# 07

# **Visualization Techniques Multivariate Data**

# Notice

- **Author**

  ♦ **João Moura Pires (jmp@fct.unl.pt)**

- **This material can be freely used for personal or academic purposes without any previous authorization from the author, provided that this notice is kept with.**

- **For commercial purposes the use of any part of this material requires the previous authorisation from the author.**

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data

- **Data that does not generally have an explicit spatial attribute**

- **Point-Based Techniques**

  - **Project records from an n-dimensional data space to an arbitrary k-dimensional display space, such that data records map to k-dimensional points. (e.g. Scatterplots)**

- **Line-Based Techniques**

  - ◆ **Points corresponding to a particular record or di- mension are linked together with straight or curved lines. (e.g. Line Graphs, Parallel Coordinates)**

- **Region-Based Techniques**

  - ◆ **Filled polygons are used to convey values, based on their size, shape, color, or other attributes. (e.g. Bar Charts/Histograms)**

FACULDADE DE CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Table of Contents

- **Introduction**

- **Point-Based Techniques**

- **Line-Based Techniques**

- **Region-Based Techniques**

- **Combinations of Techniques**

# Introduction

# Multivariate Data

- **Data that does not generally have an explicit spatial attribute**

# Multivariate Data

- **Data that does not generally have an explicit spatial attribute**

- **Point-Based Techniques**

  - **Project records from an n-dimensional data space to an arbitrary k-dimensional display space, such that data records map to k-dimensional points. (e.g. Scatterplots)**

# Multivariate Data

- **Data that does not generally have an explicit spatial attribute**

- **Point-Based Techniques**

  - **Project records from an n-dimensional data space to an arbitrary k-dimensional display space, such that data records map to k-dimensional points. (e.g. Scatterplots)**

- **Line-Based Techniques**

  - **Points corresponding to a particular record or dimension are linked together with straight or curved lines. (e.g. Line Graphs, Parallel Coordinates)**

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data

- **Data that does not generally have an explicit spatial attribute**

- **Point-Based Techniques**

  - **Project records from an n-dimensional data space to an arbitrary k-dimensional display space, such that data records map to k-dimensional points. (e.g. Scatterplots)**

- **Line-Based Techniques**

  - **Points corresponding to a particular record or dimension are linked together with straight or curved lines. (e.g. Line Graphs, Parallel Coordinates)**

- **Region-Based Techniques**

  - **Filled polygons are used to convey values, based on their size, shape, color, or other attributes. (e.g. Bar Charts/Histograms)**

# Point-Based Techniques

# Multivariate Data: Point-Based Techniques

- **Scatterplots** and **Scatterplot Matrices**

  - **Their success stems from our innate abilities to judge relative position within a bounded space**

# Multivariate Data: Point-Based Techniques

- **Scatterplots** and **Scatterplot Matrices**

  - Their **success** stems from our innate **abilities to judge relative position within a bounded space**

- **As the dimensionality of the data increases, the choices for visual analysis consist of:**

  - **dimension subsetting** (user selection or algorithm based suggestion);

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques

- **Scatterplots** and **Scatterplot Matrices**

    - **Their success stems from our innate abilities to judge relative position within a bounded space**

- **As the dimensionality of the data increases, the choices for visual analysis consist of:**

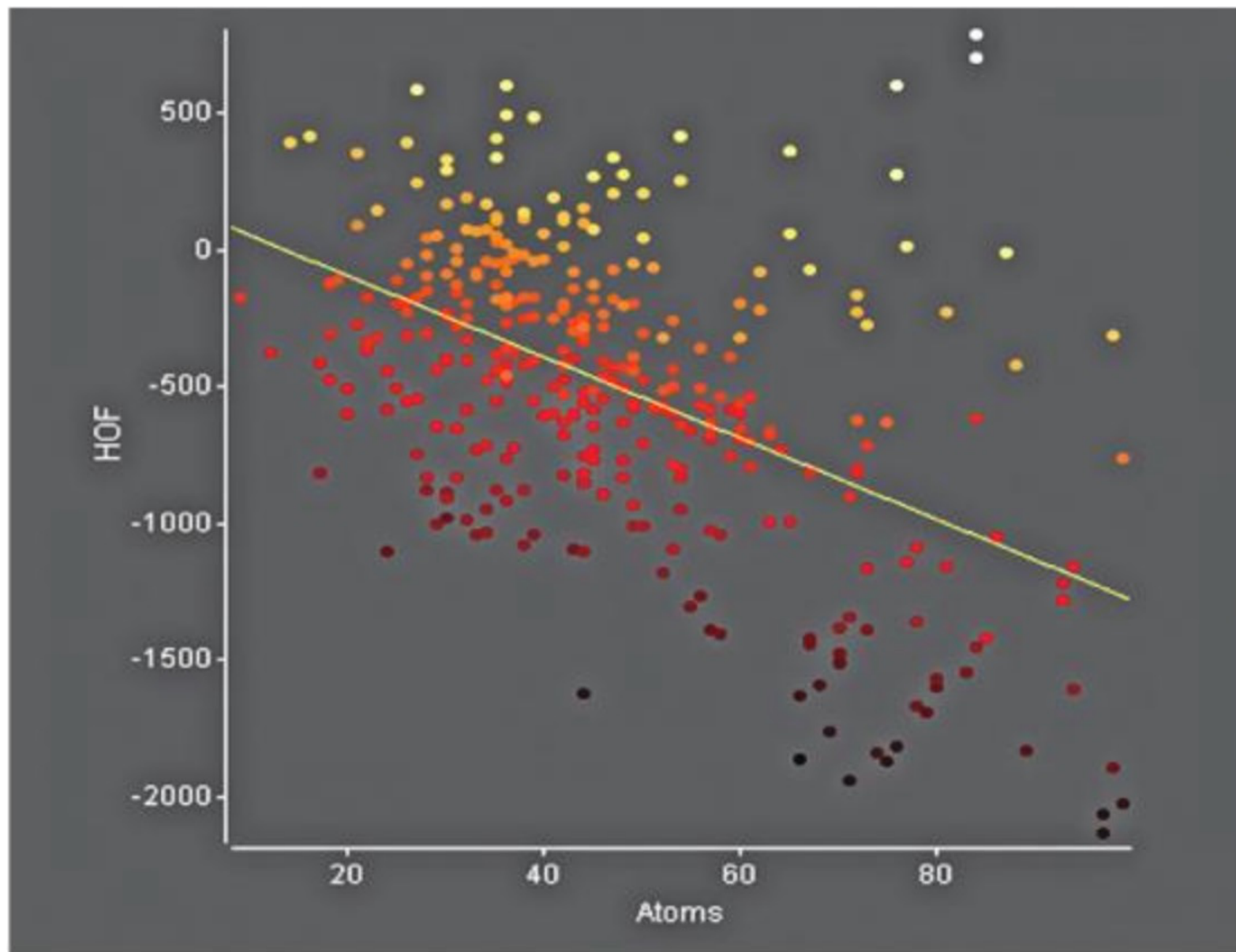    - **dimension subsetting** (user selection or algorithm based suggestion)**;**

    - **dimension embedding** (mapping dimensions to other graphical attributes besides position, such as color, size, and shape)**;**

# Multivariate Data: Point-Based Techniques

- **Scatterplots** and **Scatterplot Matrices**

  - Their **success** stems from our innate **abilities to judge relative position within a bounded space**

- As the **dimensionality** of the data increases, the choices for visual analysis consist of:

  - **dimension subsetting** (user selection or algorithm based suggestion)**;**

  - **dimension embedding** (mapping dimensions to other graphical attributes besides position, such as color, size, and shape)**;**

  - **multiple displays** (either superimposed or juxtaposed - e. g. scatterplot matrix);

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques

- **Scatterplots** and **Scatterplot Matrices**

  - Their **success** stems from our innate **abilities to judge relative position within a bounded space**

- **As the dimensionality of the data increases, the choices for visual analysis consist of:**

  - **dimension subsetting** (user selection or algorithm based suggestion)**;**

  - **dimension embedding** (mapping dimensions to other graphical attributes besides position, such as color, size, and shape)**;**

  - **multiple displays** (either superimposed or juxtaposed - e. g. scatterplot matrix);

  - **dimension reduction** (to transform the high-dimensional data to data of lower dimension).

FACULDADE DE CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques

**Scatterplots**



$x$-coordinate: number of atoms;
$y$-coordinate: heat information;

$y = mx + b$; $m = -12.5$ and $b = 50$

Color of each point: Gibs energy

# Multivariate Data: Point-Based Techniques

- **Scatterplots**

# Multivariate Data: Point-Based Techniques



A scatterplot matrix with the diagonal plot showing a histogram of each dimension.
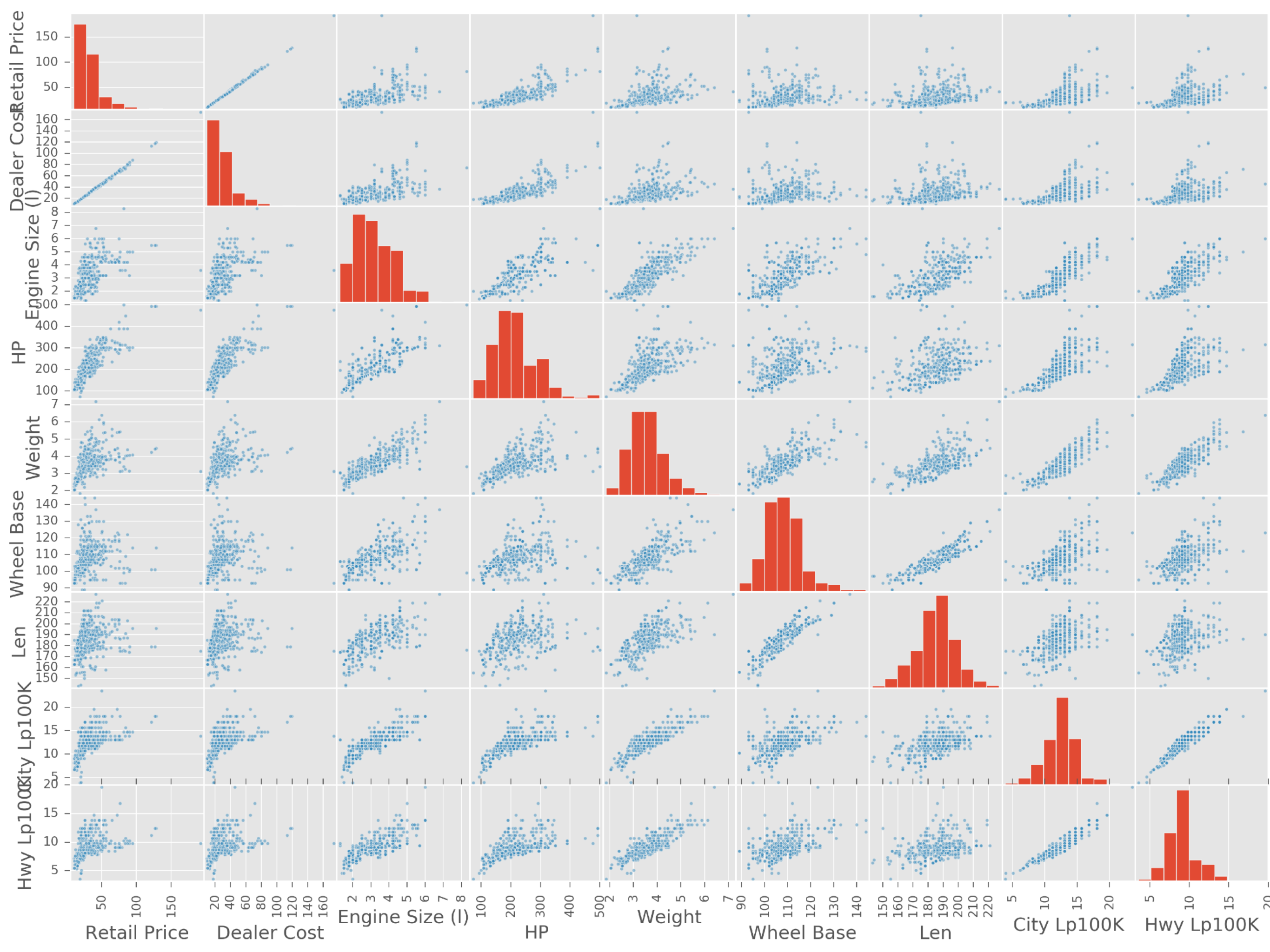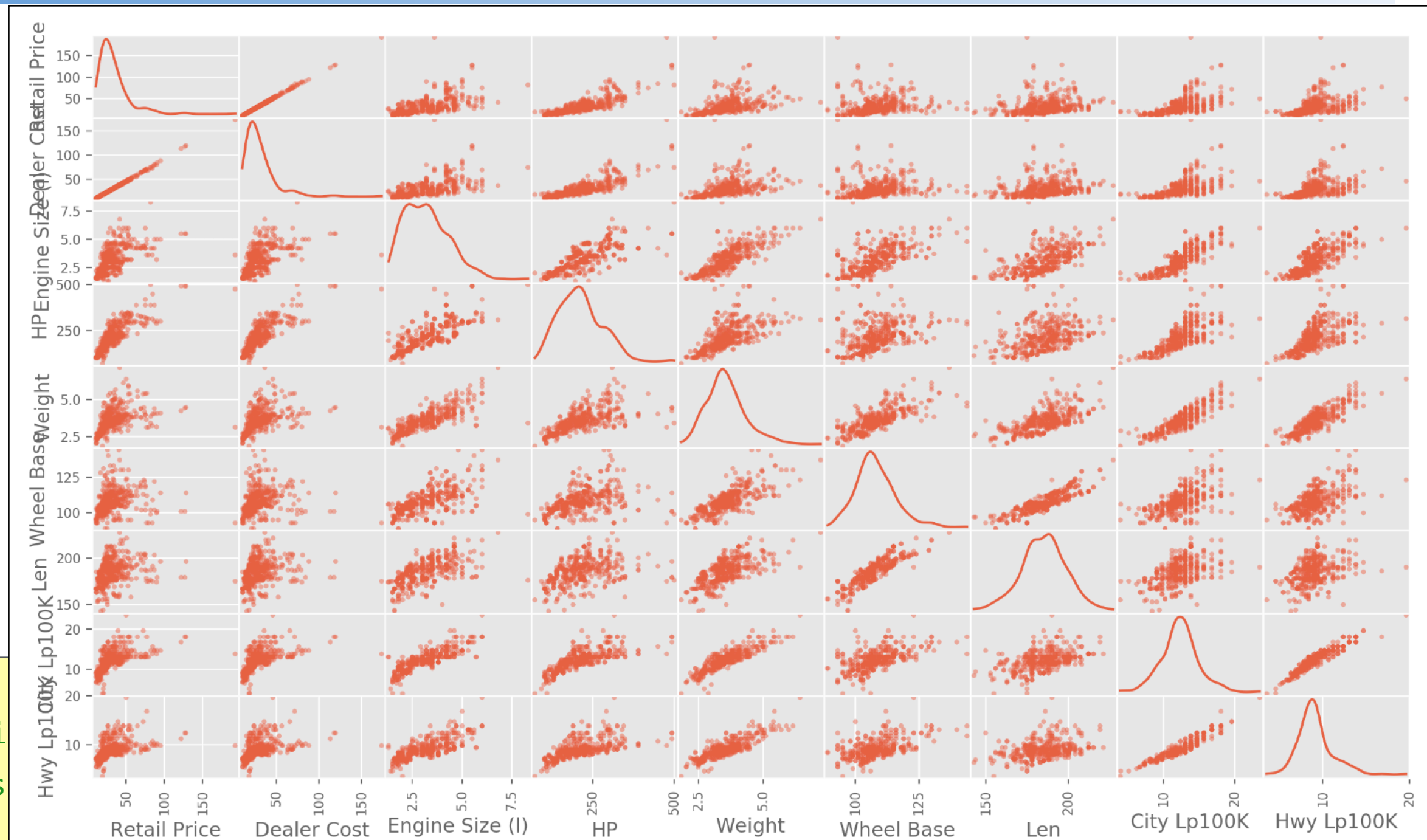Note that the points and histogram regions in red indicate selected data.

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

A scatterplot matrix with the diagonal plot showing a histogram of each dimension. Note that the points and histogram regions in red indicate selected data.

# Scatter Matrix (in Python)



```python
...
# data is the data frame with all variable

# snc is the subset of numerical variables of interest


# Let's check how these variables relate to ecah other
scatter matrix(data[snc],figsize=(12,12))
```

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Scatter Matrix (in Python)



```
...

# data i

# snc is

# Let's check how these variables relate to ecah other
scatter_matrix(data[snc],figsize=(12,12), diagonal='kde')
```

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Scatter Matrix (in Tableau)

# Scatter Matrix (in Tableau)

# Matrix of Scatter Matrix



Class
- Minivcan
- Normal
- Pickup
- Sports
- SUV
- Wagon

>26 nulls

# Matrix of Scatter Matrix



Class
- Minivcan
- Normal
- Pickup
- Sports
- SUV
- Wagon

>26 nulls

# Multivariate Data: Point-Based Techniques

- **In situations where the dimensionality of the data exceeds the capabilities of the visualization technique. It is necessary to investigate ways to reduce the data dimensionality, while at the same time preserving, as much as possible, the information contained within.**

- **Principal Component Analysis (PCA) -** read more **and see this** implementation

- **Multidimensional Scaling (MDS) -** read more **and** more

- **Non-linear dimension reduction techniques:**

  - **Self-organizing Maps (SOMs) -** read more

  - **Local Linear Embeddings (LLE) -** read more

FACULDADE DE CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multidimensional scaling (MDS)

■ Projecting **M** points in **N** dimensions into **L** dimensions (L = 2 or 3) display space.

# Multidimensional scaling (MDS)

- Projecting **M** points in **N** dimensions into **L** dimensions (L = 2 or 3) display space.

- The key goal is to **attempt to maintain the N-dimensional features and characteristics** of the data through the projection process, e.g., relationships that exist in the original data must also exist after projection.

# Multidimensional scaling (MDS)

- Projecting **M** points in **N** dimensions into **L** dimensions (L = 2 or 3) display space.

- The key goal is to **attempt to maintain the N-dimensional features and characteristics** of the data through the projection process, e.g., relationships that exist in the original data must also exist after projection.

  - The projection may also **unintentionally introduce artifacts** that may appear in the visualization and are not present in the data.

# Multidimensional scaling (MDS)

- Projecting **M** points in **N** dimensions into **L** dimensions (L = 2 or 3) display space.

- The key goal is to **attempt to maintain the N-dimensional features and characteristics** of the data through the projection process, e.g., relationships that exist in the original data must also exist after projection.

  - The projection may also **unintentionally introduce artifacts** that may appear in the visualization and are not present in the data.

- Repeat

  - Create an Similarity **M x M** Matrix (**D**) (could be distance)

  - Create a coordinates Matrix **M** x **L** and fill randomly or other method (ex: PCA)

  - Compute an M x M matrix (**L**) based on L coordinates. And compute **S** the difference between **D** and **L**.

  - Shift the positions of points in L in a direction that will reduce their individual stress levels

- Until **S** is small of not changed significantly

# Multidimensional scaling (MDS)

■ There are many possible variants on this algorithm, including:

  ◆ **Different similarity** and **stress measures**;

  ◆ **Different initial** and **termination conditions**;

  ◆ **Different position update strategies**.

# Multidimensional scaling (MDS)

- There are many possible variants on this algorithm, including:

    - **Different similarity** and **stress measures**;

    - **Different initial** and **termination conditions**;

    - **Different position update strategies**.

- As in any optimization process, there is the potential to fall into a local minimal configuration that still has a high level of stress.

# Multidimensional scaling (MDS)

- There are many possible variants on this algorithm, including:

  - **Different similarity** and **stress measures**;

  - **Different initial** and **termination conditions**;

  - **Different position update strategies**.

- As in any optimization process, there is the potential to fall into a local minimal configuration that still has a high level of stress.

  - Common strategies to alleviate this include occasionally **adding a random jump** in the position of a point to see if it will converge to a different location

# Multidimensional scaling (MDS)

■ There are many possible variants on this algorithm, including:

  ♦ **Different similarity** and **stress measures**;

  ♦ **Different initial** and **termination conditions**;

  ♦ **Different position update strategies**.

♦ As in any optimization process, there is the potential to fall into a local minimal configuration that still has a high level of stress.

  ♦ Common strategies to alleviate this include occasionally **adding a random jump** in the position of a point to see if it will converge to a different location

♦ Obviously, **the results are not unique**: minor changes in the starting conditions can lead to dramatically different results.

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques

- **Iris flower data set**



Iris setosa

Iris versicolor

Iris virginica

FACULDADE DE CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques



Iris Data (red=setosa,green=versicolor,blue=virginica)

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Point-Based Techniques

■ **Iris data set projected using MDS**

# Multivariate Data: Point-Based Techniques

- **RadViz:** is a force-driven point layout technique that is based on Hooke's Law for equilibrium.

- For an N-dimensional data set, **N anchor points** are placed on the circumference of the circle to represent the fixed ends of the **N springs** attached to each data point.

- **Different placement and ordering of the anchors will give different results**, and that points that are quite distinct in N dimensions may map to the same location in 2D.

# Multivariate Data: Point-Based Techniques

- **RadViz:** different views of the same data set in RadViz, using manual reordering of dimensions.

# Multivariate Data: Point-Based Techniques

■ **RadViz:** different views of the same data set in RadViz, using manual reordering of dimensions.

# Multivariate Data: Point-Based Techniques

- **RadViz:** different views of the same data set in RadViz, using manual reordering of dimensions.

# Multivariate Data: Point-Based Techniques

■ **Vectorized RadViz, or VRV,** constructs multiple dimensions from individual dimensions by a

flattening process, breaking each dimension into many



Vectorized RadViz, formed by splitting each dimension into multiple dimensions to create a binary representation for each data record. In this case, each cluster set is separated into multiple dimensions, where each dimension represents a cluster in each cluster set [372].

# Multivariate Data: Point-Based Techniques

Dimension representing the number of cylinders can be broken down into 5 new dimensions:
- having 1 or 2 cylinders;
- having 3 or 4 cylinders;
- having 5 or 6;
- having 7;
- having 8.



Vectorized RadViz, formed by splitting each dimension into multiple dimensions to create a binary representation for each data record. In this case, each cluster set is separated into multiple dimensions, where each dimension represents a cluster in each cluster set [372].

FACULDADE DE CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Line-Based Techniques

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Line-Based Techniques



(a) superimposed

**Line Graphs**

(b) stacked

(c) ordered superimposed

(d) ordered stacked

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Line-Based Techniques

- **Andrews curves**

$$f(t) = \frac{d_1}{\sqrt{2}} + d_2 \sin(t) + d_3 \cos(t) + d_4 \sin(2t) + d_5 \cos(2t) + \dots.$$



(a)                                             (b)

An example of Andrews curves using two different dimension orders: (a) based on the original order of the dimensions (sepal length, sepal width, petal length, petal width); (b) based on the original order of the dimensions in reverse order.

# Multivariate Data: Line-Based Techniques

- **Parallel Coordinates**



An example of a 7-dimensional data set visualized with parallel coordinates. A single data point is represented as the darkened polyline.

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Parallel Coordinates (||-coords or PCP)

- **Inselberg in 1985**



Figure 3: Constructing parallel coordinates with five dimensions represented by $N = 5$ vertical lines. Points in the plane are represented by lines joining the corresponding coordinates at the respective axes. Typically, only the line segments between the axes are drawn (represented by the bold polyline).

State of the Art of Parallel Coordinates
J. Heinrich and D. Weiskopf

# Parallel Coordinates (||-coords or PCP)



Figure 4: The line with slope $m = 1$ in the data domain is mapped to the ideal point $\bar{\ell}_\infty$ in parallel coordinates (top). The vertical line $\overline{P}_m^\infty : x = \frac{d}{1-m}$ in parallel coordinates is represented by the ideal point $P_m^\infty$ with slope $m$ in the data domain. Both domains are considered projective planes.

State of the Art of Parallel Coordinates
J. Heinrich and D. Weiskopf

# Parallel Coordinates (||-coords or PCP)



Figure 5: Common patterns in Cartesian coordinates (top) and their dual representation in parallel coordinates (bottom). The envelope of lines is highlighted for the ellipse–hyperbola duality.

State of the Art of Parallel Coordinates
J. Heinrich and D. Weiskopf

# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)
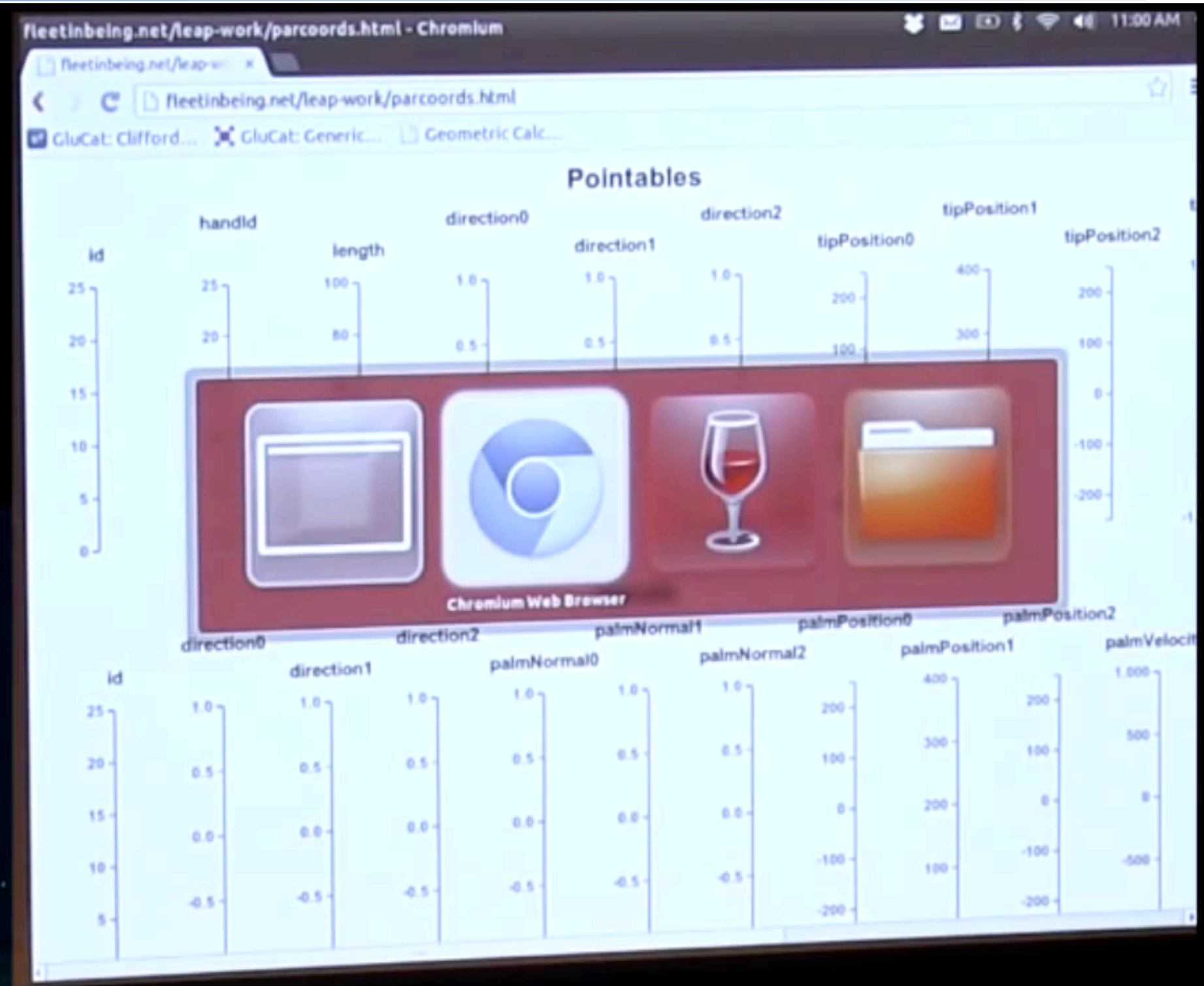


Kai Chang
Visually Exploring Multidimensional Data

# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)

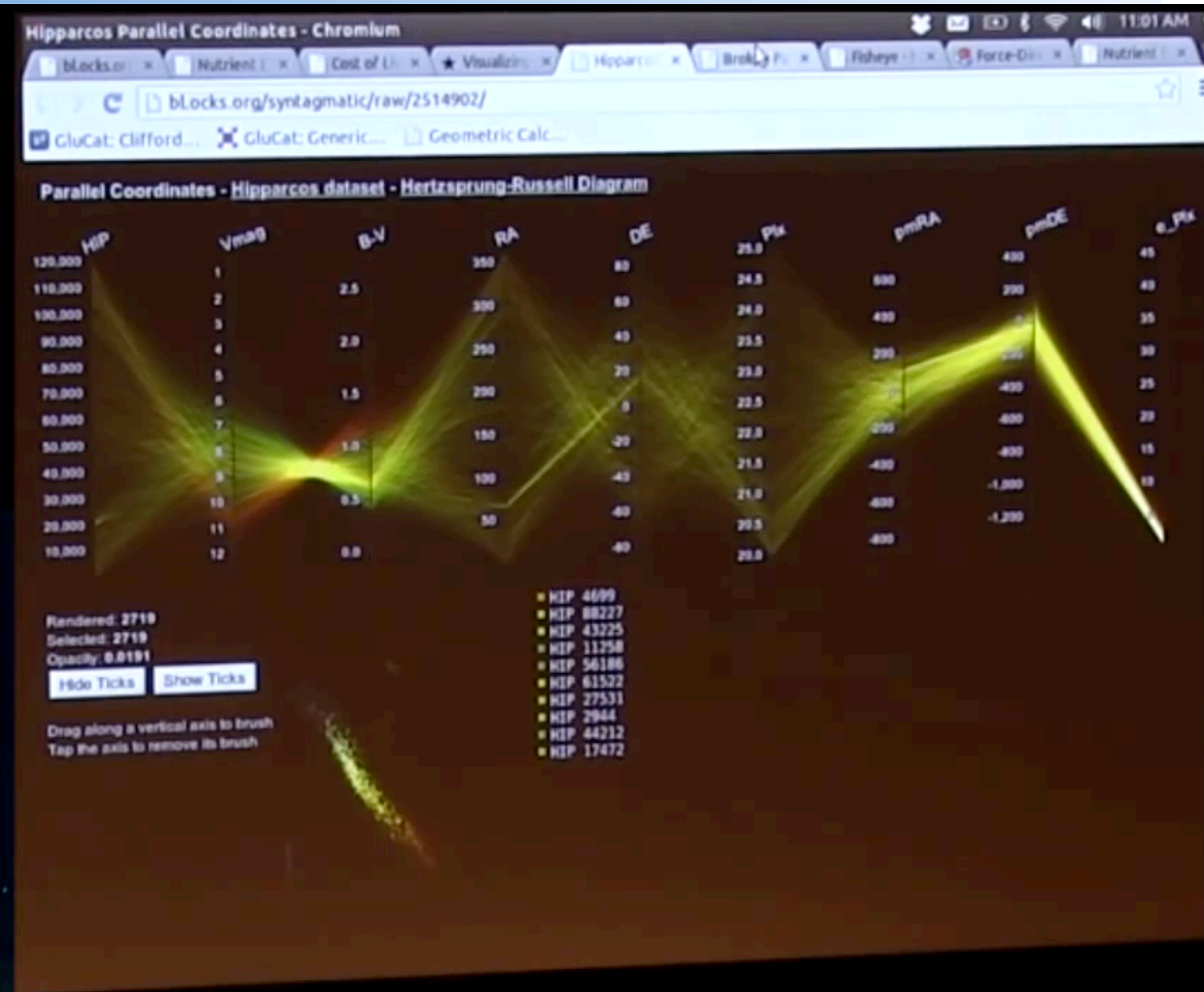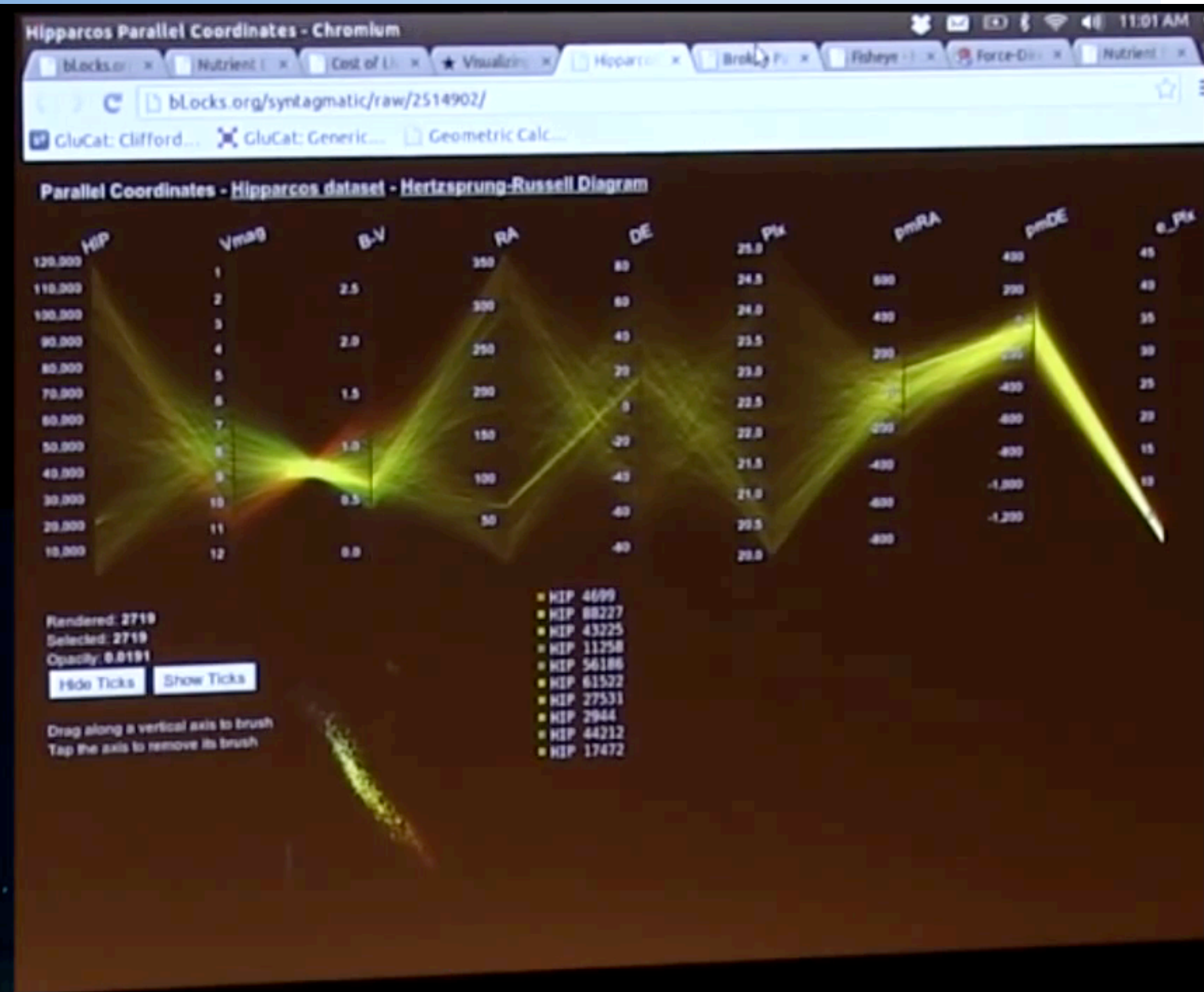# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)

# Parallel Coordinates (||-coords or PCP)

- **Check https://eagereyes.org/techniques/parallel-coordinates**

- **Check https://syntagmatic.github.io/parallel-coordinates/**

- **See the video: https://youtu.be/ypc7Ul9LkxA**

- **http://www.xdat.org/**

- **Check http://www.parallelcoordinates.de/paco/#**

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Parallel Coordinates (||-coords or PCP)

- **Very special videos !**

- **Tutorial by Alfred Inselberg at <u>iV 2016</u> (at Lisbon) (<u>FB</u> and <u>Twitter</u>)**

  - **<u>Part1</u>**

  - **<u>Part2</u>**
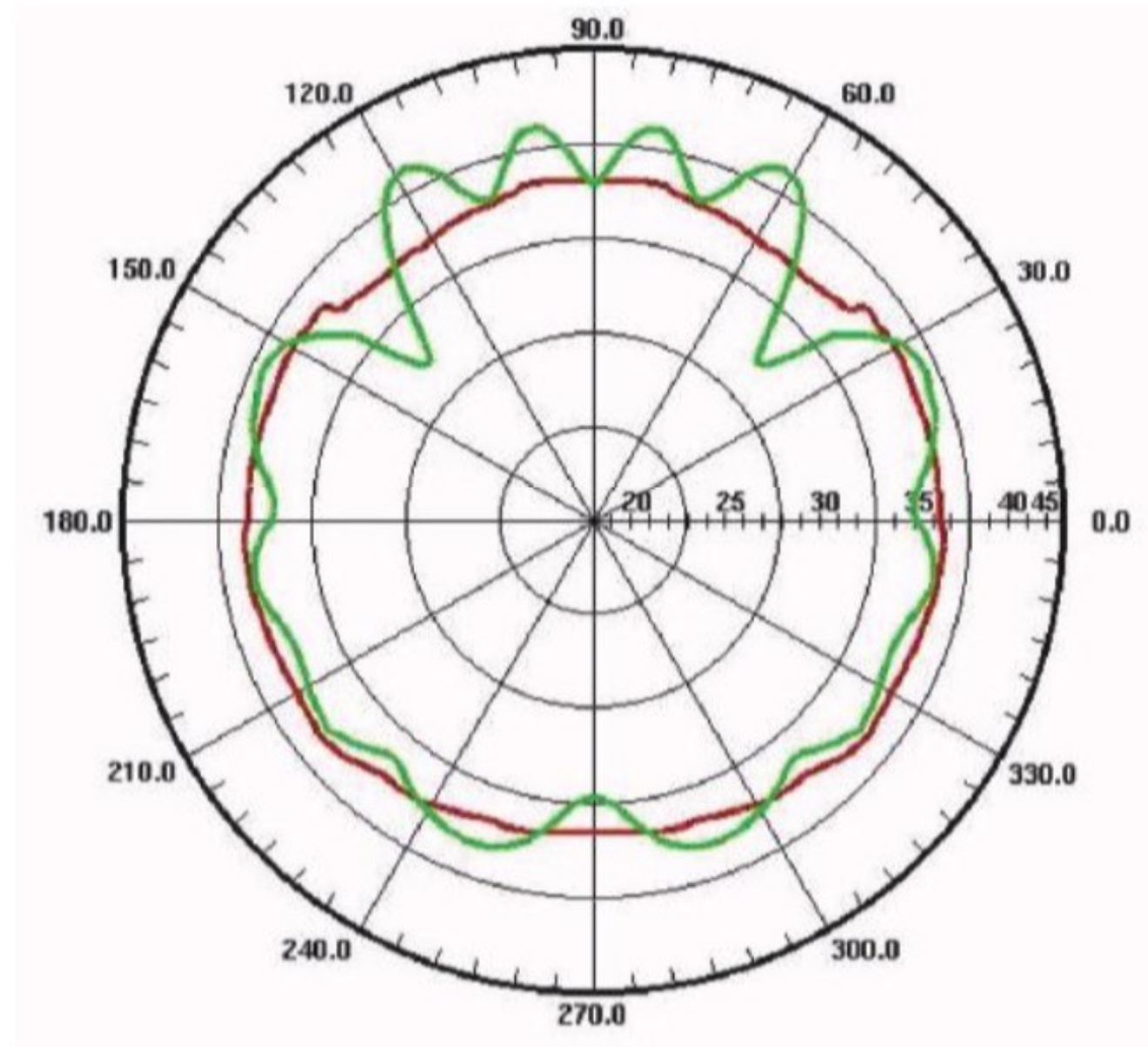
  - **<u>Part3</u>**

State of the Art of Parallel Coordinates
J. Heinrich and D. Weiskopf

# Multivariate Data: Line-Based Techniques

- **Radial Axis Techniques**

  - **circular line graph**;

  - **polar graphs**: point plots using polar coordinates;

  - **circular bar charts**: like circular line graphs, but plotting bars on the base line;

  - **circular area graphs**: like a line graph, but with the area under line filled in with a color or texture;

  - **circular bar graphs**: with bars that are circular arcs with a common center point and base line.
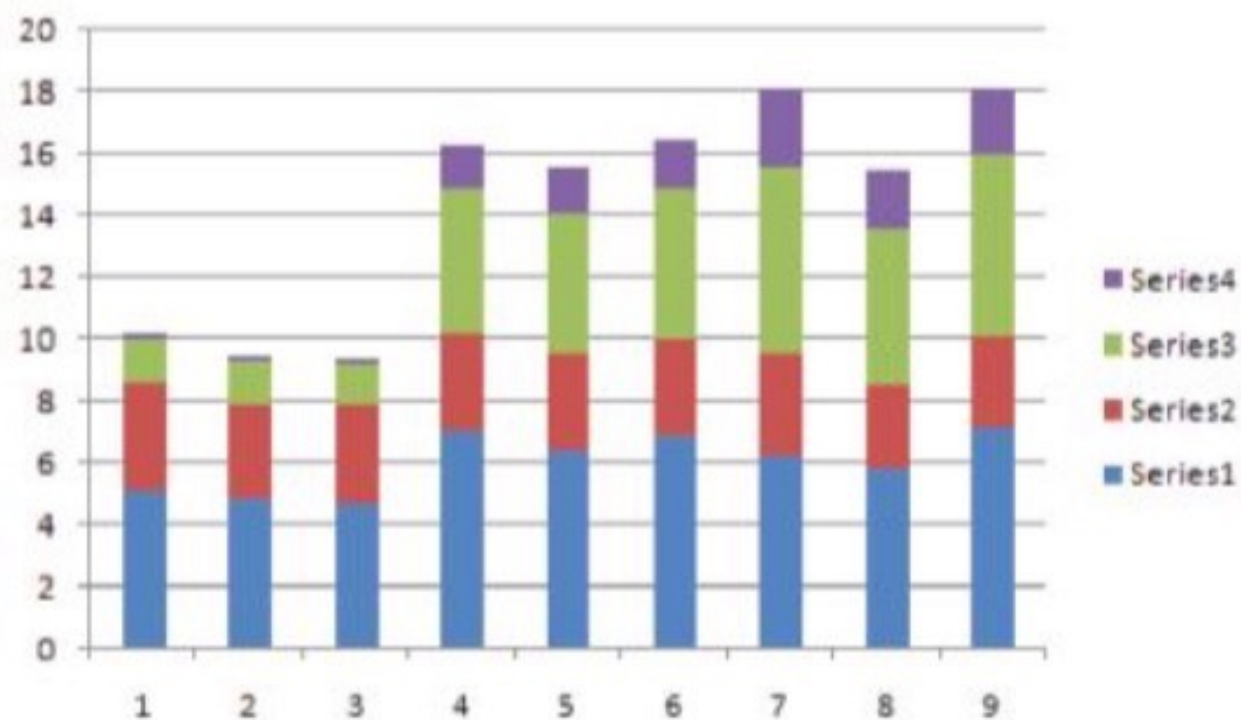
FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Line-Based Techniques



An example of a circular line graph. (Image courtesy http://www.cemframework .com/img/PolarPlot1.png.)
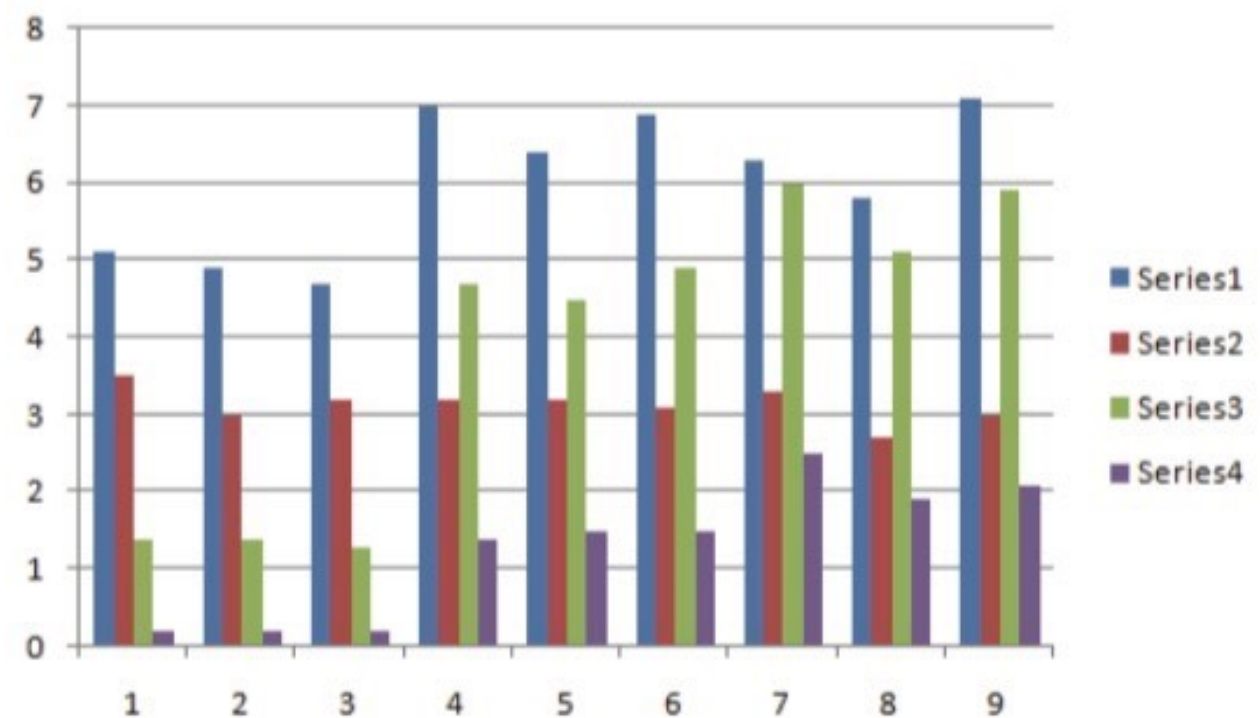
# Region-Based Techniques

# Multivariate Data: Region-Based Techniques
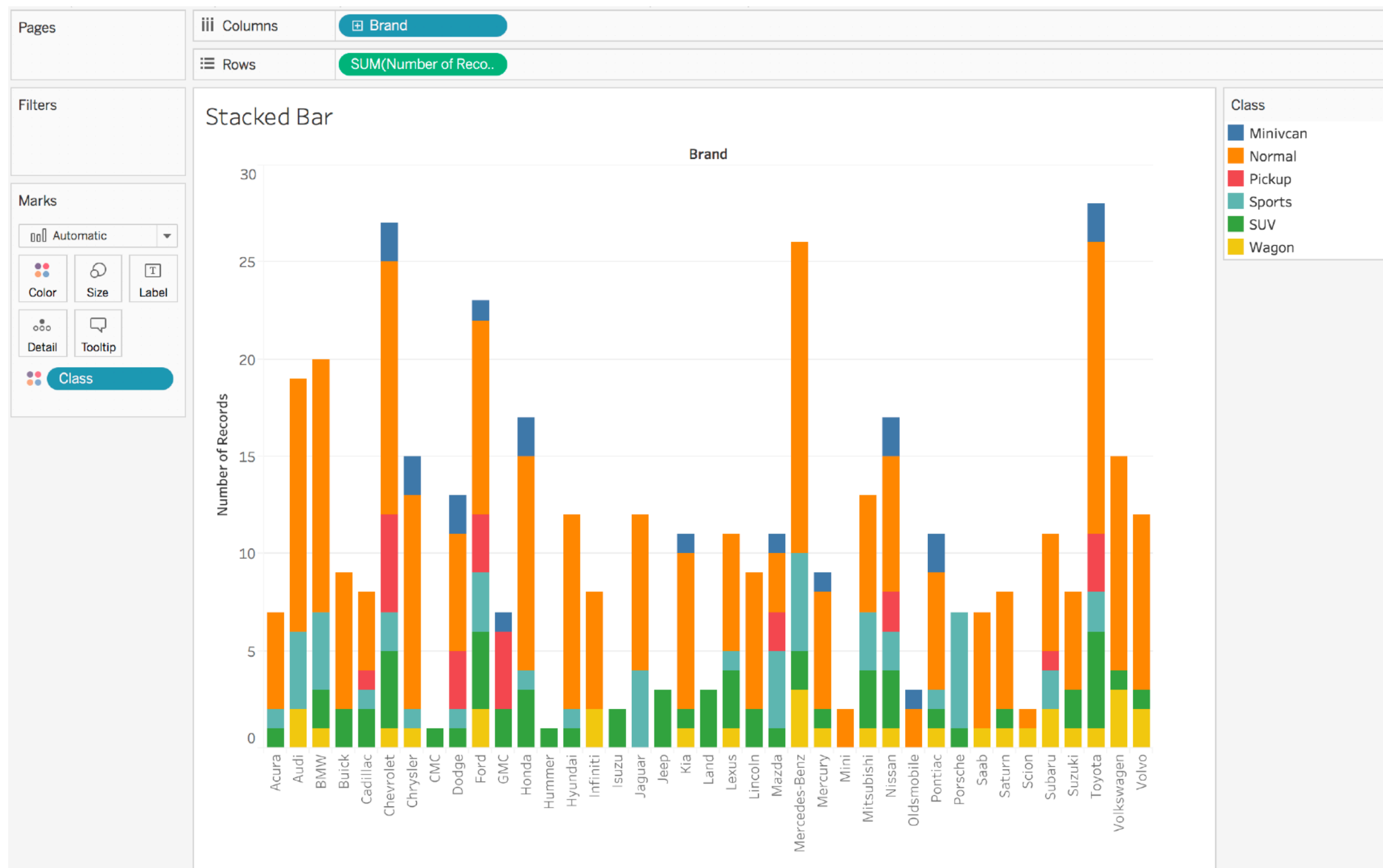
- **Bar Charts/Histograms**



(a) Stacked bar chart.



(b) Clustered bar chart.

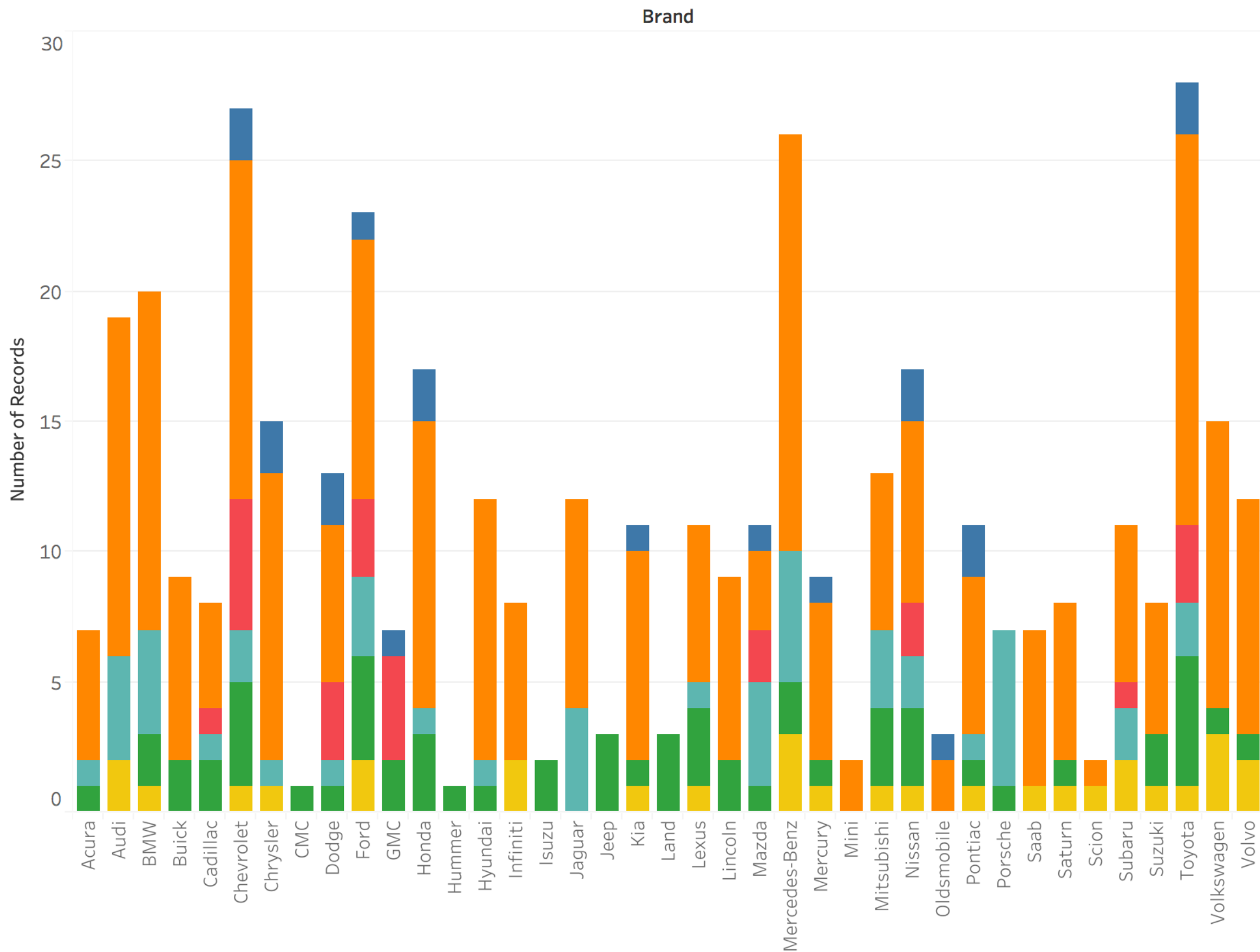# Multivariate Data: Region-Based Techniques

- **Bar Charts**

**Columns** | ⊞ Brand

**Rows** | SUM(Number of Reco..

## Stacked Bar

**Brand**

Class

- ■ Minivcan
- ■ Normal
- ■ Pickup
- ■ Sports
- ■ SUV
- ■ Wagon

# Stacked Bar 100%



A 100% stacked bar chart titled "Brand" showing the "% of Total Number of Records" (y-axis from 0% to 100%) for each automobile brand along the x-axis. The legend "Class" identifies the categories: Minivcan (blue), Normal (orange), Pickup (red), Sports (teal), SUV (green), and Wagon (yellow).

Brands along x-axis: Acura, Audi, BMW, Buick, Cadillac, Chevrolet, Chrysler, CMC, Dodge, Ford, GMC, Honda, Hummer, Hyundai, Infiniti, Isuzu, Jaguar, Jeep, Kia, Land, Lexus, Lincoln, Mazda, Mercedes-Benz, Mercury, Mini, Mitsubishi, Nissan, Oldsmobile, Pontiac, Porsche, Saab, Saturn, Scion, Subaru, Suzuki, Toyota, Volkswagen, Volvo.

# Multivariate Data: Region-Based Techniques

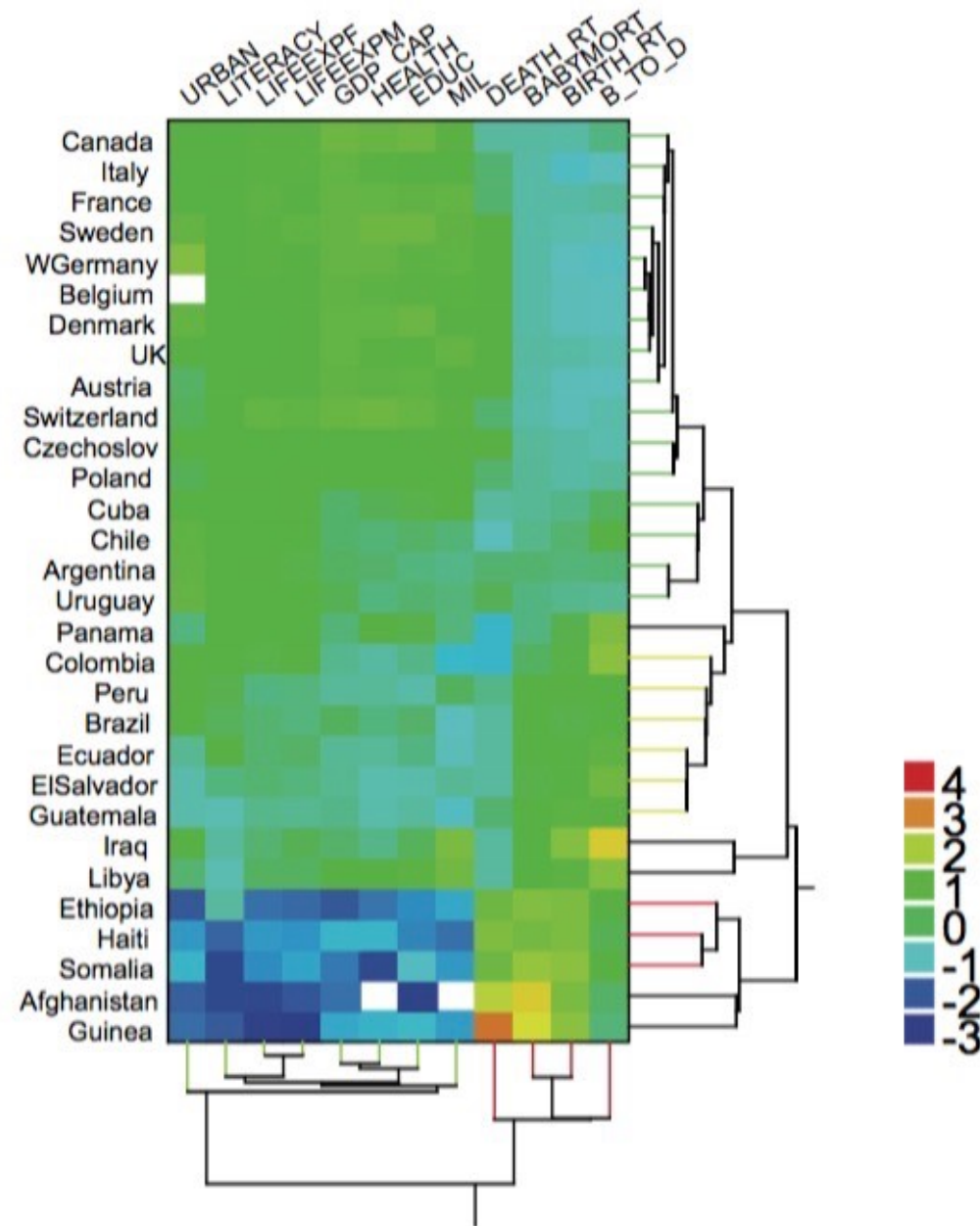- **Tabular Displays**

    ◆ **Heatmaps** are created by displaying the table of record values **using color rather than text**. All **data values are mapped to the same normalized color space**, and each is rendered as a colored square or rectangle.



A heatmap showing social statistics for several countries from a U.N. survey. Rows and columns have been reordered via clustering. (Image courtesy Leland Wilkinson [459].)
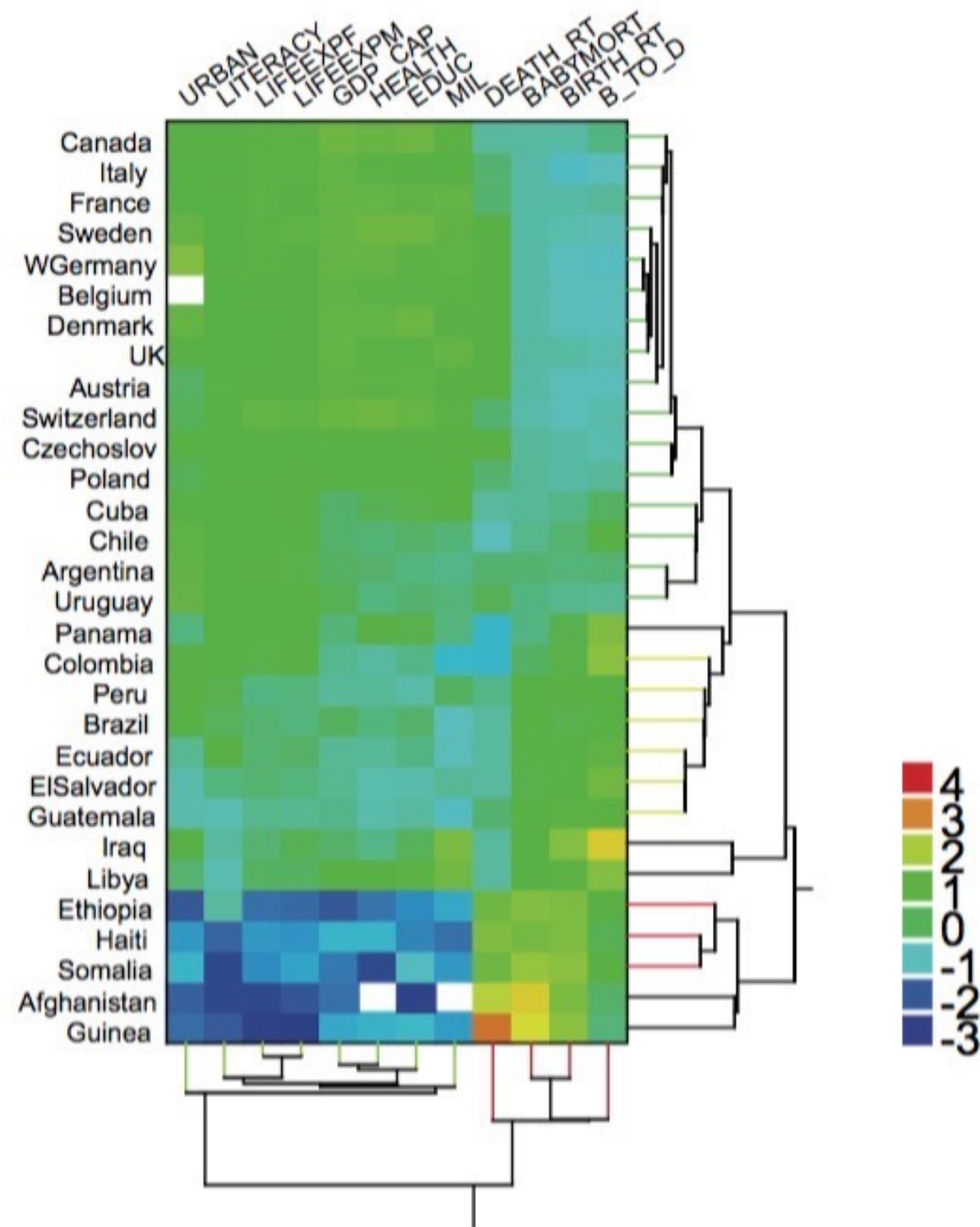
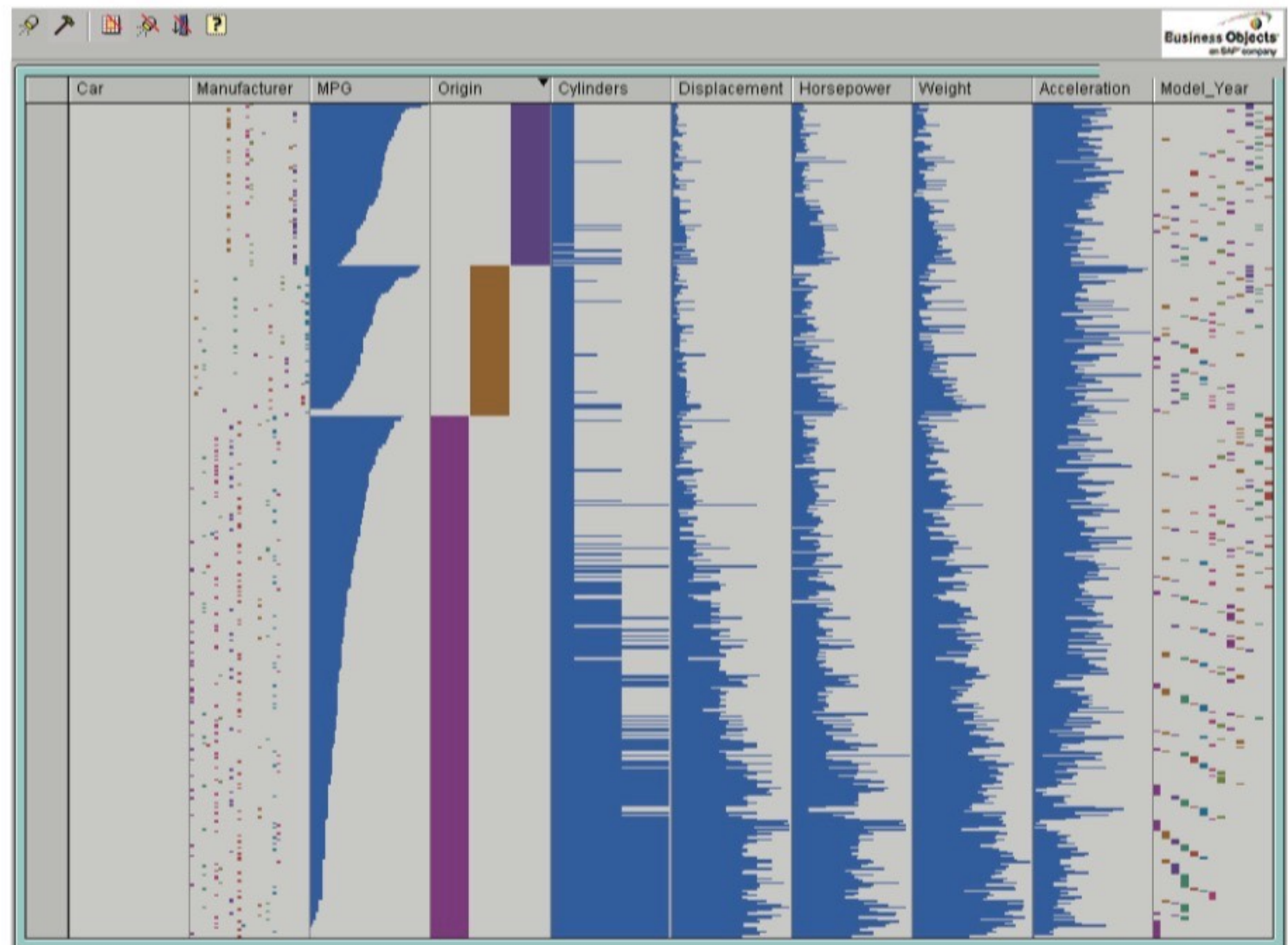# Multivariate Data: Region-Based Techniques



A heatmap showing social statistics for several countries from a U.N. survey. Rows and columns have been reordered via clustering. (Image courtesy Leland Wilkinson [459].)

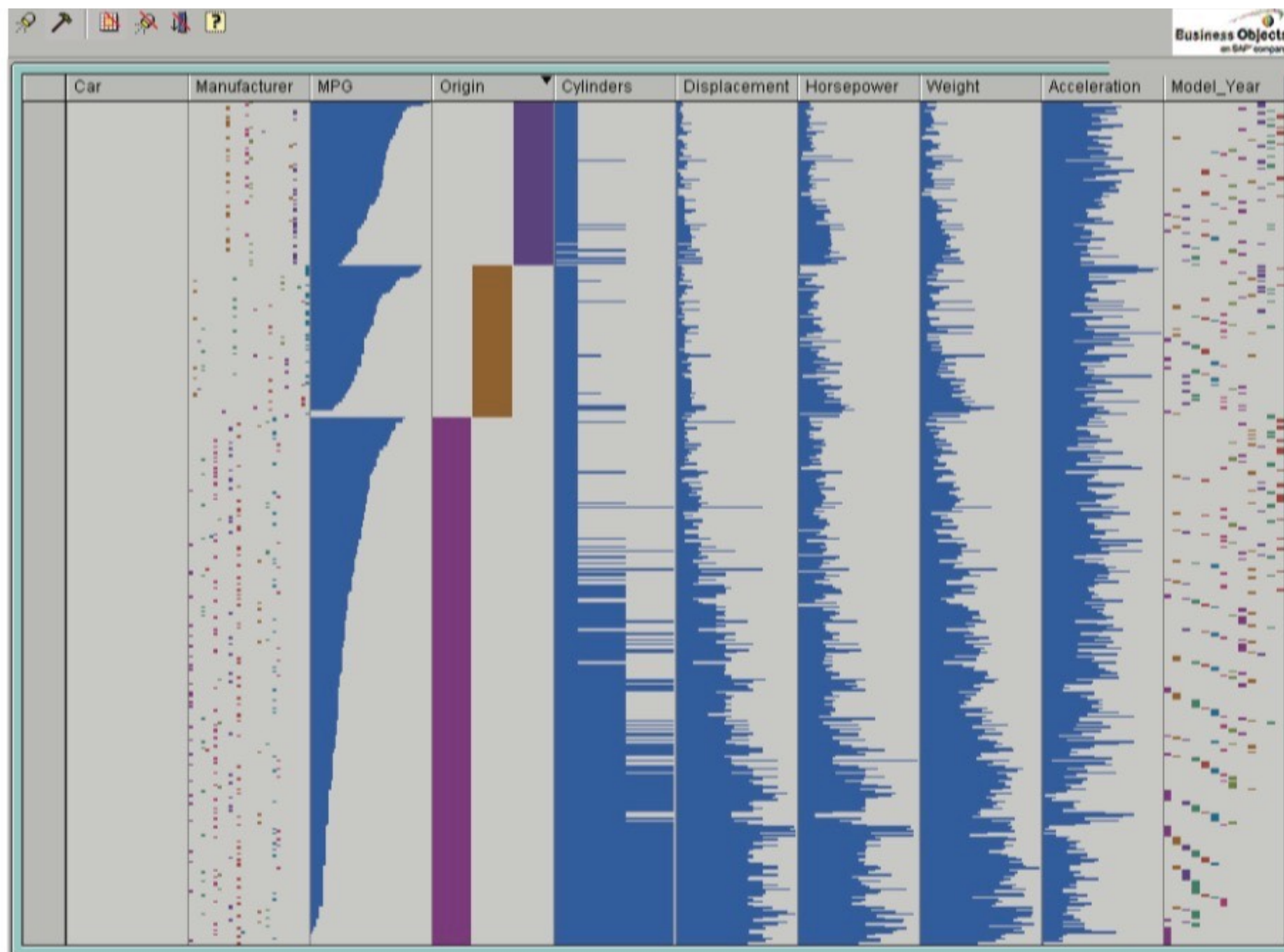# Multivariate Data: Region-Based Techniques

# Multivariate Data: Region-Based Techniques

- **table lens** combines all these ideas and includes a **level-of-detail mechanism** for providing panning and zooming capabilities to display whole table views, while still providing some detail through local table lenses



An example of Inxight Table Lens showing the cars data set sorted first by car origin and then by MPG.
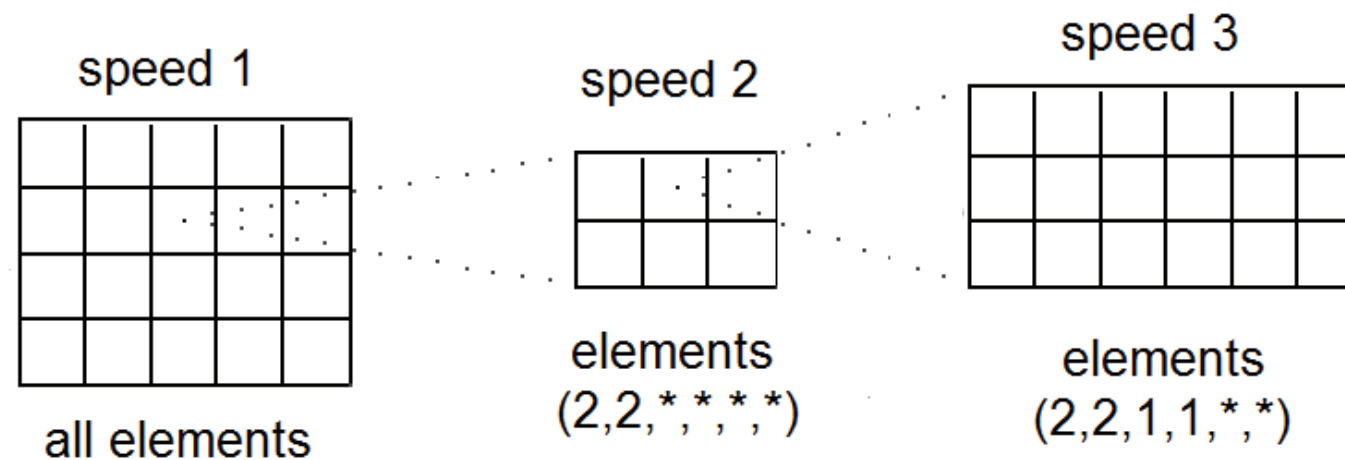
# Multivariate Data: Region-Based Techniques



An example of Inxight Table Lens showing the cars data set sorted first by car origin and then by MPG.

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Multivariate Data: Region-Based Techniques

■ **Dimensional Stacking**

♦ Begin with data of dimension **2N** + 1 (for an even number of dimensions there would be an additional implicit dimension of cardinality one).

♦ Select a **finite cardinality/discretization** for each dimension.

♦ Choose **one** of the dimensions **to be the dependent variable**. The rest will be considered independent

♦ Create ordered pairs of the independent dimensions (**N pairs**) and assign to each pair a unique value (speed) from 1 to N.

♦ The pair corresponding to speed 1 will create a virtual image whose size coincides with the cardinality of the dimensions (the first dimension in the pair is oriented horizontally, the second vertically).

# Multivariate Data: Region-Based Techniques

- **Dimensional Stacking**

  - ♦ Create ordered pairs of the independent dimensions (**N pairs**) and assign to each pair a unique value (speed) from 1 to N.

  - ♦ The pair corresponding to speed 1 will create a virtual image whose size coincides with the cardinality of the dimensions (the first dimension in the pair is oriented horizontally, the second vertically).
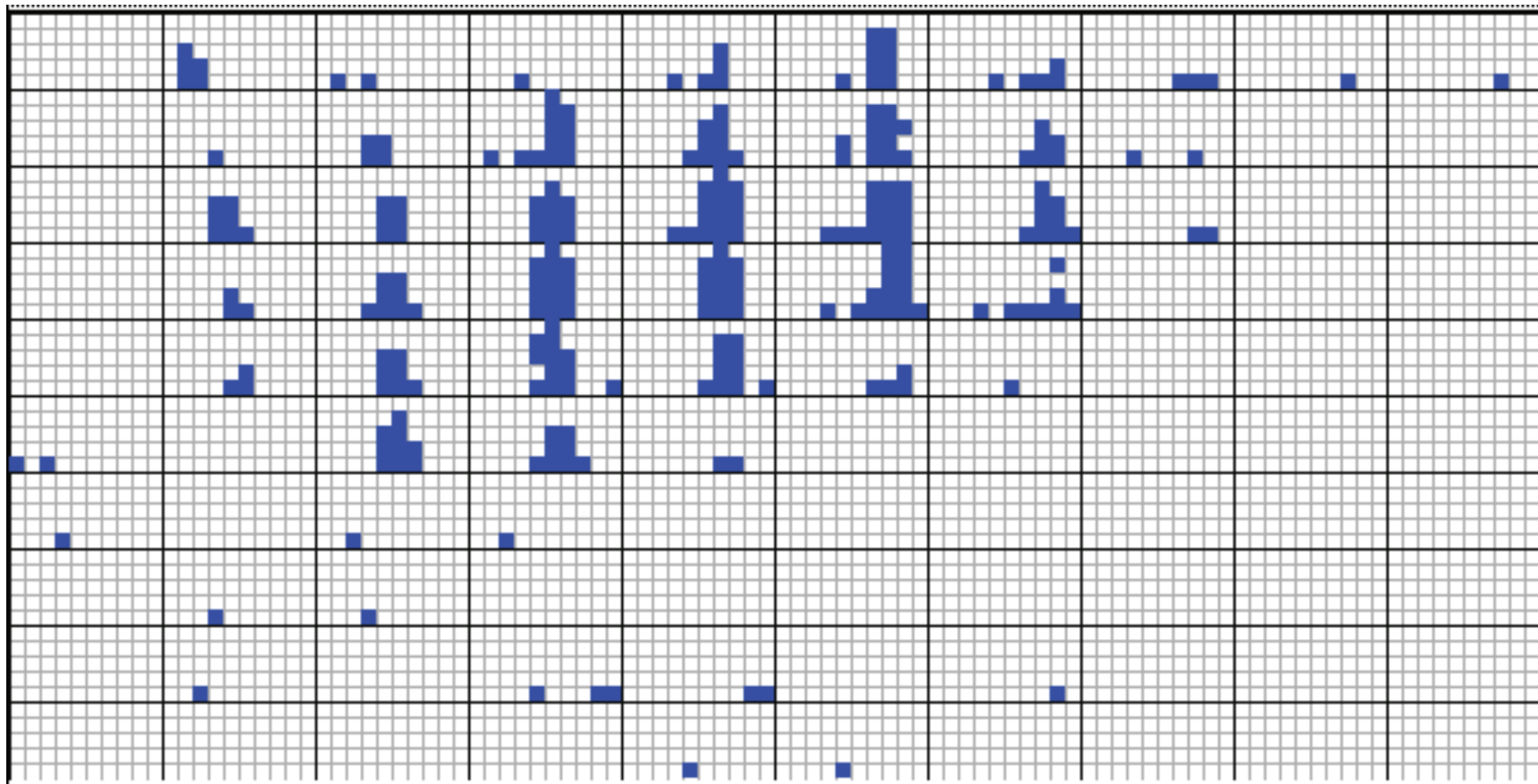


Conceptualization of dimensional stacking; collapsing six dimensions into two dimensions.

**d1,. . . , d6 have cardinalities 4, 5, 2, 3, 3, and 6, respectively**

# Multivariate Data: Region-Based Techniques

- **Dimensional Stacking**



An example of 4D data visualized using dimensional stacking. The data consists of drill-hole data, with three spatial dimensions, and the ore grade as the fourth dimension.

# Combinations of Techniques

# Multivariate Data: Combinations of Techniques

- **Glyphs and Icons**

- **Dense Pixel Displays**

- **Many others**

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
**UNIVERSIDADE NOVA** DE LISBOA

# Multivariate Data: Combinations of Techniques

- **Glyphs and Icons**



Figure 8.20.    Examples of multivariate glyphs (from [445]).

# Further Reading and Summary

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA

# Further Reading

- **Recommend Readings**

  - Interactive Data Visualization: Foundations, Techniques, and Applications, Matthew O. Ward et all, 2015, pages 285-314.

- **Supplemental readings:**

  - Visualization Analysis & Design , Tamara Munzner, Chapter 7

# What you should know

- **Point based techniques**

  - ♦ Classical point base techniques have a limited dimensionality - Scatter based

  - ♦ Dimension reduction or selection for data viz

- **Line based**

  - ♦ Classical line based

  - ♦ Radial Axis Techniques

  - ♦ Parallel coordinates techniques and related stuff

- **Region based**

  - ♦ Reordering the data in graphical tables

- **Combination Techniques**

  - ♦ Dense

  - ■ Glyphs

FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA