

TIES443

Lecture 9

Visualization

Mykola Pechenizkiy

Course webpage: <http://www.cs.jyu.fi/~mpechen/TIES443>

November 17, 2006

Department of Mathematical Information Technology
University of Jyväskylä

Topics for today

- Purpose of visualization in DM/KDD
- Visualization of data
 - 1D, 2D, 3D, 3+D
- Visualization of data mining results
 - Model visualization, model performance visualization
- Visualization of data mining processes
- Application ...
 - Interactive data mining and visualization
 - Time-series visualization
 - Data mining as part of visualization system or vice versa?

Sources

- Books:
 - "[Information Visualization in DM and KD](#)" Fayyad *et al.*
 - "[Information Visualization: Beyond the Horizon](#)": 2nd ed, C. Chen
 - "[The Visual Display of Quantitative Information](#)", 2nd ed, E. Tufte
- People:
 - Edward Tufte, Ben Shneiderman, Daniel A. Keim, Marti Hearst, Tamara Munzner, and other
- Excellent handouts and recommended reading on information visualization by Tamara Munzner
 - <http://www.cs.ubc.ca/~tmm/courses/cpsc533c-04-spr/>
- Some papers:
 - Visual Data Mining Framework for Convenient Identification of Useful Knowledge
 - http://www.cs.uic.edu/~kzhao/Papers/06_ICDM_Zhao_Visual.pdf
- Lots of other sources ...

Visualization & related fields

- User Interface Studies
- Computer Vision
- Cognitive Science
- Computer-Aided Design
- Geometric Modelling
- Computer Graphics
- Signal Processing
- Approximation Theory

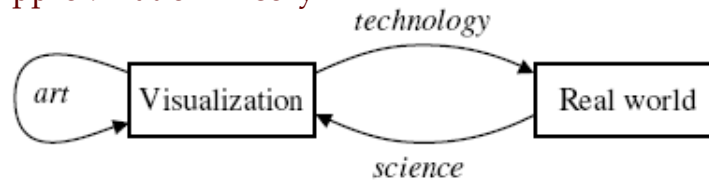


Fig. by Jack van Wijk from "[Views on Visualization](#)"

Purpose of Visualization in DM/KDD

- **Data Visualization**
 - Qualitative overview of large and complex data sets
 - Data summarization and (interactive) exploration
 - Regions of interest identification
 - Support humans in looking for structures, features, patterns, trends, and anomalies in data
 - Perceptual capabilities of the human visual system
- **Model (and its performance) Visualization**
- **Process Visualization**
- **Disadvantages?**
 - requires human eyes
 - can be misleading

"The purpose of computing is insight, not numbers"
Richard Hamming

Visualization Techniques Categorization

- **Based on task at hand**
 - To explore data
 - To confirm a hypothesis
- **Based on structure of underlying data set**
- **Based on the dimension of the display**
 - Stationary
 - Animated
 - Interactive
- **Geometric vs. Symbolic**
 - Results of physical models, simulations, computations
 - Nonnumeric data
 - Networks, relations as graphs; icons, arrays
- **1D-2D-3D, stereoscopic, virtual reality vs. 3+D**
- **Static vs. Dynamic**

Some Basic Visualization Terminology

- Visualization
- Interaction
- Graphical attributes
- Mapping
- Rendering
- Field
- Scalar, Vector, Tensor
- Isosurface
- Glyph
- ...

Perception in Visualization

- What graphical primitives do humans detect quickly?
- What graphical attributes can be accurately measured?
- How many distinct values for a particular attribute can be used without confusion?
- How do we combine primitives to recognize complex phenomena?
- Human perception and information theory

The most of following slides for *data* visualization were adopted from the presentation “Visualization and Data Mining” by G. Piatetsky-Shapiro (available from www.kdnuggets.com)

Data Visualization Outline

- Graphical excellence and lie factor
- Representing data in 1,2, and 3-D
- Representing data in 4+ dimensions
 - Parallel coordinates
 - Scatterplots
 - Stick figures

Napoleon Invasion of Russia, 1812

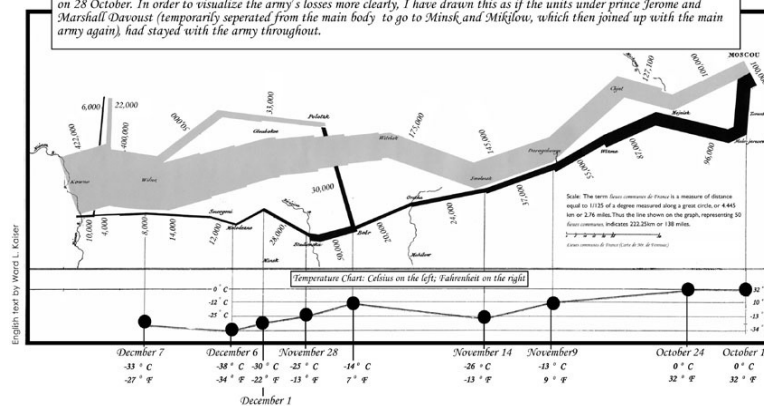


Napoleon



Map representing the losses over time of French army troops during the Russian campaign, 1812-1813. Constructed by Charles Joseph Minard, Inspector General of Public Works retired. Paris, 20 November 1869

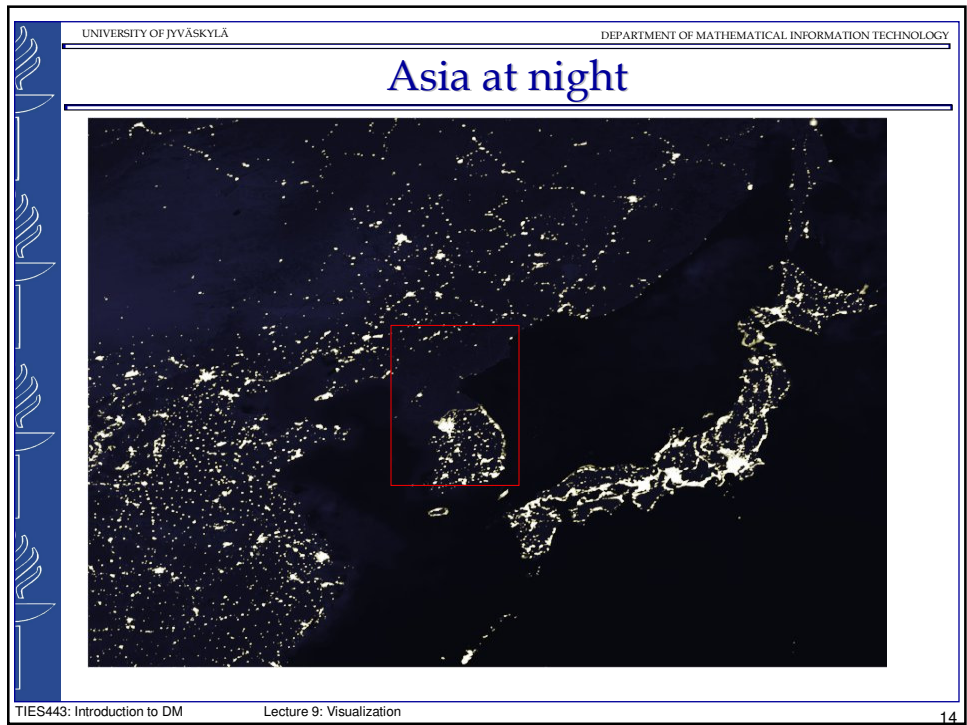
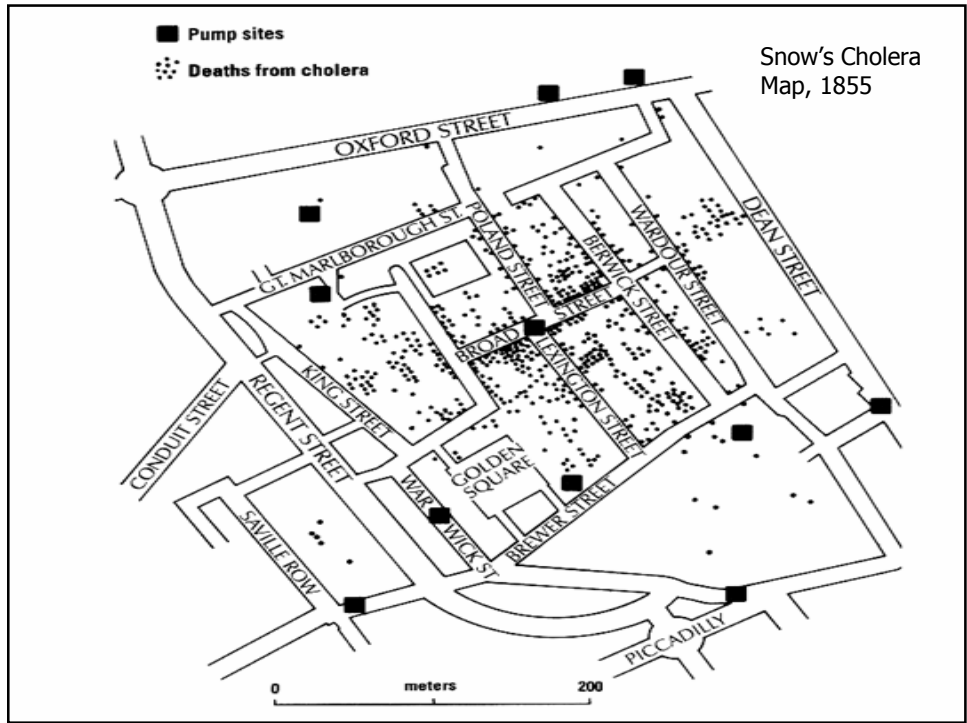
The number of men present at any given time is represented by the width of the grey line: one mm. indicates ten thousand men. Figures are also written besides the lines. Grey designates men moving into Russia; black, for those leaving. Sources for the data are the works of messrs. Thiers, Segur, Fezensac, Chambray and the unpublished diary of Jacob, who became an Army Pharmacist on 28 October. In order to visualize the army's losses more clearly, I have drawn this as if the units under prince Jerome and Marshall Davoust (temporarily separated from the main body to go to Minsk and Mikilow, which then joined up with the main army again), had stayed with the army throughout.



Editor's note: dates & temperatures are only referenced for the retreat from Moscow © 2001, ODT Inc. All rights reserved.

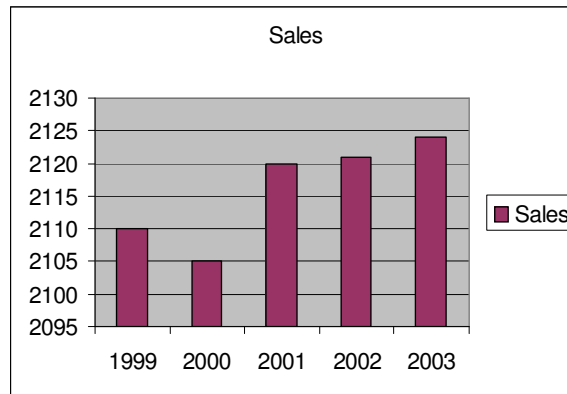
Figure 58. Minard's map of Napoleon's Russian campaign. This graphic has been translated from French to English and modified to most effectively display the temperature data.

© www.odt.org, from <http://www.odt.org/Pictures/minard.jpg>, used by permission



Bad Visualization: misleading Y-axis

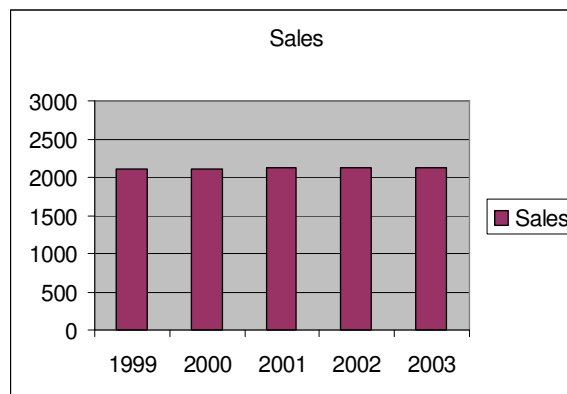
Year	Sales
1999	2110
2000	2105
2001	2120
2002	2121
2003	2124



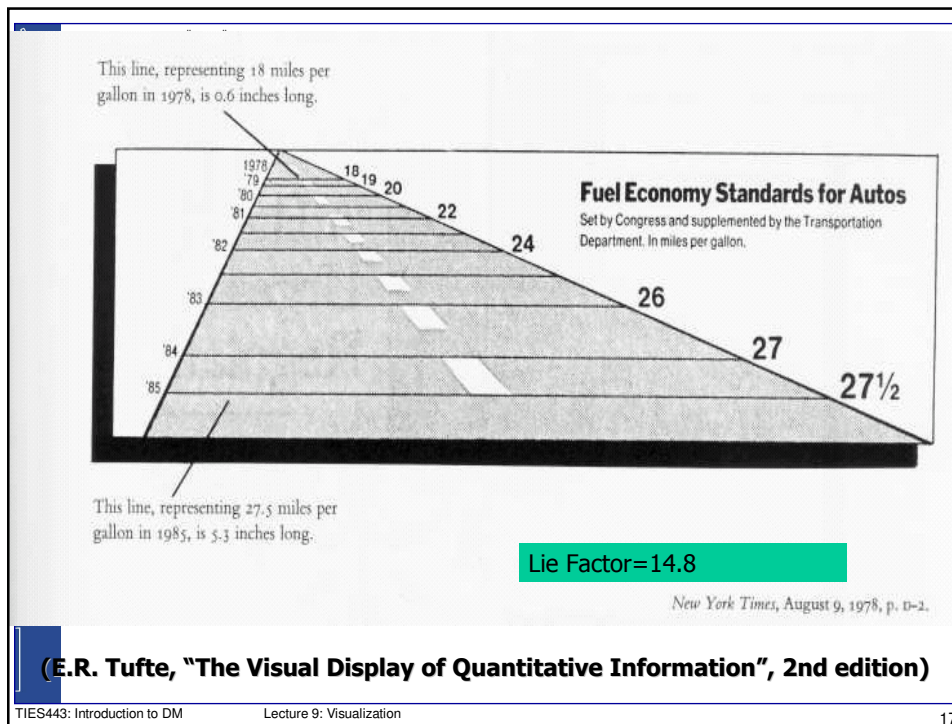
Y-Axis scale gives **WRONG** impression of big change

Better Visualization

Year	Sales
1999	2110
2000	2105
2001	2120
2002	2121
2003	2124



Axis from 0 to 2000 scale gives correct impression of small change



UNIVERSITY OF JYVÄSKYLÄ DEPARTMENT OF MATHEMATICAL INFORMATION TECHNOLOGY

Lie Factor

$$\text{Lie Factor} = \frac{\text{size of effect shown in graphic}}{\text{size of effect in data}} =$$

$$= \frac{(5.3 - 0.6)}{18} = \frac{0.6}{(27.5 - 18.0)} = \frac{7.833}{0.528} = 14.8$$

Tufte requirement: $0.95 < \text{Lie Factor} < 1.05$

(E.R. Tufte, "The Visual Display of Quantitative Information", 2nd edition)

TIES443: Introduction to DM Lecture 9: Visualization 18

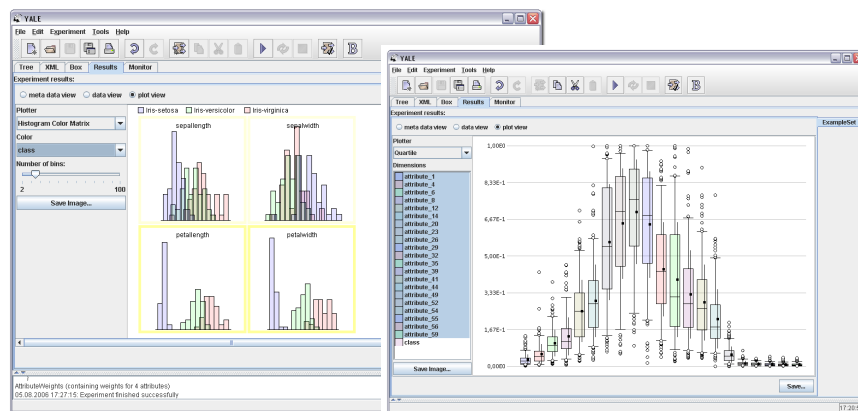
Tufte's Principles of Graphical Excellence

- Give the viewer
 - the greatest number of ideas
 - in the shortest time
 - with the least ink in the smallest space.
- Tell the truth about the data!

(E.R. Tufte, "The Visual Display of Quantitative Information", 2nd edition)

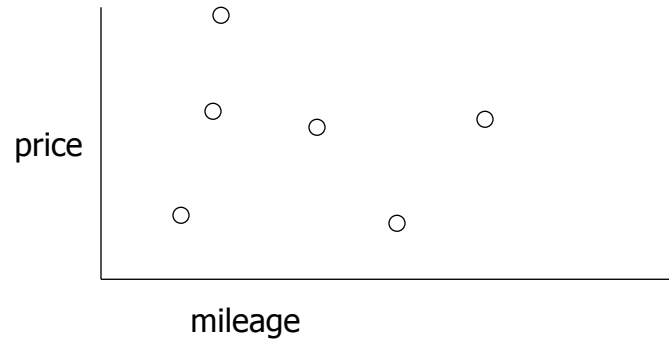
1-D (Univariate) Data

- Representations

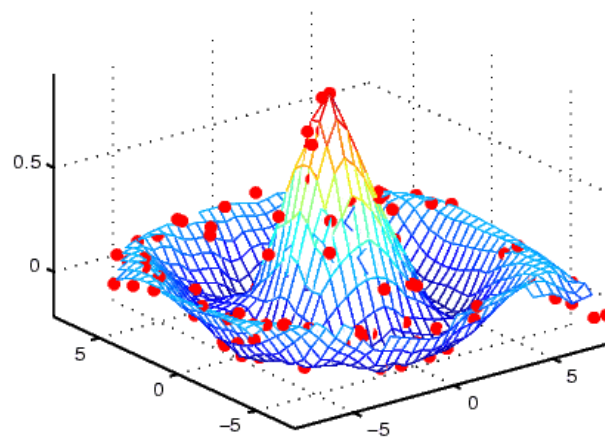


2-D (Bivariate) Data

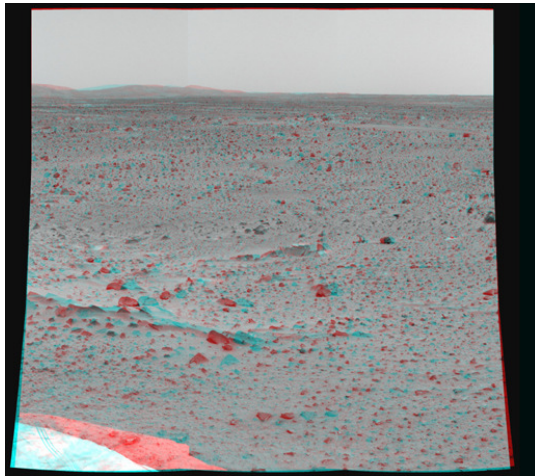
- Scatter plot, ...



3-D Data (projection)



3-D image (requires 3-D blue and red glasses)



Taken by Mars Rover Spirit, Jan 2004

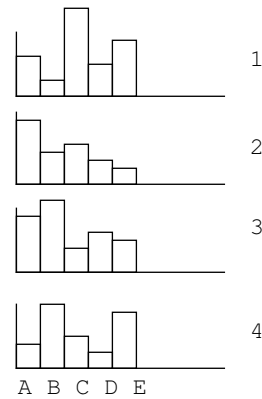
Visualizing in 4+ Dimensions

- Scatterplots
- Parallel coordinates
- Chernoff faces
- Stick figures
- Star glyphs
- ...

Multiple Views

Give each variable its own display

	A	B	C	D	E
1	4	1	8	3	5
2	6	3	4	2	1
3	5	7	2	4	3
4	2	6	3	1	5



Problem: does not show correlations

Scatterplot Matrix

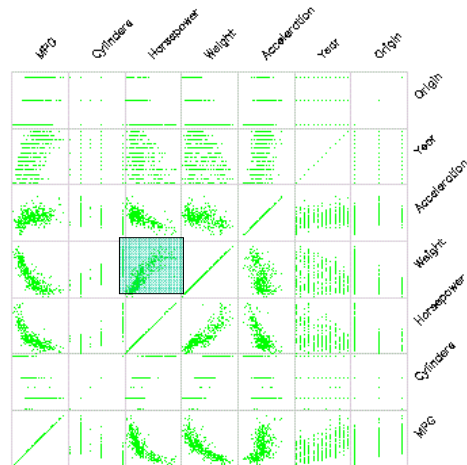
Represent each possible pair of variables in their own 2-D scatterplot (car data)

Q: Useful for what?

A: linear correlations (e.g. horsepower & weight)

Q: Misses what?

A: multivariate effects



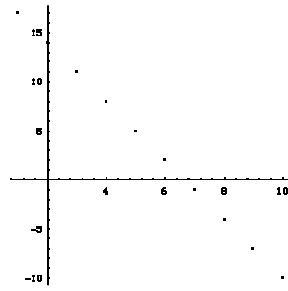
For 13 dimensions (columns) :

Number of histograms = 13

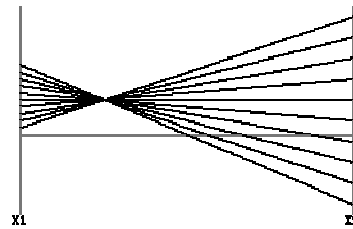
Number of scatterplots = $C(13,2) = 78$

Parallel Coordinates

- Encode variables along a horizontal row
- Vertical line specifies values



Dataset in a Cartesian coordinates



Same dataset in parallel coordinates

Invented by
Alfred Inselberg
while at IBM, 1985

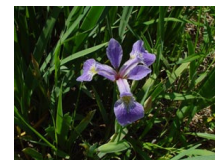


Example: Visualizing Iris Data



Iris setosa

sepal length	sepal width	petal length	petal width
5.1	3.5	1.4	0.2
4.9	3	1.4	0.2
...
5.9	3	5.1	1.8

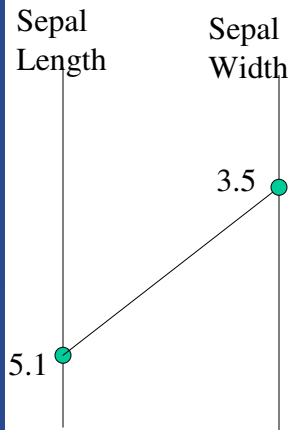


Iris versicolor



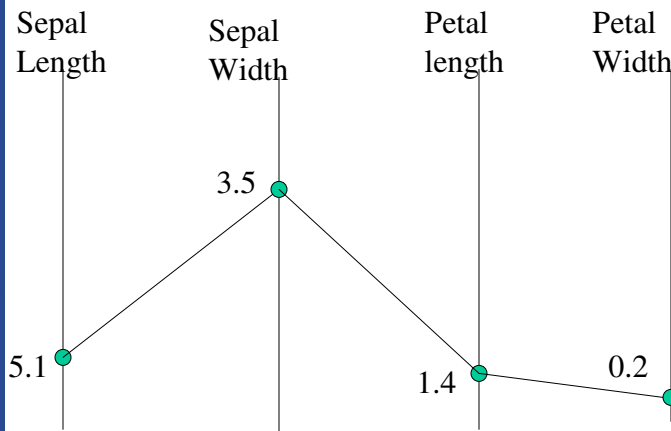
Iris virginica

Parallel Coordinates: 2 D

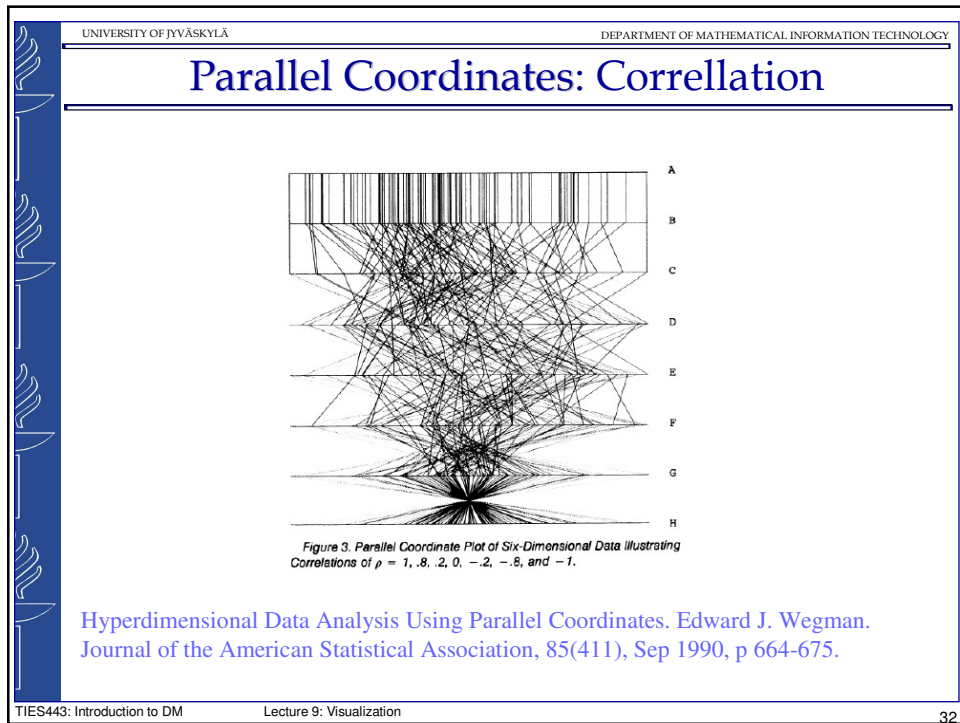
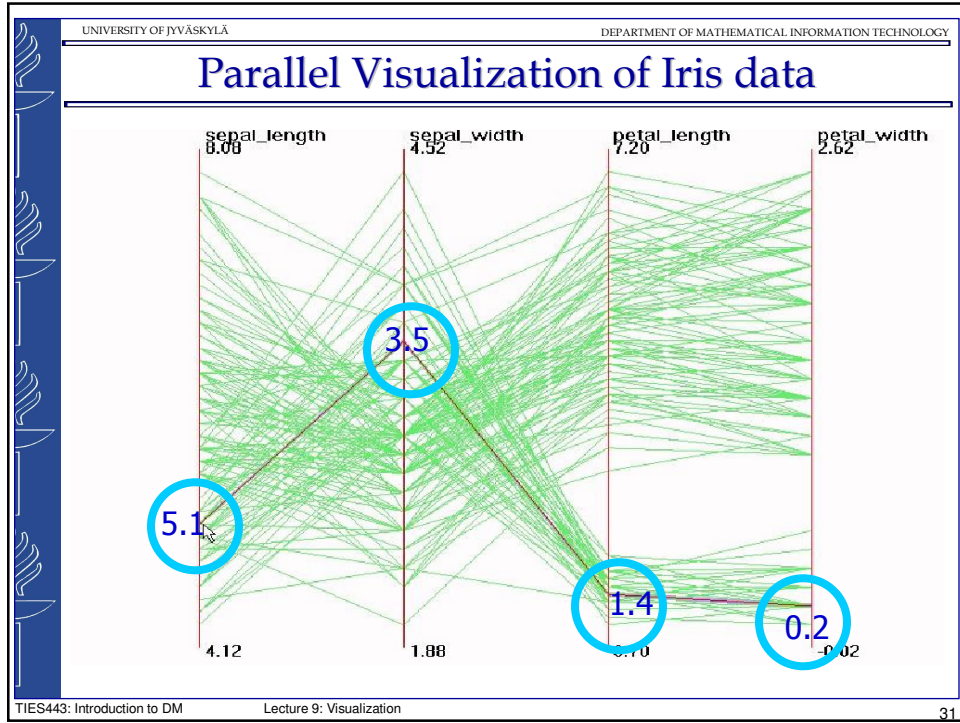


sepal length	sepal width	petal length	petal width
5.1	3.5	1.4	0.2

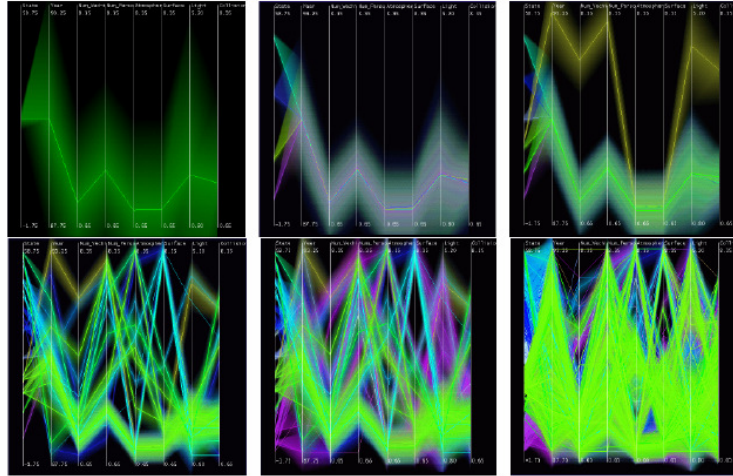
Parallel Coordinates: 4 D



sepal length	sepal width	petal length	petal width
5.1	3.5	1.4	0.2



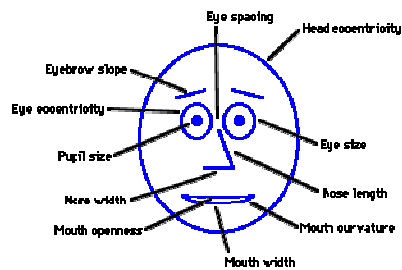
Hierarchical Parallel Coordinates



Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Fua, Ward, and Rundensteiner, IEEE Visualization 99

Chernoff Faces

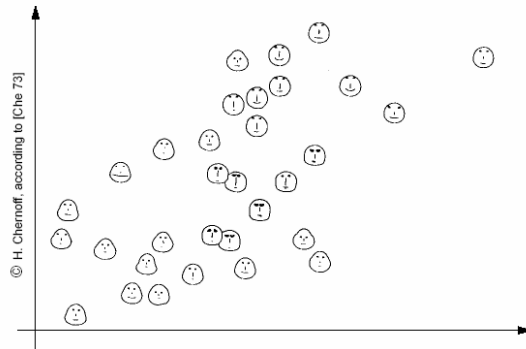
Encode different variables' values in characteristics of human face



Cute applets: <http://www.cs.uchicago.edu/~wiseman/chernoff/>
<http://hesketh.com/schampeo/projects/Faces/chernoff.html>

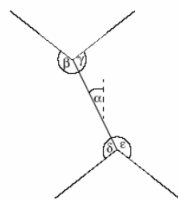
Chernoff faces, example

Chernoff-Faces [Che 73, Tuf 83]

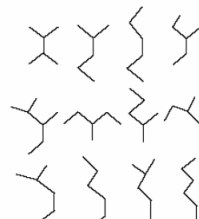


Stick Figures

- Two variables are mapped to X , Y axes
- Other variables are mapped to limb lengths and angles
- Texture patterns can show data characteristics

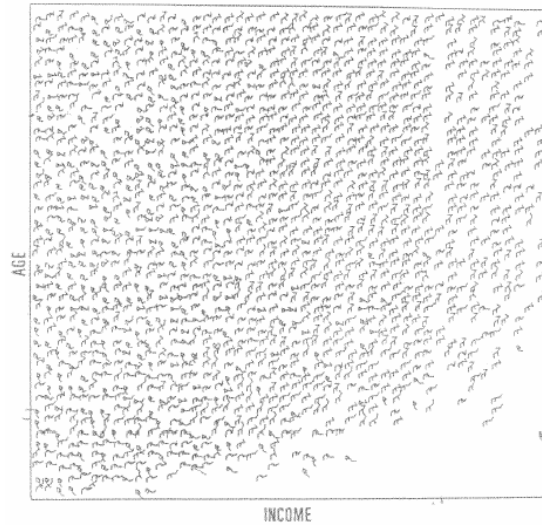


Stick Figure Icon



A Family of Stick Figures

Stick figures, example



used by permission of G. Grinstein, University of Massachusetts at Lowell

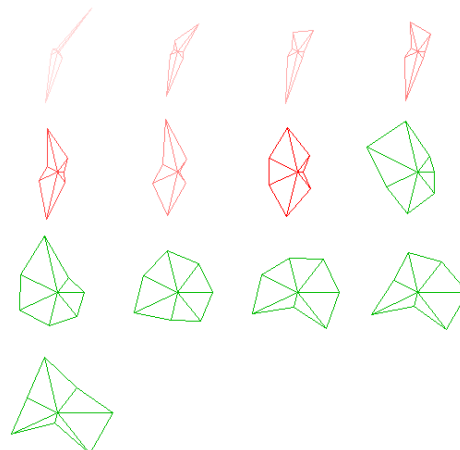
census data
showing
age, income, sex,
education, etc.

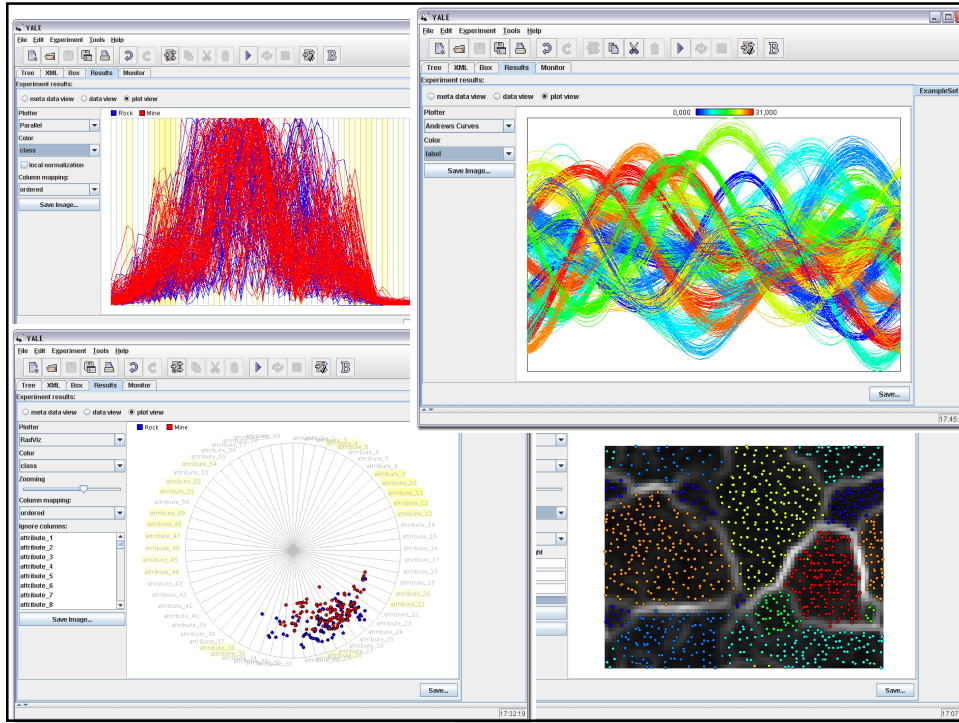
Closed figures
correspond to
women and we can
see more of them
on the left.

Note also a young
woman with high
income

Star Glyphs from Xmdv

- Two variables are mapped to X, Y axes
- Other variables are mapped to limb lengths and angles
- Texture patterns can show data characteristics





UNIVERSITY OF JYVÄSKYLÄ DEPARTMENT OF MATHEMATICAL INFORMATION TECHNOLOGY

Model Visualization

Weka Classifier Tree Visualizer: 11:49:05 - trees_j48_j48 (iris)

```

graph TD
    Root((petalwidth)) -- <= 0.6 --> Node1[iris-setosa (50.0)]
    Root -- > 0.6 --> Node2((petalwidth))
    Node2 -- <= 1.7 --> Node3((petalwidth))
    Node2 -- > 1.7 --> Node4((petalwidth))
    Node3 -- <= 4.9 --> Node5[iris-versicolour (48.0)]
    Node3 -- > 4.9 --> Node6((petalwidth))
    Node6 -- <= 1.5 --> Node7[iris-virginica (3.0)]
    
```

Neural Network

The diagram shows a neural network with three layers of nodes. The input layer has three nodes, the hidden layer has three nodes, and the output layer has three nodes. Connections are shown between nodes in adjacent layers.

3D Scatter Plot

A 3D scatter plot showing data points in a 3D space. The axes are labeled x, y, and z. The points are colored and clustered.

3D Bar Chart

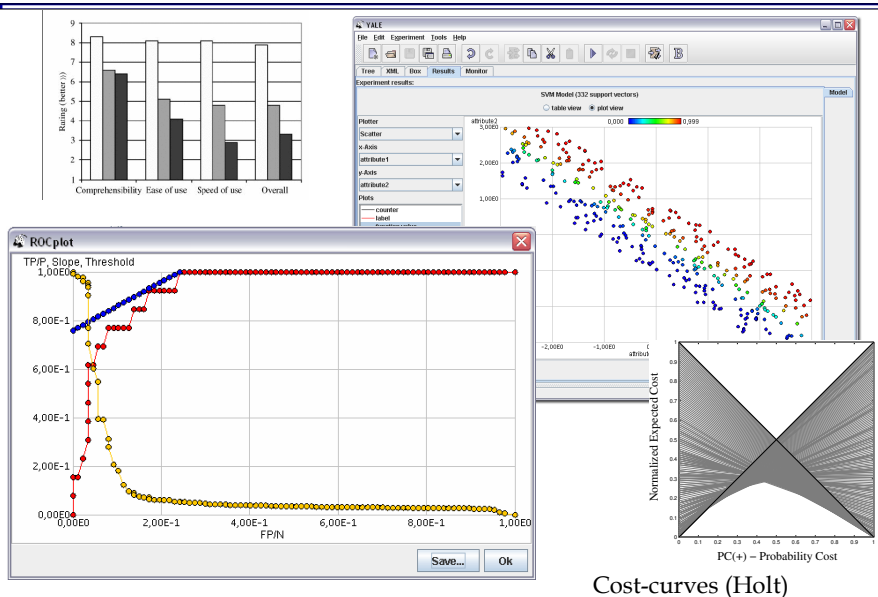
A 3D bar chart showing data points in a 3D space. The axes are labeled x, y, and z. The bars are colored and clustered.

TIES443: Introduction to DM Lecture 9: Visualization 40

DM Model Visualization: Why

- **Understanding the model**
 - allow a user to discuss and explain the logic behind the model
 - also with colleagues, customers, and other users
 - more than just comprehension; it also involves business context
 - financial indicators like profit and cost
 - visualization of the DM output in a meaningful way
 - allowing the user to interact with the visualization
 - simple "what if" questions can be answered
 - 3 whales: representation, interaction, and integration
- **Trusting the model**
 - good quantitative measures of "trust"
 - the overall goal of "visualizing trust" - understanding the limitations of the model
- **Comparing Different Models using Visualization**
 - as Input-Output Mappings, as Algorithms, as Processes

Model Visualization: Performance Evaluation



Visualization of DM/KDD Process

The screenshot displays the Weka KnowledgeFlow Environment. The main window shows a workflow diagram with nodes such as 'Iris', 'CrossValidationFoldMaker', 'AttributeSelection', 'SMO', 'DataVisualizer', 'AttributeSummarizer', 'ScatterPlotMatrix', 'ClassifierPerformanceEvaluator', and 'TextViewer'. A 'New Operator' dialog is open, showing search constraints and a list of operators. The status bar at the bottom indicates '05.08.2006 19:06:33: Experiment finished successfully'.

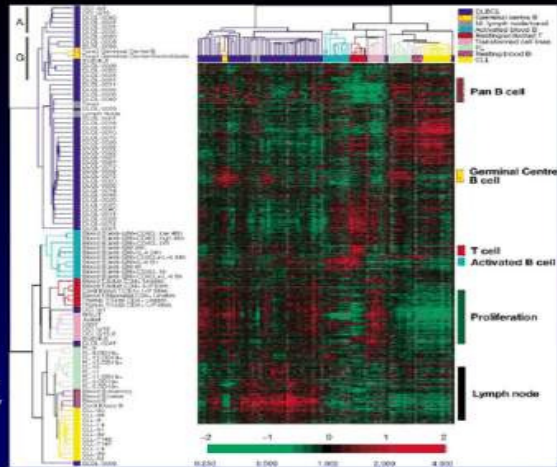
Applications ...

- Bioinformatics
- Visual genomics
- Social networks
- Process mining
- Time series visualization

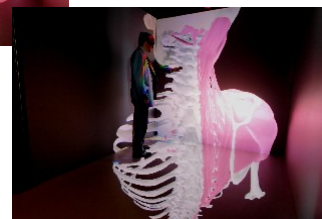
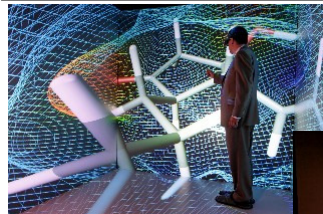
Visualization of Clusters of Microarray Data

Clusters

Taken from
Nature February, 2000
Paper by Alizadeh, A. et al
*Distinct types of diffuse large
B-cell lymphoma identified by
Gene expression profiling.*



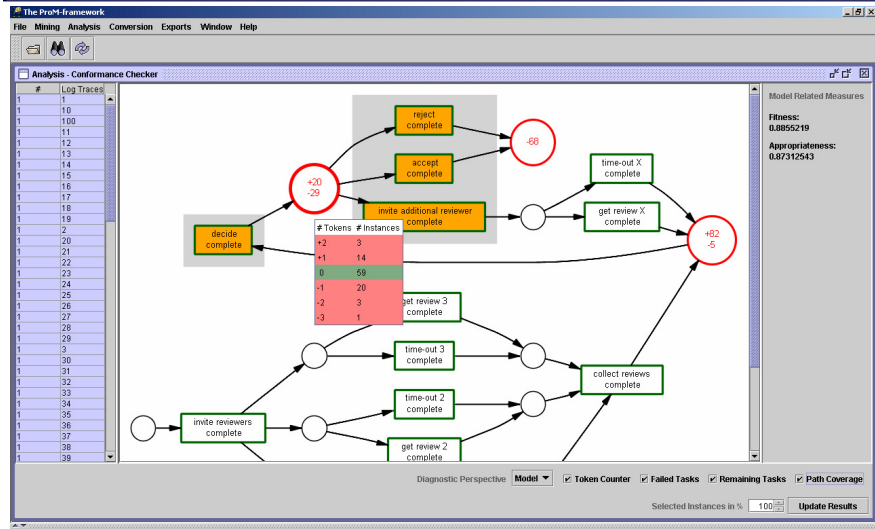
Visual Genomics



Pictures by Christoph W. Sensen

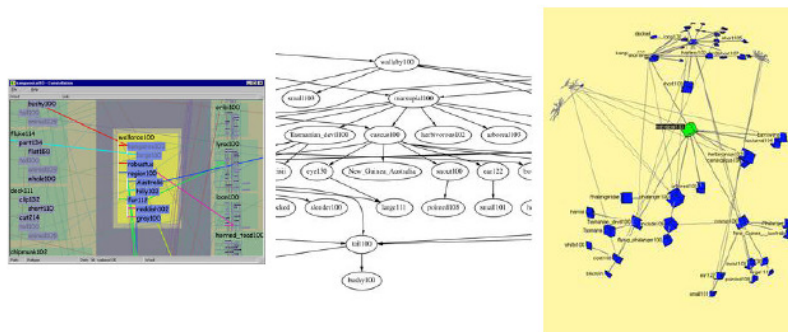
www.visualgenomics.ca/index.php?option=com_content&task=section&id=15&Itemid=143

The Result of Process Mining

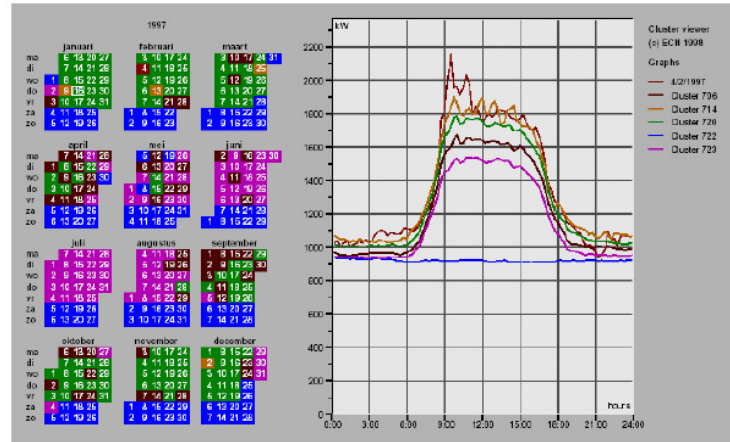


From presentation by Wil van der Aalst, TU/e www.processmining.org

Design Studies

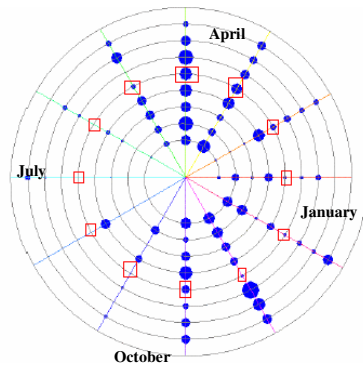


Power Consumption



van Wijk and van Selow, Cluster and Calendar based Visualization of Time Series Data, InfoVis99, <http://www.win.tue.nl/~vanwijk/clv.pdf>

Time Series Spirals

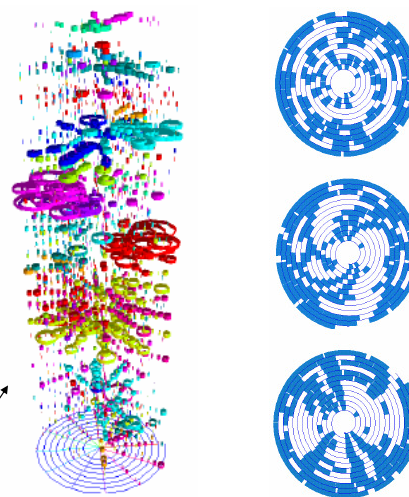


Chimpanzees Monthly Food Intake 1980-1988

The spokes are months, and spiral guide lines are years

112 types of food

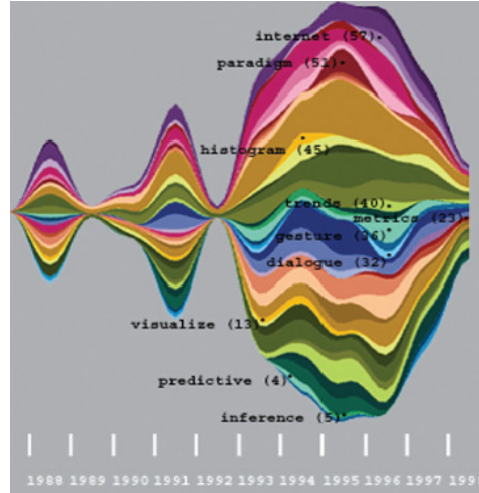
(c) Eamonn Keogh, eamonn@cs.ucr.edu



Simple and intuitive, but works only with periodic data, and if the period is known

ThemeRiver

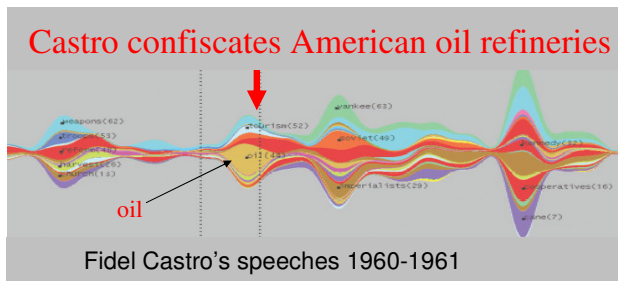
- Current width = strength of theme
- River width = global strength
- Color mapping (similar themes/ color family) same
- Time axis
- External events can be linked



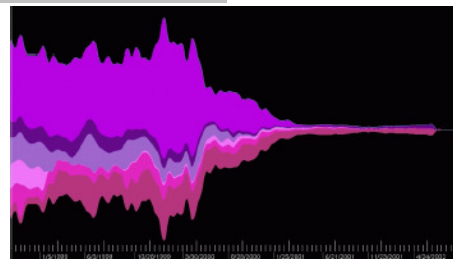
Havre, Hetzler, Whitney & Nowell
InfoVis 2000

A company's patent activity
1988 to 1998

ThemeRiver



- Simple and intuitive
- Many extensions possible
- Scalability is still an issue

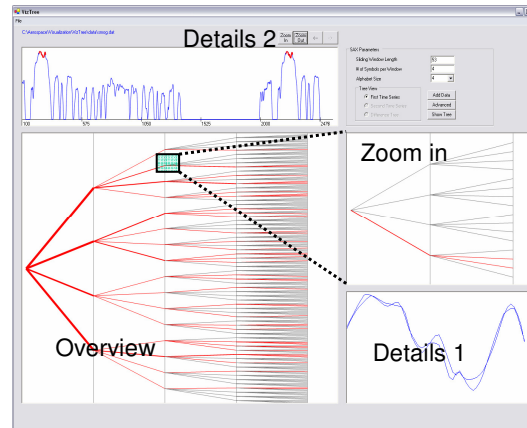


(c) Eamonn Keogh, eamonn@cs.ucr.edu

dot.com stocks 1999-2002

VizTree

The “trick” on the previous slide only works for discrete data, but time series are *real* valued.
we will discuss how we can SAX up a time series to make it discrete during the next tutorial on time series mining



“Overview, zoom & filter, details on demand”
– main principles of visualization Ben Shneiderman

(c) Eamonn Keogh, eamonn@cs.ucr.edu

Visualization Summary

- Purpose of visualization in DM/KDD
- Visualization of data - Many methods
 - 1D, 2D, 3D
 - Visualization is possible in more than 3-D
- Visualization of data mining results
 - Model visualization, model performance visualization
- Visualization of data mining processes
- Application ...
 - Interactive and integrative data mining and visualization
 - Time-series, symbolic, networks visualization

One picture may worth 1000 words!

What else did you get from this lecture?

UNIVERSITY OF JYVÄSKYLÄ DEPARTMENT OF MATHEMATICAL INFORMATION TECHNOLOGY

Additional Slides

TIES443: Introduction to DM Lecture 9: Visualization 59

UNIVERSITY OF JYVÄSKYLÄ DEPARTMENT OF MATHEMATICAL INFORMATION TECHNOLOGY

Visualization Software

Free and Open-source

- Ggobi
- Xmdv
- Some techniques are available in WEKA and YALE as you have seen already
- Many more - see www.kdnuggets.com/software/visualization.html

Matlab has also many nice integrated tools for viz.

TIES443: Introduction to DM Lecture 9: Visualization 60