text{Qtr3}, \text{Qtr4}, \text{Coffee}, \text{Espresso}, \text{Herbal Tea}, \text{Tea}\}$:

Qtr1	Qtr2	Qtr3	Qtr4	Coffee	Espresso	Herbal Tea	Tea
------	------	------	------	--------	----------	------------	-----

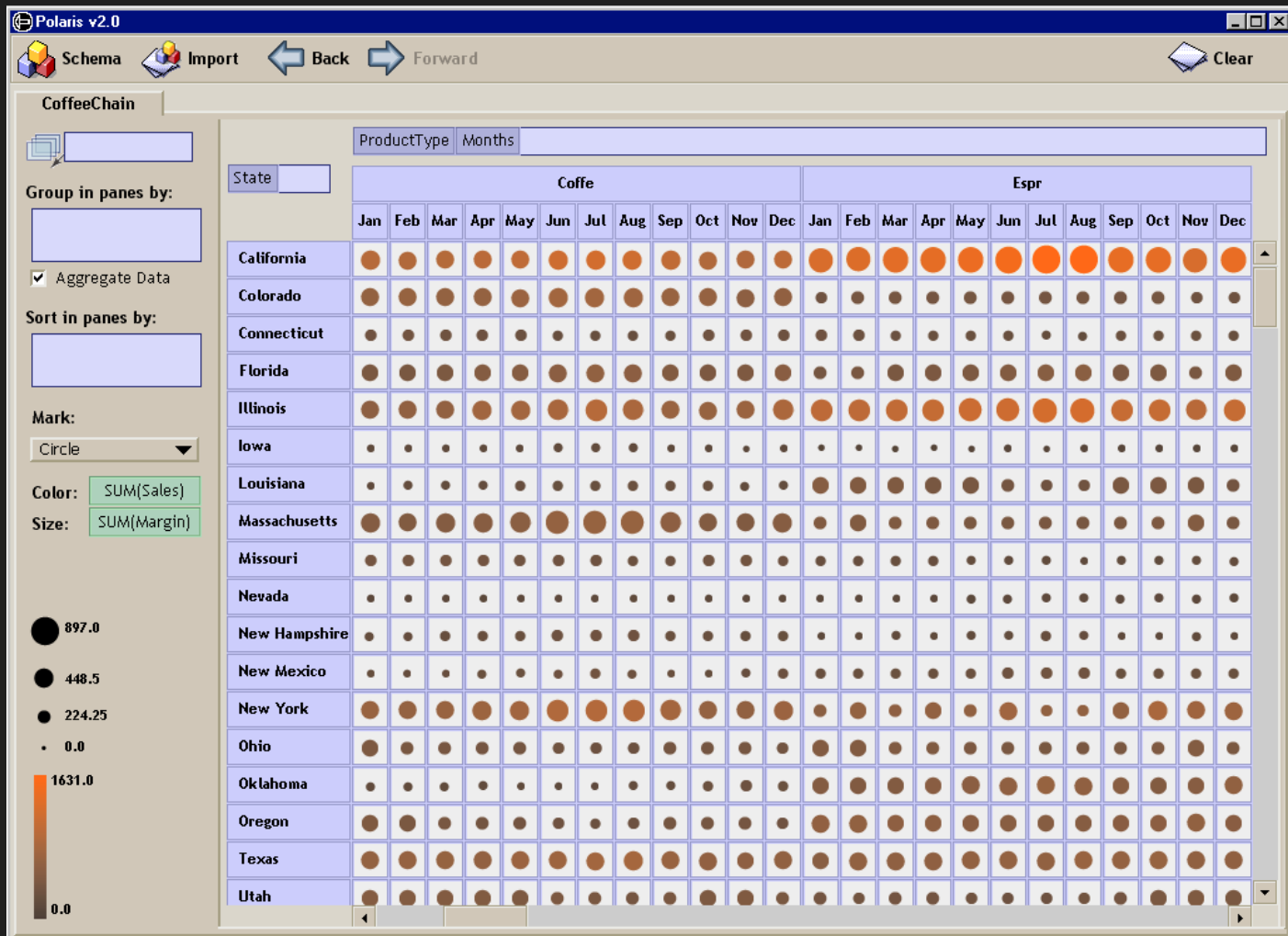
$Q + Q = \text{Profit} + \text{Sales} = \{\text{Profit}, \text{Sales}\}$:



Space of Graphics

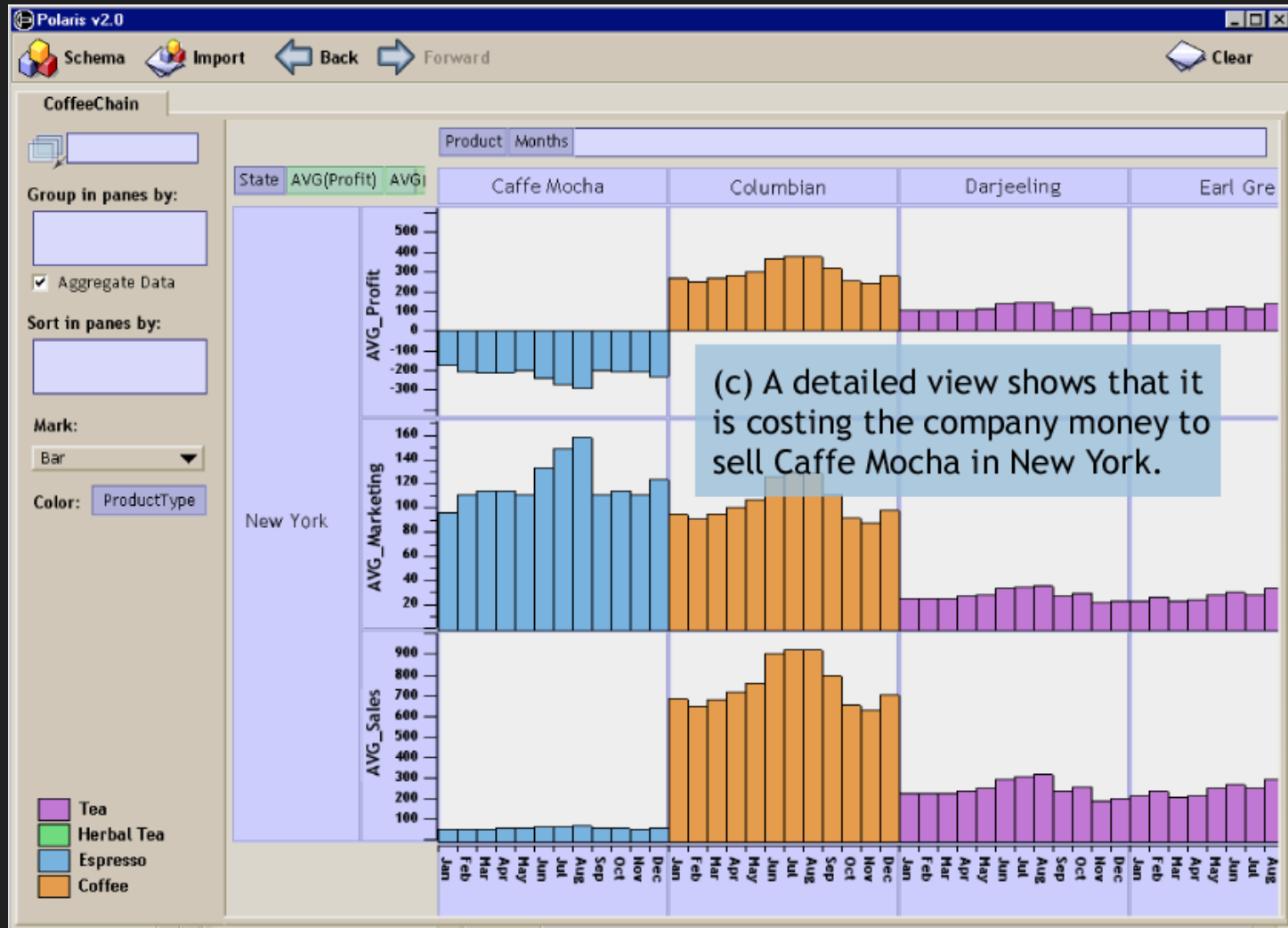
- Structured into three families
 - Ordinal-Ordinal
 - Ordinal-Quantitative
 - Quantitative - Quantitative

Ordinal-Ordinal



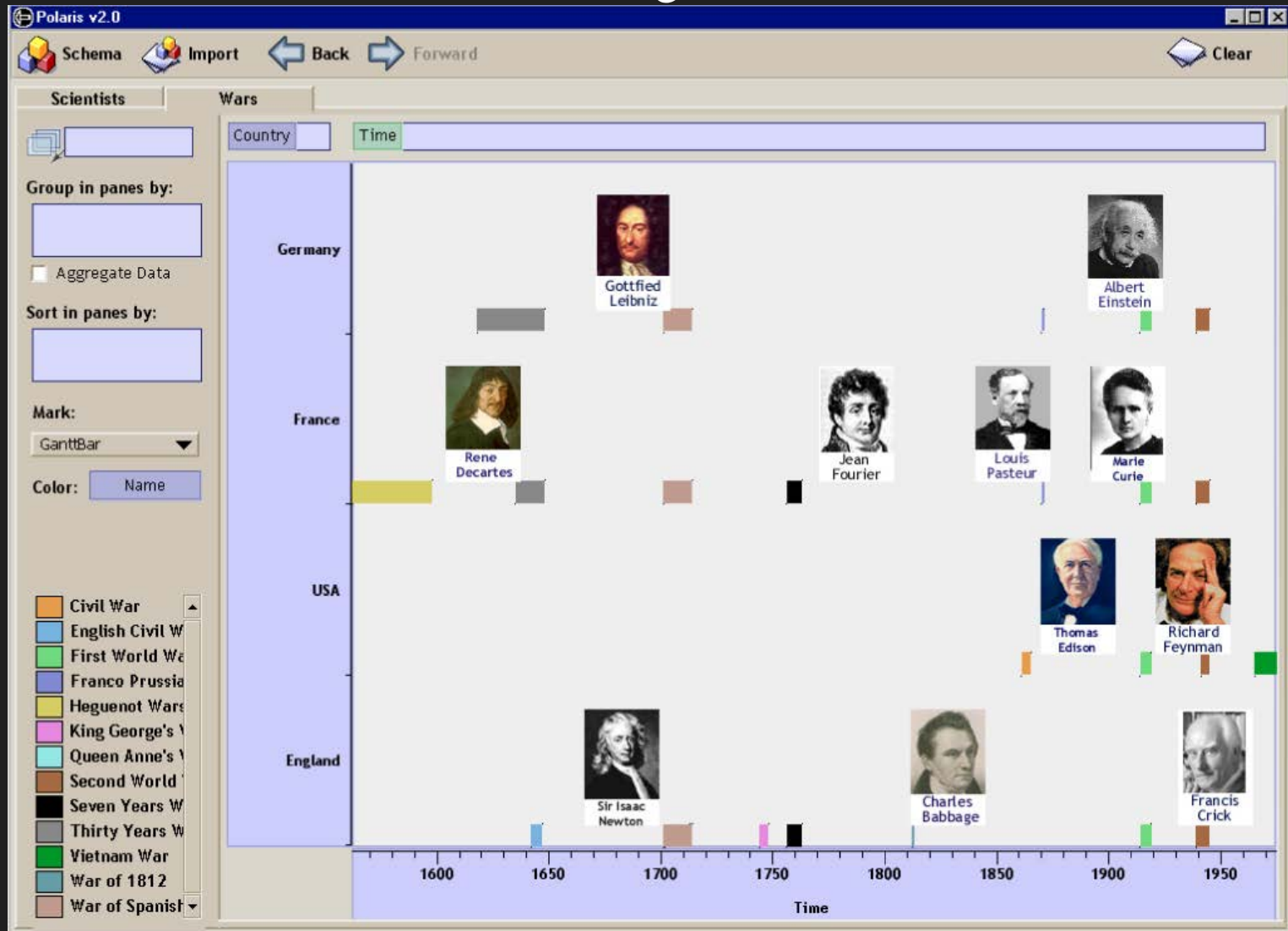
Sales and margins vs product type, month and state for the items sold

Ordinal - Quantitative



Matrix of bar charts is used to study independent variables – product and month

Ordinal - Quantitative

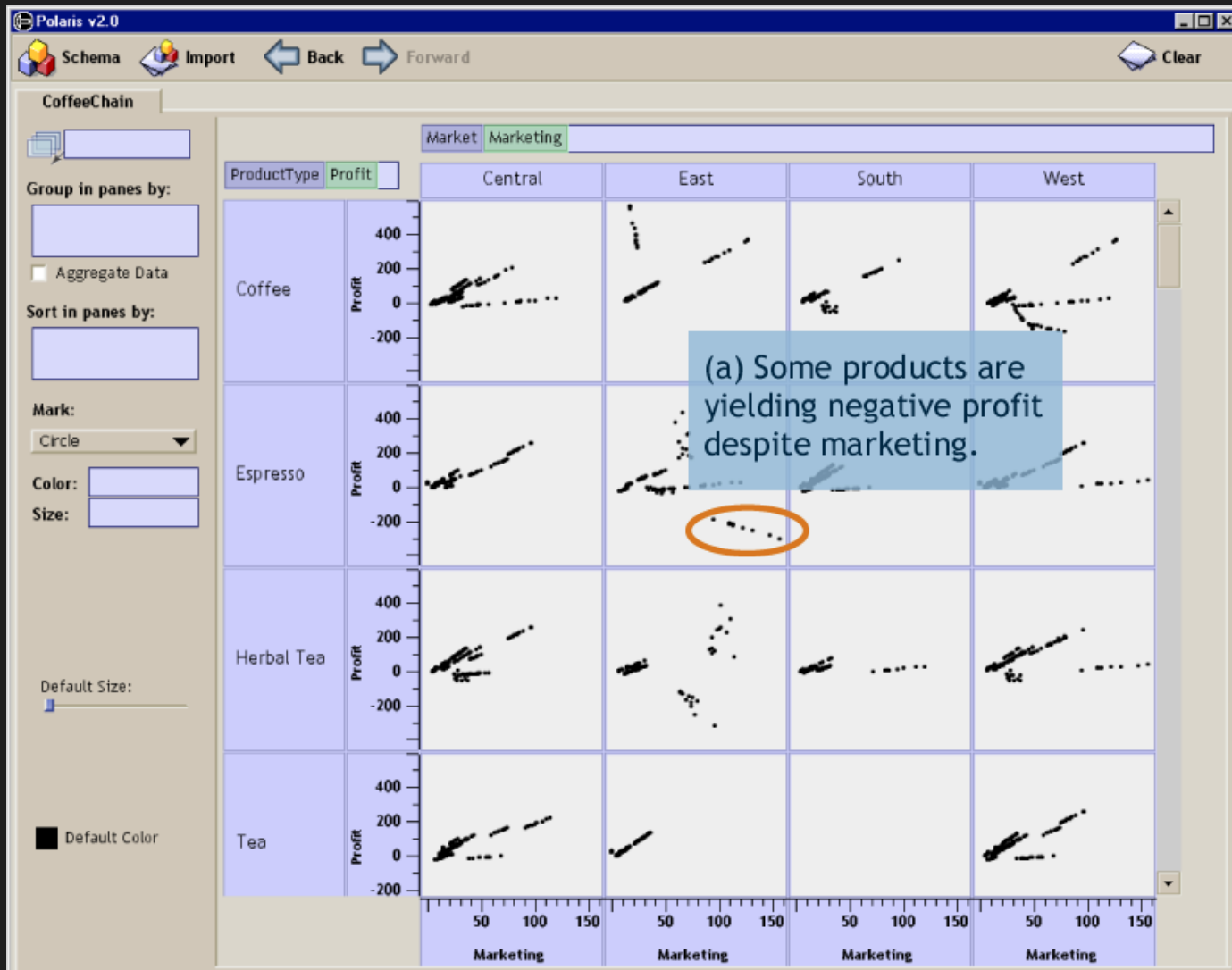


Major wars over the last five hundred years and additional layer of major scientists (country and date of birth)

Ordinal - Quantitative

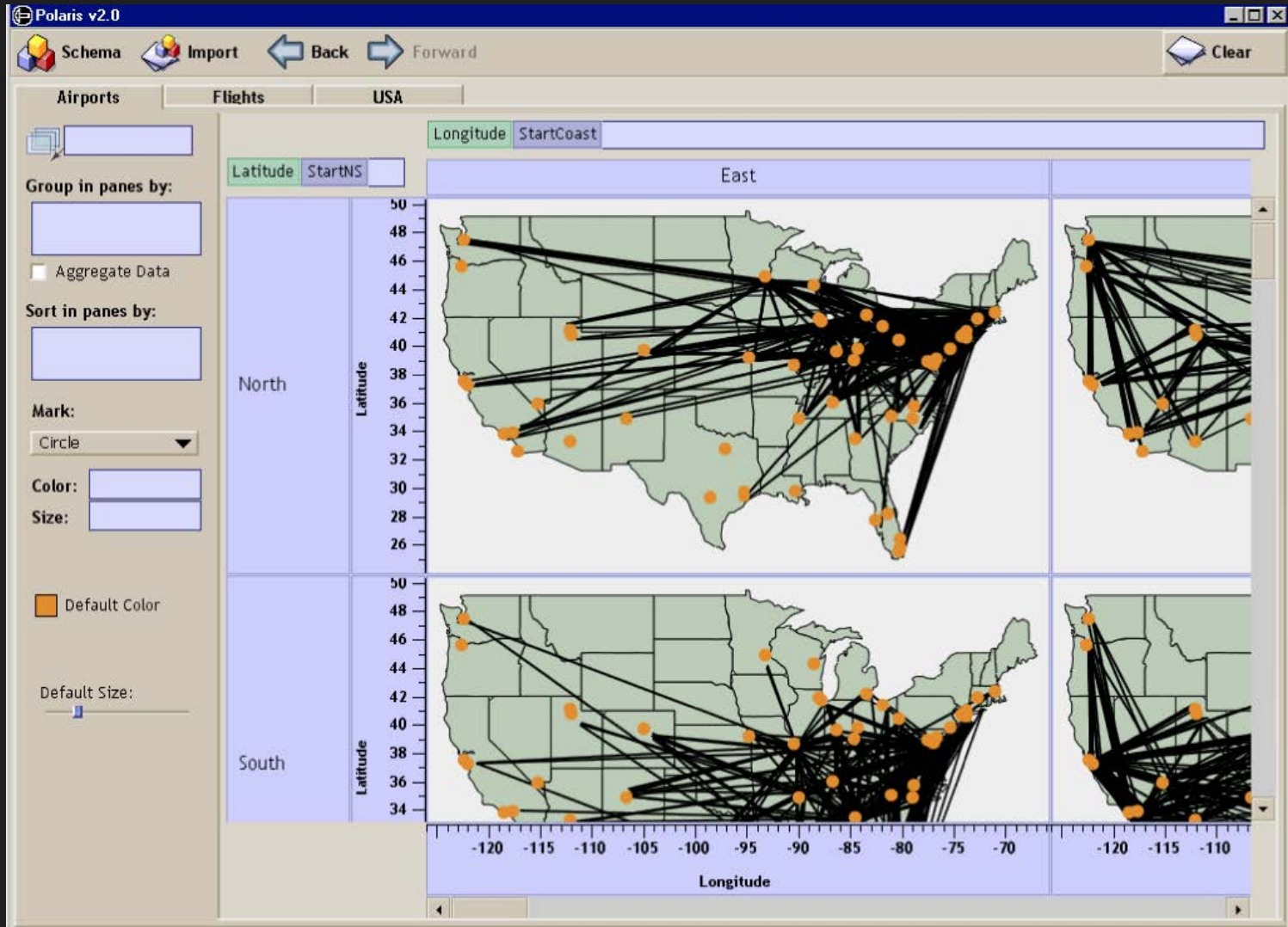


Quantitative - Quantitative



Number of attributes of different products sold by a coffee chain

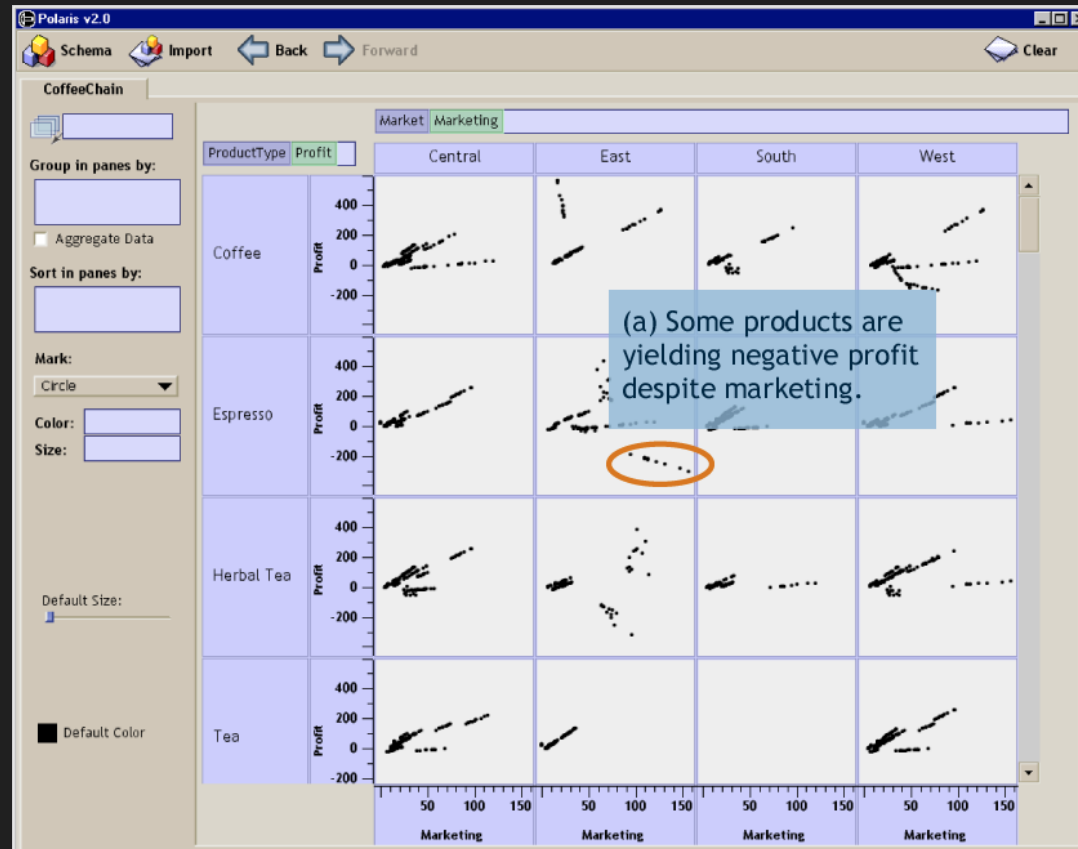
Quantitative - Quantitative



Flight scheduling varies with the region of the country the flight originated in

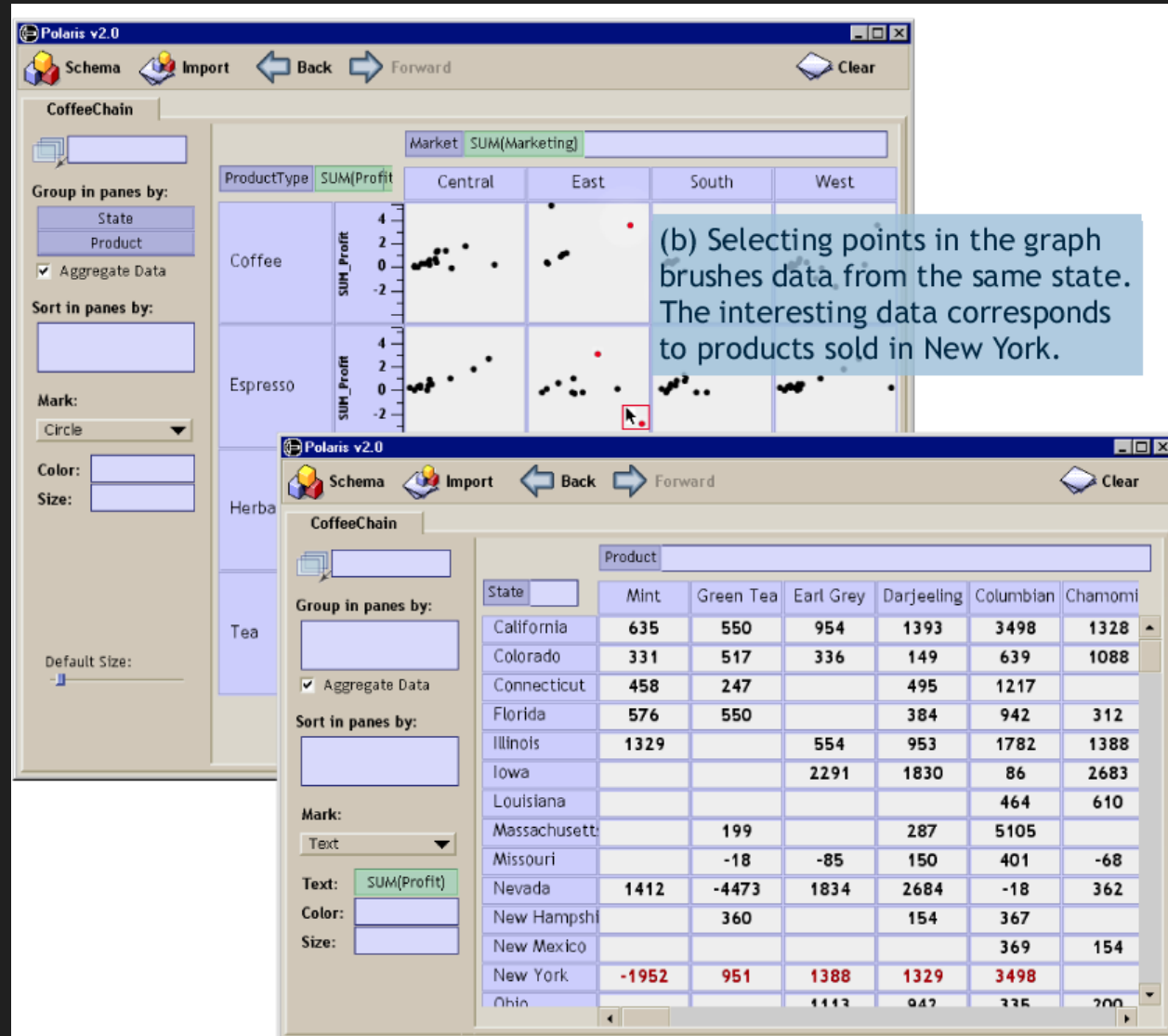
Scenarios

- CFO told to cut expenses – Scatterplots of marketing costs and profit categorized by product type and market



Scenarios

- Further investigate by creating two linked displays
- Scatter plot and text table



Scenarios

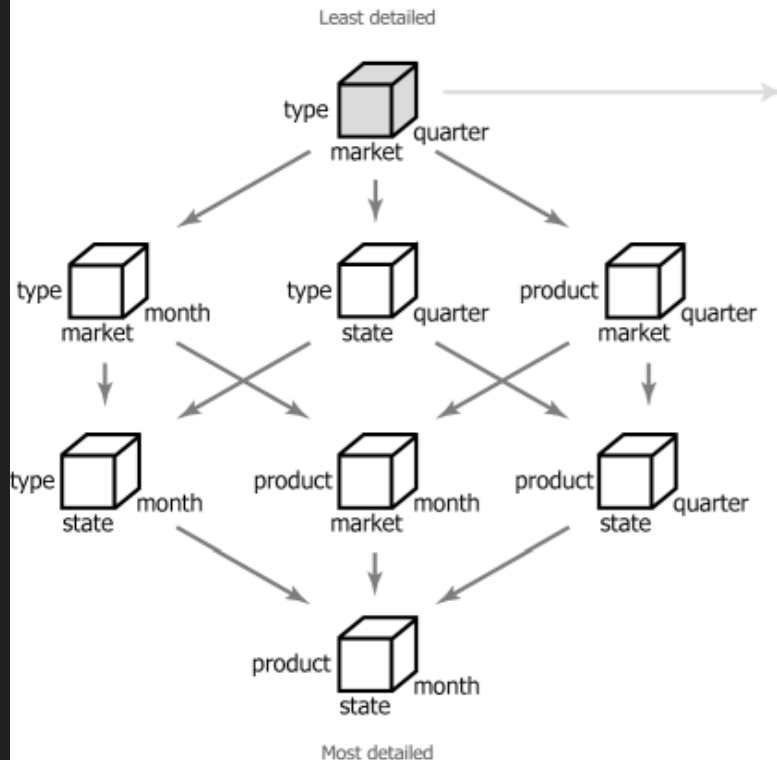
- Third display to visualize data for each product sold in New York
- Café Mocha



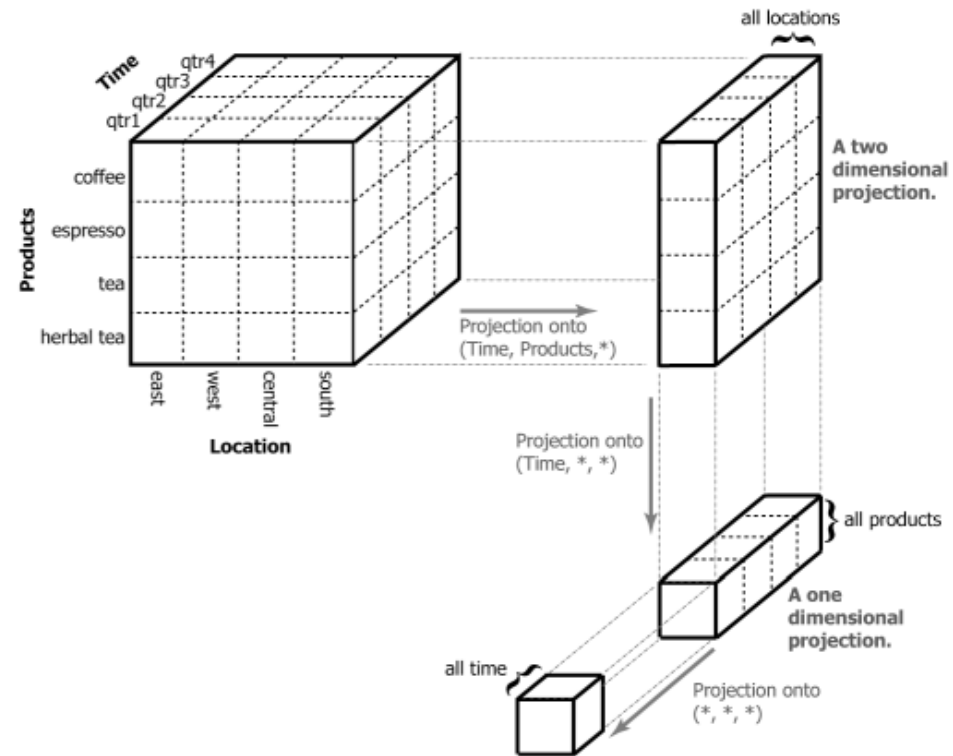
Multiscale Visualizations

- Data Abstraction through Data Cubes

(a) The lattice of data cubes

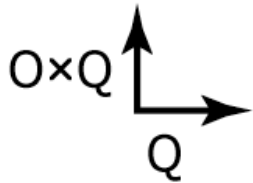


(b) Projecting a three dimensional data cube

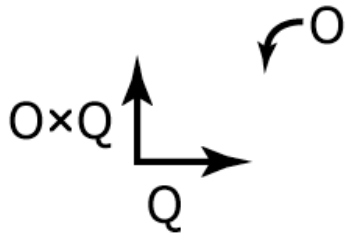


Multiscale Visualizations

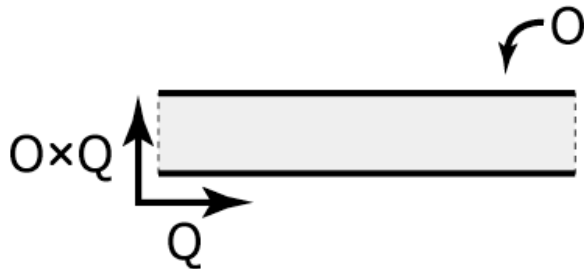
- Visual Abstraction: Polaris



The table expressions are depicted next to two orthogonal axes.
































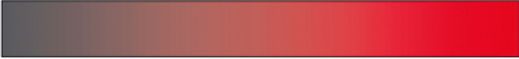
The internal level of detail is shown next to a curved arrow above the axes.



} each layer is shown as a horizontal grey band

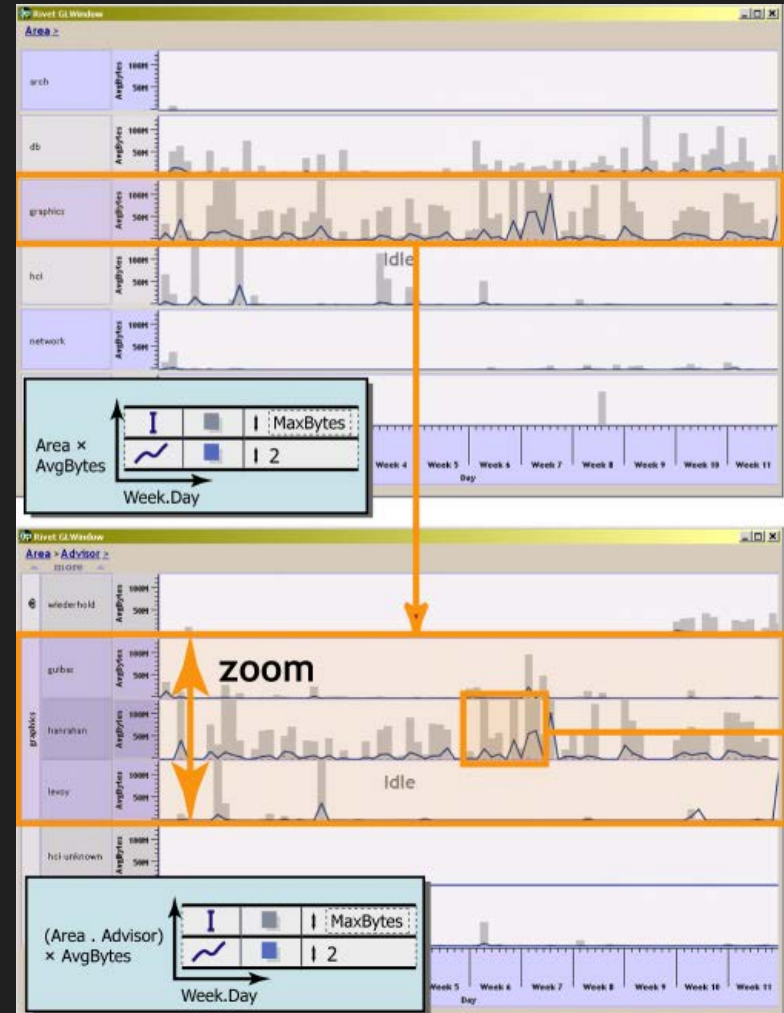
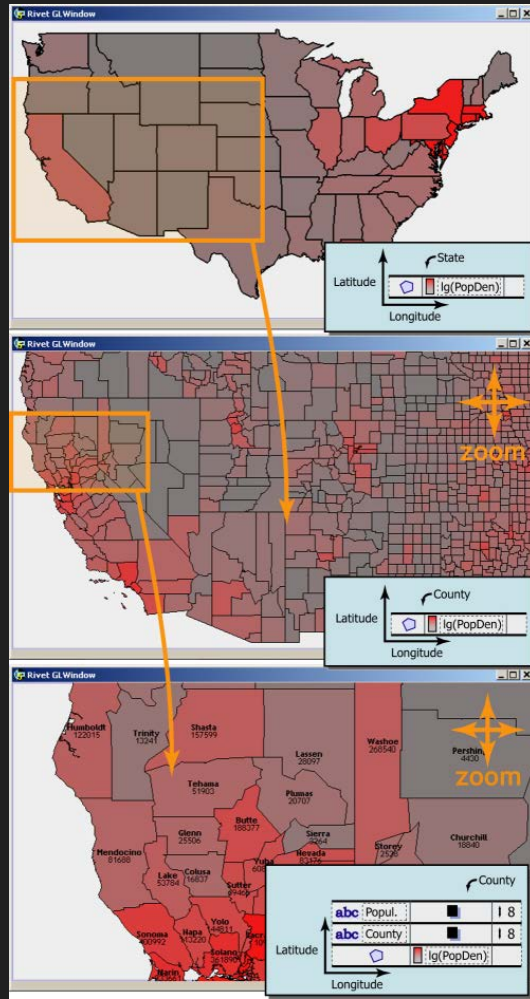
Multiscale Visualizations

- Visual Abstraction: Polaris

property	marks	ordinal/nominal mapping	quantitative mapping
shape	glyph	     	
size	rectangle, circle, glyph, text	   	
orientation	rectangle, line, text	     	
color	rectangle, circle, line, glyph, y-bar, x-bar, text, gantt bar	          ... 	

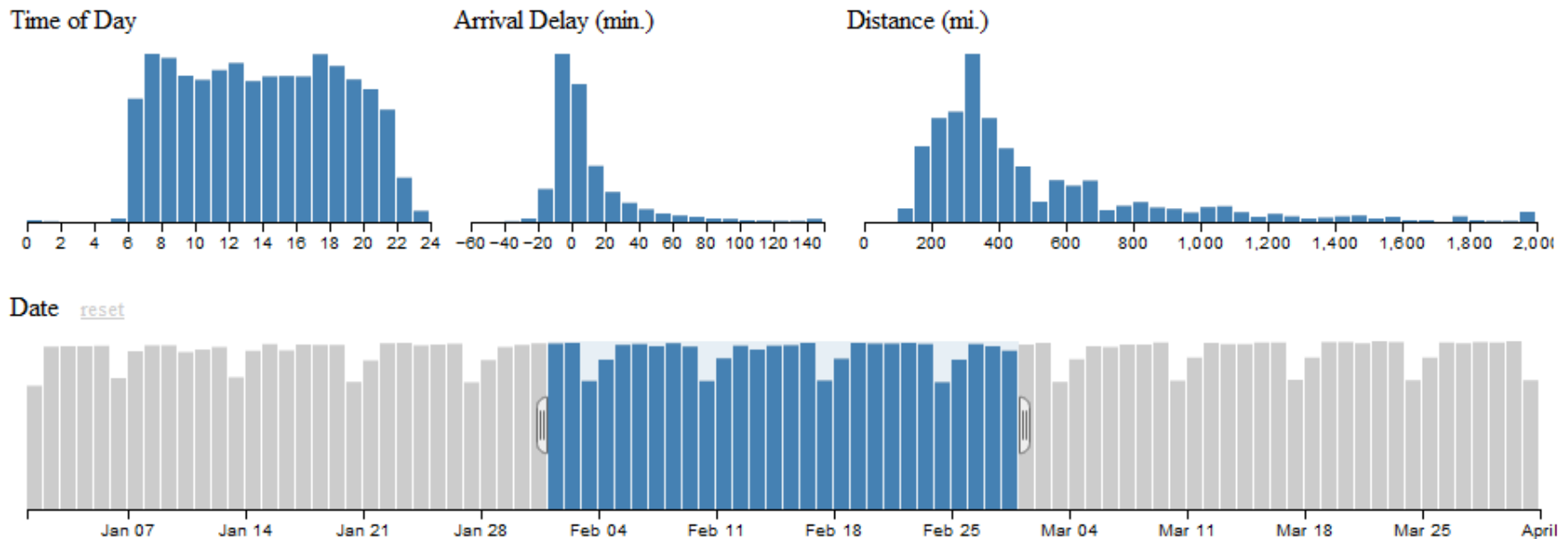
Multiscale Visualizations

- Zoom Graphs:

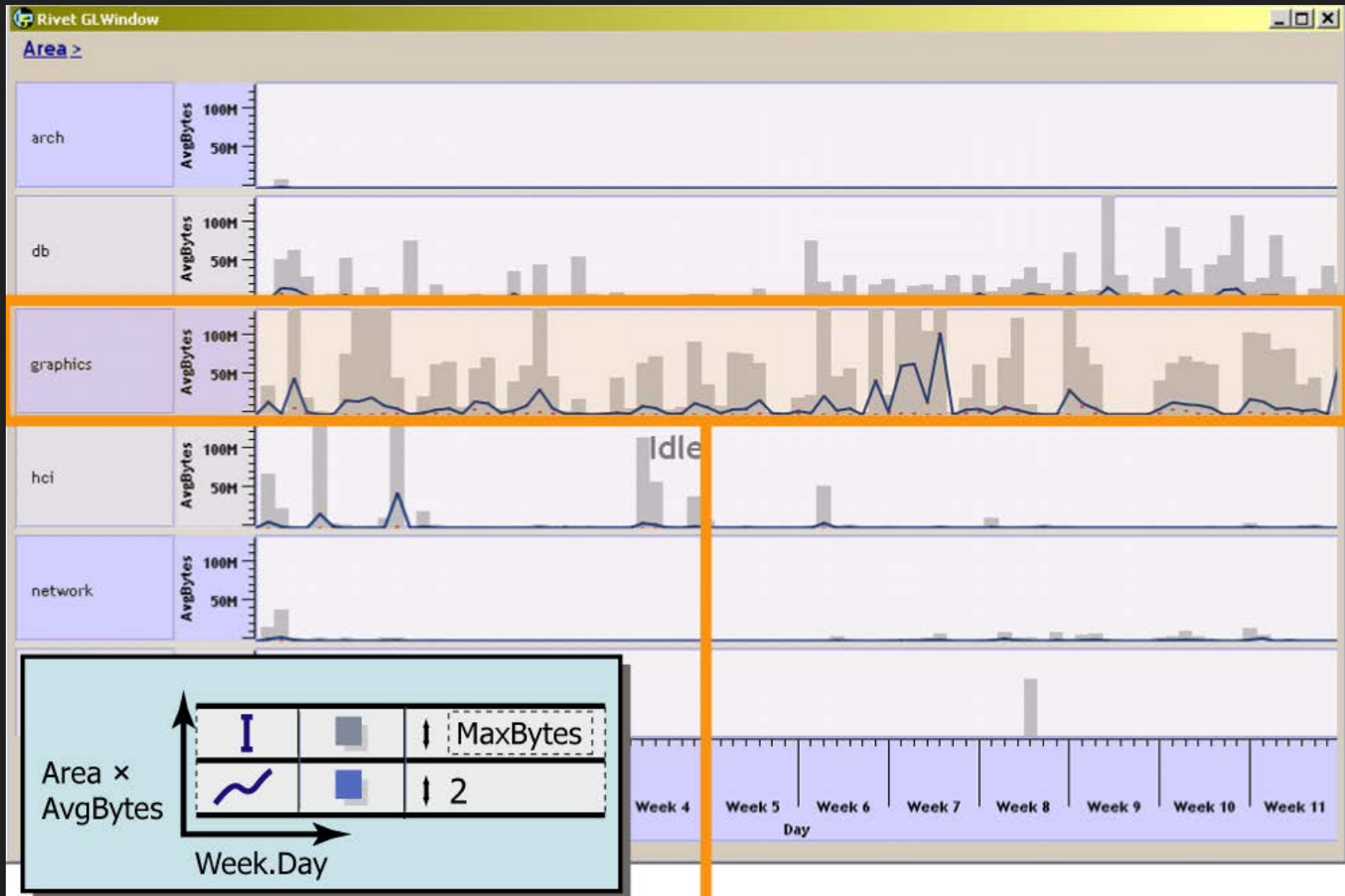


Fast Multidimensional Filtering for Coordinated Views

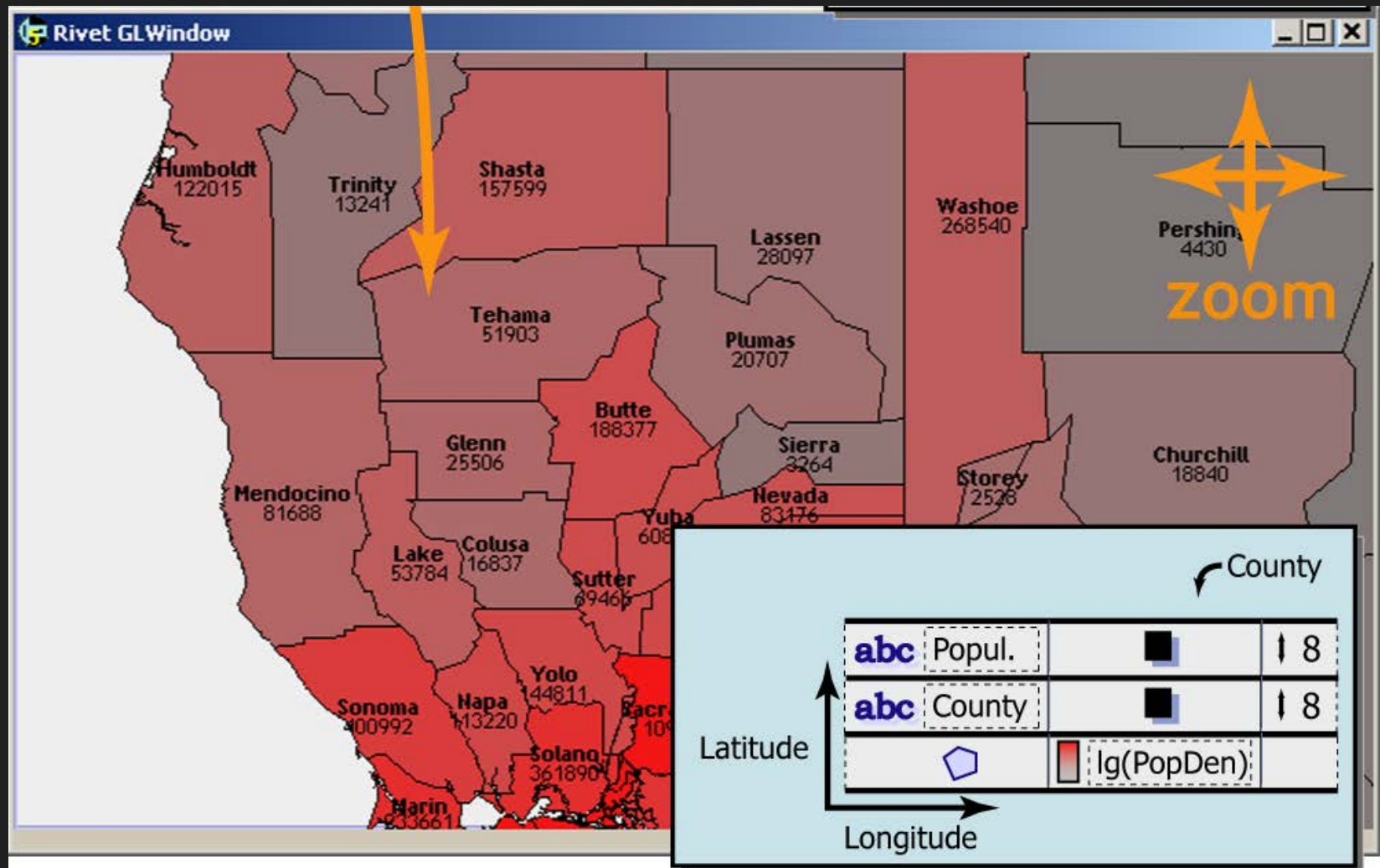
- <http://square.github.io/crossfilter/>



Multiscale Design Patterns: Chart Stacks

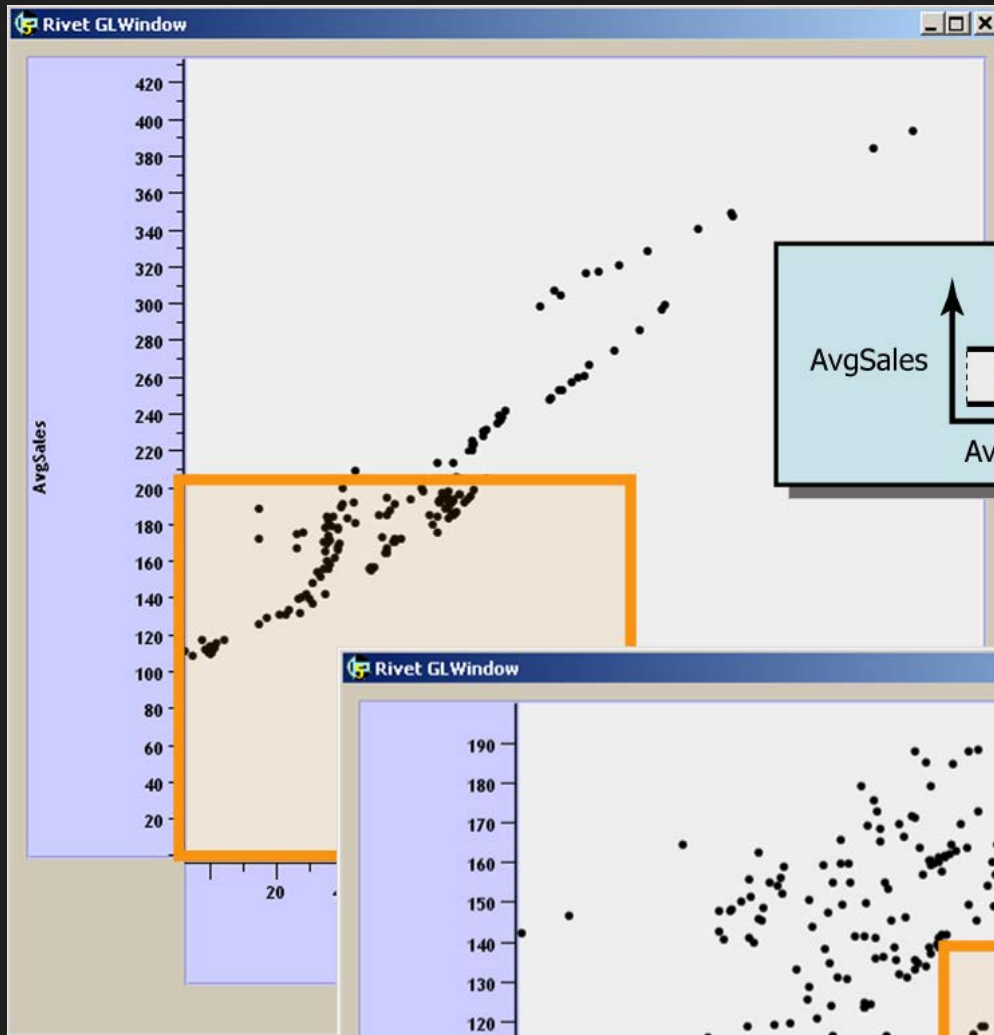


Multiscale Design Patterns: Thematic Maps

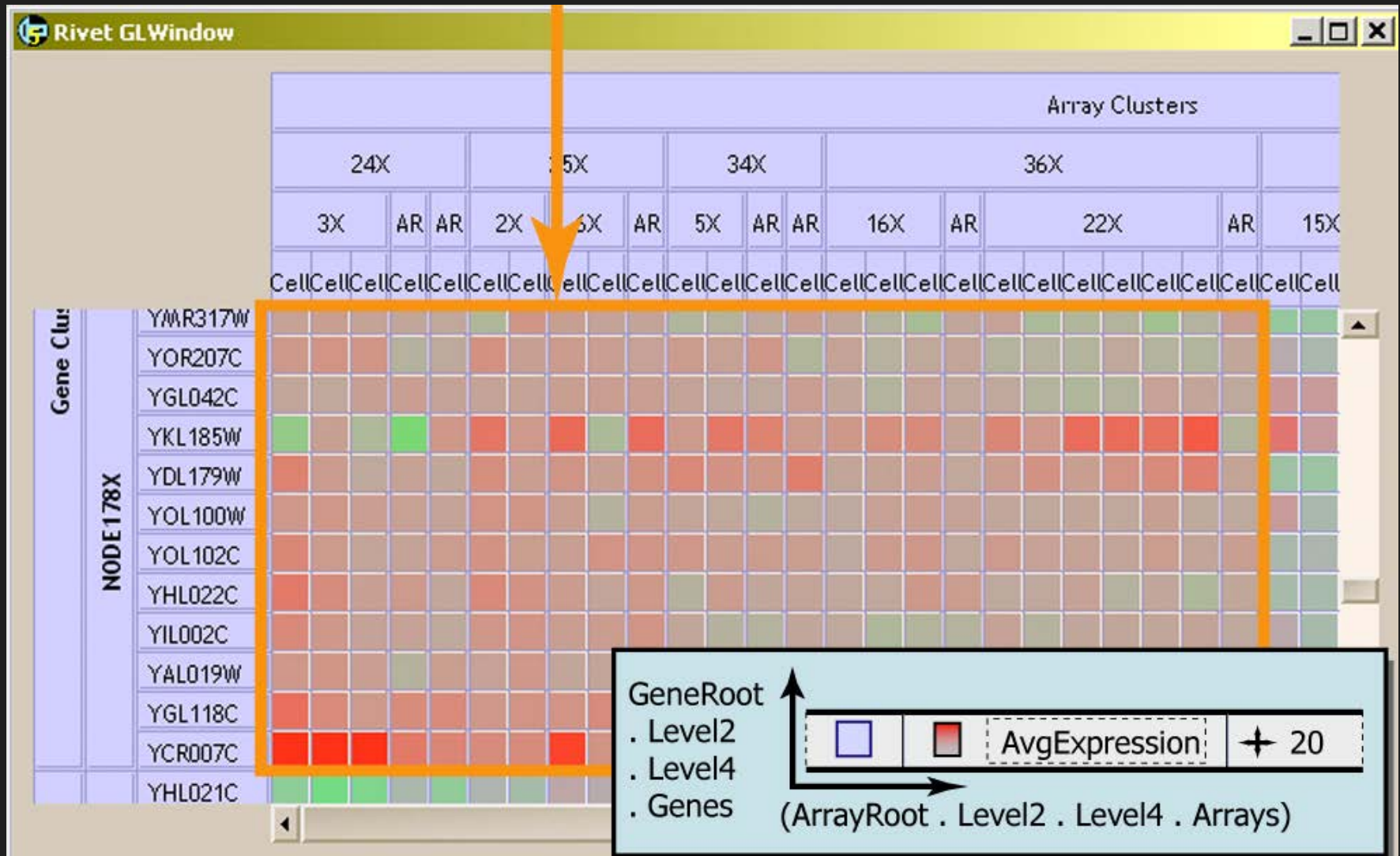


Population density

MDP: Dependent Quantitative-Dependent Quantitative Scatterplots



Multiscale Design Patterns: Matrices



Evaluation of Multivariate Trend Visualization

- Mark A. Livingston, Jonathan W. Decker:
Evaluation of Trend Localization with Multi-
Variate Visualizations. IEEE Trans. Vis. Comput.
Graph. 17(12): 2053-2062 (2011)
- Evaluates the multivariate visualization
techniques