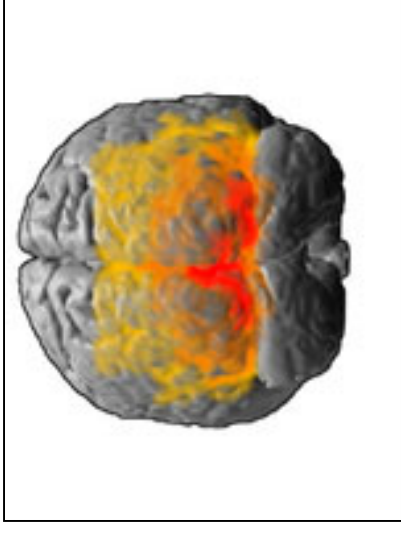# Application of data visualisation in particle physics

Steve Watts
BITlab
School of Engineering and Design
Brunel University, West London, UK

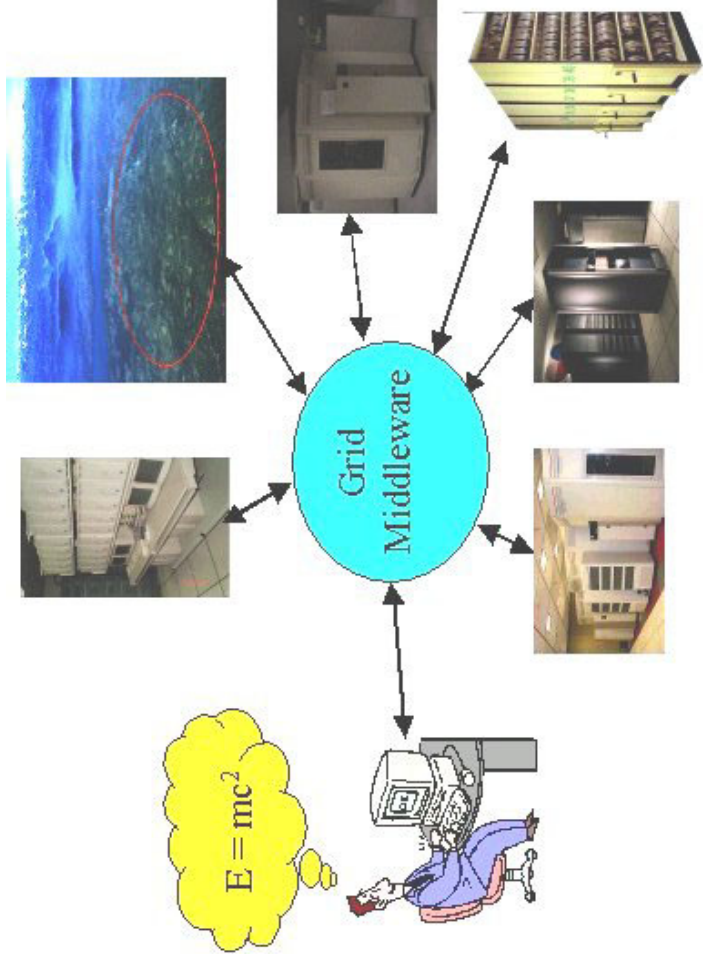University of Manchester – from Jan '07

NOTE: "The visual cortex is the most massive system in the human brain and is responsible for higher-level processing of the visual image. It lies at the rear of the brain (highlighted in the image), above the cerebellum"

There is more to data visualisation than histograms, scatterplots and x/y plots.
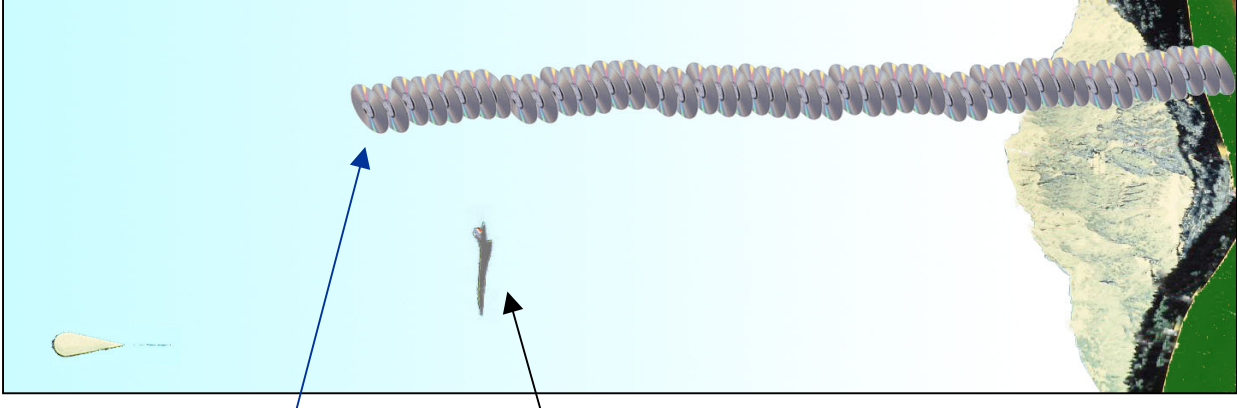
*Edinburgh 8 November 2006*

- WHAT SPARKED THE INTEREST ???
- INTRODUCTION TO DATA MINING
- DATA VISUALISATION
  (history, look at a PP dataset, parallel coordinates, brushing, pruning….)
- CLASSIFICATION ALGORITHMS
  Analysis of wine dataset (Decision Trees, SVM, kNN )
  Links to visualisation – e.g. GRAND Tour, Radviz, polyviz
- ANALYSIS of PP dataset with different techniques
- OTHER USEFUL VISUALISATION TECHNIQUES
  ( survey plots, heat maps, mosaic dsiplay)
- DIY (Do it yourself ) – some advice.
- CONCLUSIONS

CD stack with
1 year LHC data
(~ 20 km)

(Ex-)Concorde
(15 km)

Grid Middleware

$E = mc^2$

GRID computing gives massive computing power.
No excuse for not being smarter with
data analysis techniques
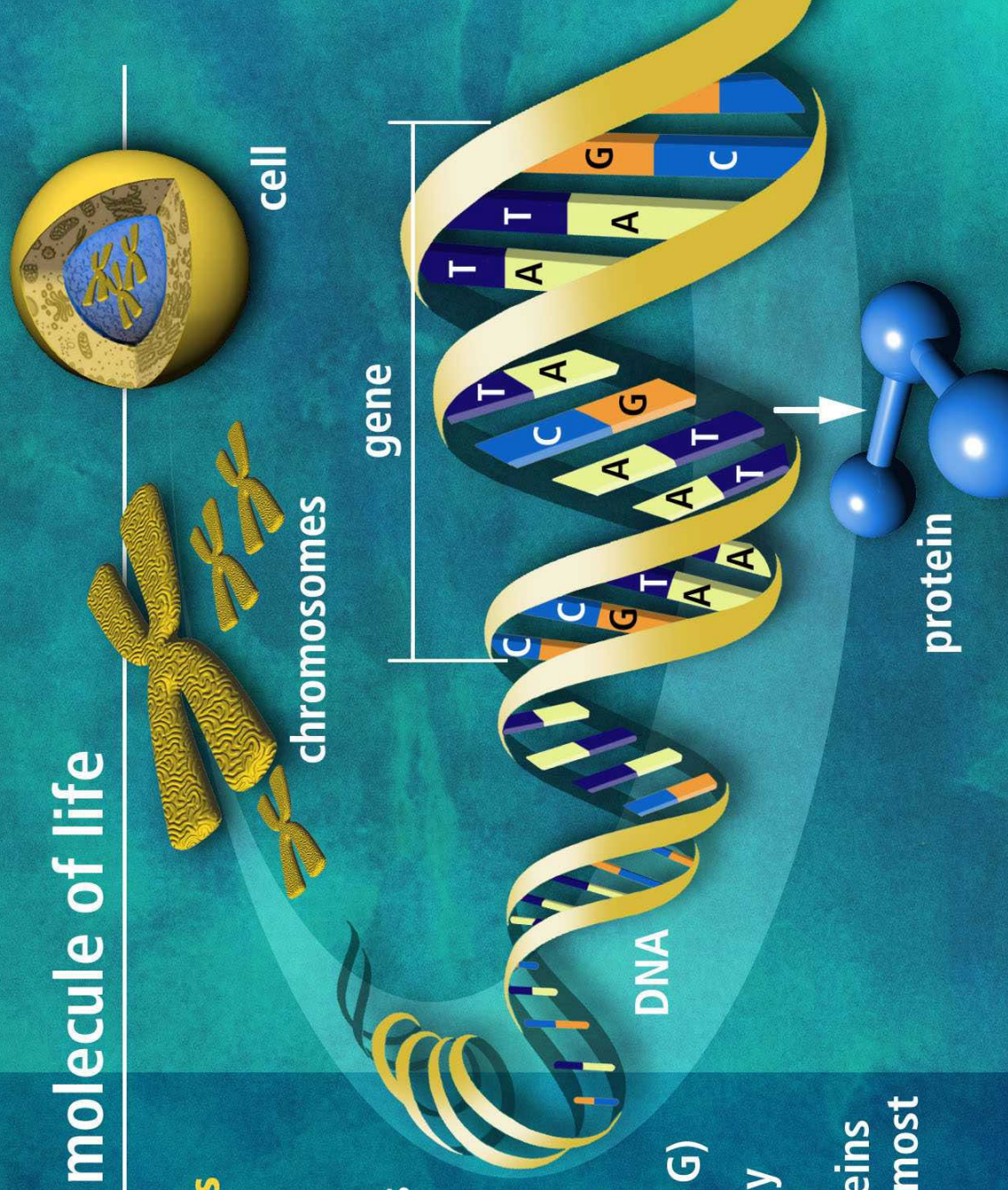
Thanks to Steve Lloyd for the nice pictures…..

MANY AREAS OF MODERN SCIENCE, ENGINEERING, HUMANITIES AND ARTS HAVE A DATA OCEAN TO SWIM IN !!



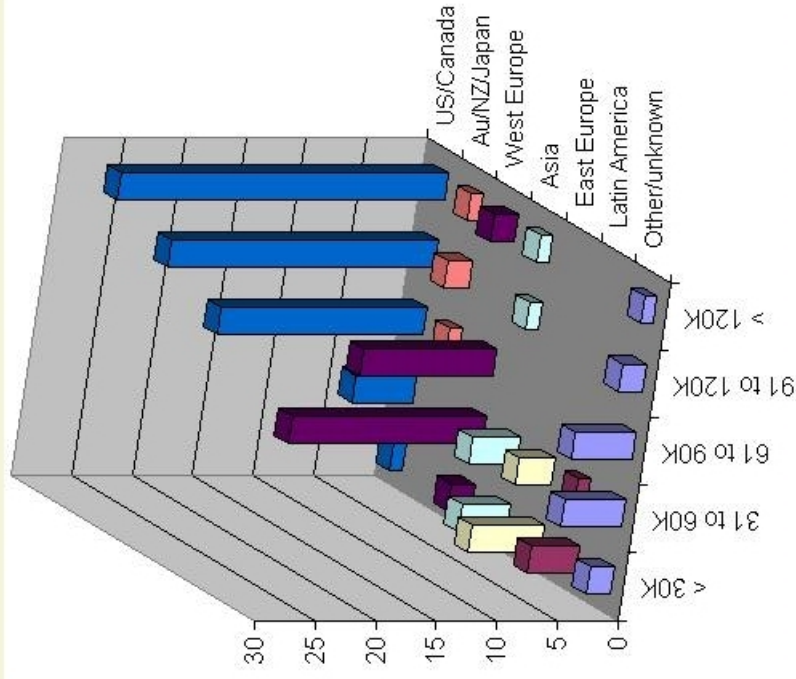# DNA the molecule of life

**Trillions of cells**

Each cell:

- 46 human chromosomes
- 2 meters of DNA
- 3 billion DNA subunits (the bases: A, T, C, G)
- Approximately 30,000 genes code for proteins that perform most life functions

cell

chromosomes

gene

DNA

protein

**Data mining, also known as knowledge-discovery in databases (KDD), is the practice of automatically searching large stores of data for patterns. To do this, data mining uses computational techniques from statistics and pattern recognition.**

en.wikipedia.org/wiki/Datamining

http://www.kdnuggets.com/polls/2006/data_miner_income_by_region.htm

The following graph shows the breakdown of income (salary) by region, with number of respondents on the vertical axis (excluding students). (Note: the poll asks for income to include data miners who work for a company as well as self-employed).
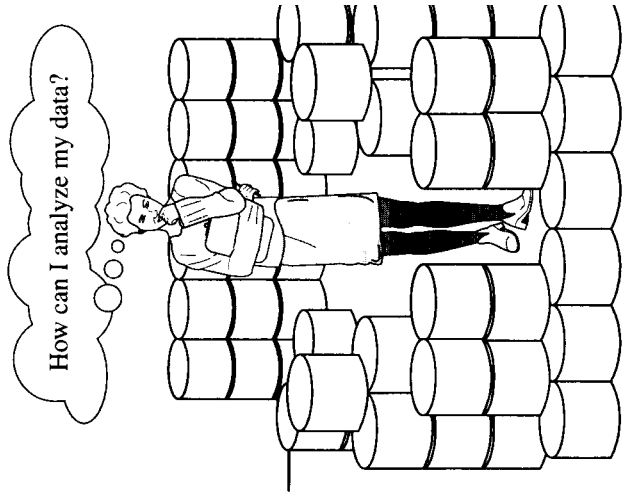


The message is clear!

**Figure 1.2** We are data rich, but information poor.



**Figure 1.4** Data mining as a step in the process of knowledge discovery.

Data Mining:
Concepts and Techniques
2nd ed. Jiawei Han
Micheline Kamber

**WEKA**
The University
of Waikato

Data Mining: Practical Machine Learning Tools and Techniques (Second Edition)

Ian H. Witten, Eibe Frank

Morgan
Kaufmann
June 2005
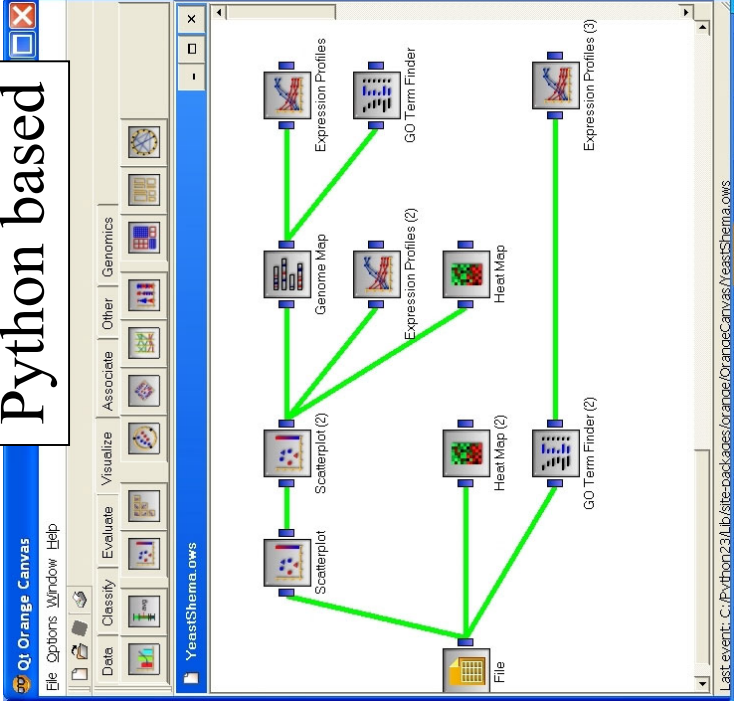525 pages
Paper
ISBN
0-12-088407-0

Eibe Frank and Ian Witten

The Elements of Statistical Learning
Trevor Hastie
Robert Tibshirani
Jerome Friedman

Data Mining, Inference, and Prediction

Springer Series in Statistics

Springer

Statistical
Learning
KEY TEXT

WEKA Machine Working Workbench
Waikato Environment for Knowledge Analysis

"We are drowning in information and starving for knowledge"

Rutherford D Roger

"Information is not knowledge"
Albert Einstein

Java based

Python based



Knowledge Flow
Similar to ORANGE
intrface

http://www.cs.waikato.ac.nz/~ml/weka/

http://www.ailab.si/orange

KEY TERMS
Preprocess Data
Classify
Cluster
Associate
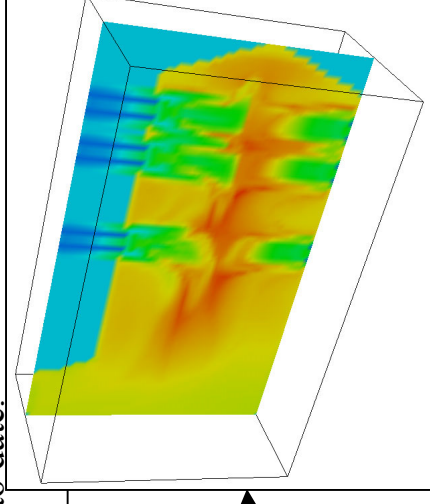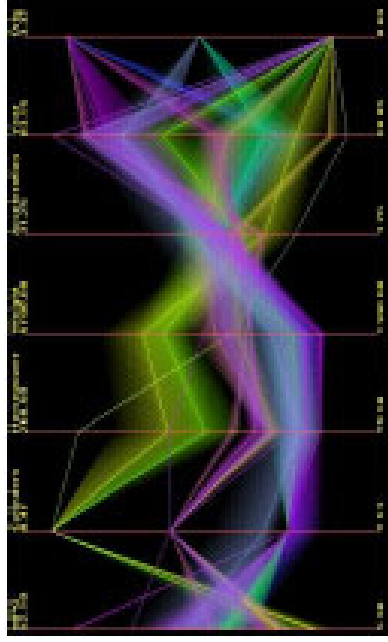Visualisation

# Examples of scientific visualisation

**Contributors: Matt Challacombe and Eric Schwegler**

This image shows a view of electrostatic potential iso-surfaces and a wireframe representation of the p53 tumor suppressor tetramerization monomer. Mutations in the p53 tumor suppressor are the most frequently observed genetic alterations in human cancer. The structure of the monomer's electrostatic potential has been rendered on an SGI workstation using iso-surfaces corresponding to -0.06 and +0.06 au. The electrostatic potential is widely implicated in molecular recognition, binding, and the enhanced diffusion of charged substrates.These results have been obtained from first principles electronic structure calculations using linear scaling Hatree-Fock theory recently developed at the University of Minnesota. Involving 3836 basis functions, this calculation was performed in 3 cpu days on an IBM RS6000 model 590 workstation, and is the largest Gaussian-based *ab initio* calculation performed to date.



The **pseudocolor plot** (right) is used to map temperature to color on the same planar slice.

AVS Express
Paraview - free !
Tecplot
IBM Data Explorer
VisIt - free

# Information Visualisation

Displaying information to help the user understand it better. Abstraction of data.



Information Visualization image shown courtesy of Matt Ward of Worcester Polytechnic Institute (WPI).

Example above I would categorise as **Data Visualisation**

The London Tube map I would categorise as **Information Visualisation** – recommend you read Edward Tufte
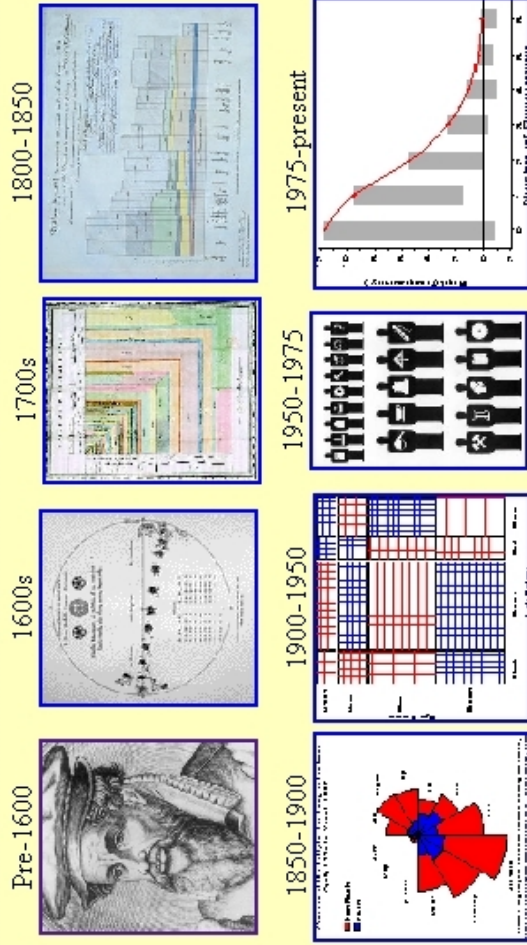


SciVis  - late '80s
InfVis – late '90's

This is a vast new field - especially important for **data mining**

File   Edit   View   Go   Bookmarks   Tools   Help

http://www.math.yorku.ca/SCS/Gallery/milestone/   Go

Getting Started   Latest Headlines

# Milestones in the History of Thematic Cartography, Statistical Graphics, and Data Visualization

An illustrated chronology of innovations

by Michael Friendly and Daniel J. Denis

Pre-1600

1600s

1700s

1800-1850

1850-1900

1900-1950

1950-1975

1975-present

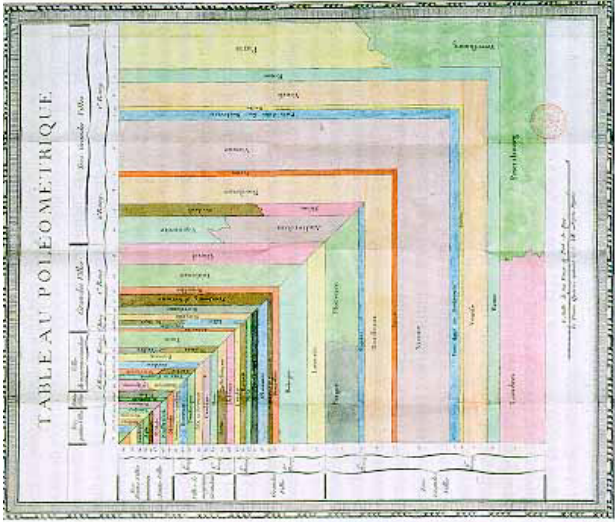| Up: Gallery | Introduction | Related | References | Term Index | Category XRef |
| --- | --- | --- | --- | --- | --- |
| Pre-1600 | 1600s | 1700s | 1800+ | 1850+ | 1900+ |
|  |  |  | 1950+ | 1975+ | Search |

See also:

- This document in PDF form, with active links. (You need Adobe Acrobat Reader)
- Chapter on the Milestone Project in C. Weihs and W. Gaul (eds.), *Classification– The Ubiquitous Challenge*, Springer, 2005.
- Images from the JSM 2002 Technical Poster Session [Thanks to Andy Mauromoustakos!]:
   ◦ Image1 (864 x 648; 123K);
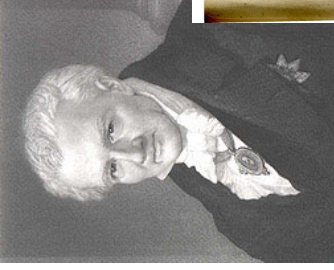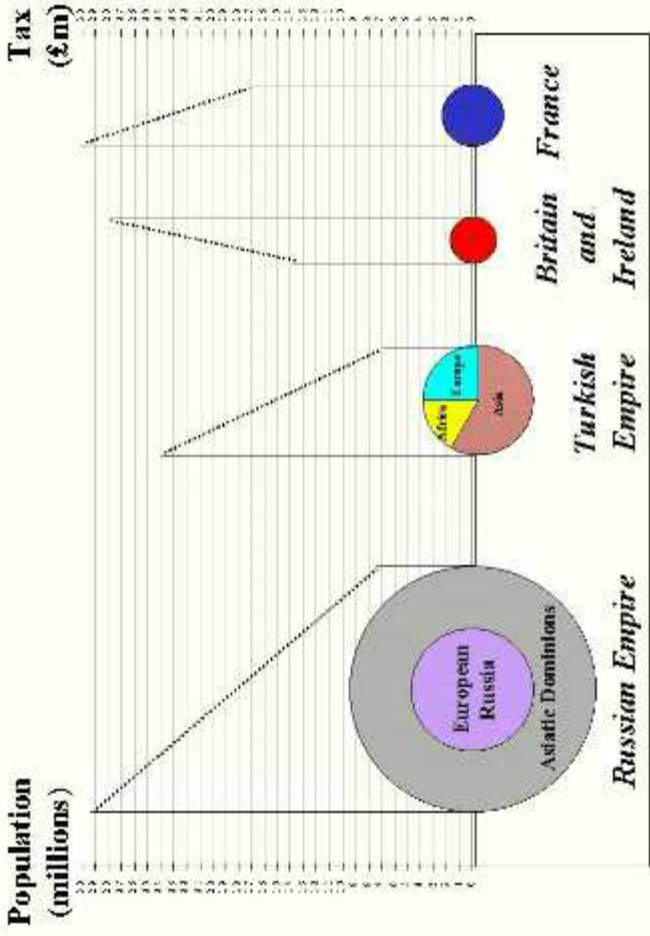   ◦ Image2 (864 x 648; 124K).

This web version is dedicated to Arthur H. Robinson (1915-2004), who inspired and encouraged our interest; to Antoine de Falguerolles, who initiated it, and to *les Chevaliers des*
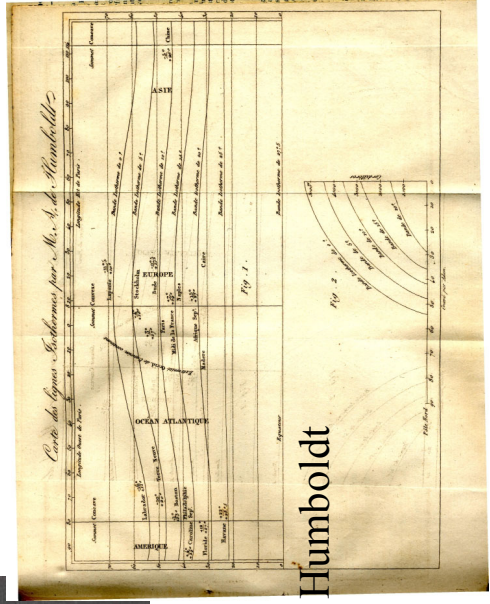
Done

Tax
(£m)

Population
(millions)

France

Britain
and
Ireland

Turkish
Empire

Russian Empire

European
Russia

Asiatic Dominions

William Playfair's chart
Invented pie-chart. 1801

Humboldt's
isotherms
1817



Alexander Von Humboldt

Luke Howard 1800
1st use of coordinate paper
in a research paper

TABLEAU POLEOMETRIQUE

Charles de Fourcroy
1782
"proportional squares"

Prof. Dalitz FRS – 1925-2006



All UT Constraints

34

Different statistical approaches

Tim Gershon, IoP Particle Physics 2006

Status of the CKM matrix

© Spanier

$m^2(\pi^0\eta)$ / GeV²

$m^2(\pi^0\eta)$ / GeV²

Dalitz Plot 1954

Some data visulisations from particle physics

Excluding event displays!

# 1975-present

## 9. 1975-present: High-D data visualization

It is harder to provide a succinct overview of the most recent developments in data visualization, because they are so varied, have occurred at an accelerated pace, and across a wider range of disciplines. It is also more difficult to highlight the most significant developments (and because we have focused on the earlier history), so there are presently areas and events unrepresented here.

With this disclaimer, a few major themes stand out:

- the development of a variety of highly interactive computer systems and more importantly,

- new paradigms of direct manipulation for visual data analysis (linking, brushing, selection, focusing, etc.)

- new methods for visualizing high-dimensional data (grand tour, scatterplot matrix, parallel coordinates plot, etc.);

- the invention of new graphical techniques for discrete and categorical data (fourfold display, sieve diagram, mosaic plot, etc.), and analogous extensions of older ones (diagnostic plots for generalized linear models, mosaic matrices, etc.) and,

- the application of visualization methods to an ever-expanding array of substantive problems and data structures.

These developments in visualization methods and techniques arguably depended on advances in theoretical and technological infrastructure. Some of these are: (a) large-scale software engineering; (b) extensions of classical linear statistical modeling to wider domains; (c) vastly increased computer processing speed and capacity, allowing computationally intensive methods and access to massive data problems.

In turn, the combination of these themes and advances now provides some solutions for earlier problems.

Milestones in the history of thematic cartography, statistical graphics and data visualisation – M. Friendly and D. Denis Jan 2006

Big thankyou to Michael Friendly website
http://www.math.yorku.ca/SCS/StatResource.html

**1975 to now High D data visualisation**
Some key dates…selective list ..

1985 Alfred Inselberg **Parallel Coordinates**

1985 D. Asimov  **Grand Tour**

1985 DataDescription Inc. Paul Velleman Cornell - DataDesk

1987 A. Becker and W. Cleveland **Linking and Brushing**

1998 A. Buja, D. Asimov, C. Hurley, J. McDonald  XGobi

1990 **E. Wegman Statistical analysis and parallel coord. CrystalVision.**

1991 M. Friendly Mosaic Display and Categorical data

1999 L. Wilkinson "Grammar of Graphics"

Systemization of data and graphs and graph algebras in an OO framework.

**Particle Physics Data - a problem in the analysis of a huge amount of multivariate data**

What do we use ?  Histograms and scatterplots.  Sometimes use colour

Can one use the latest computer graphics technology or ideas that statisticians and computer scientists have dreamt up in the last decade…?

To illustrate, will use the "pollen dataset" to show use of parallel coordinates, brushing and pruning. The Grand Tour comes later……………………………
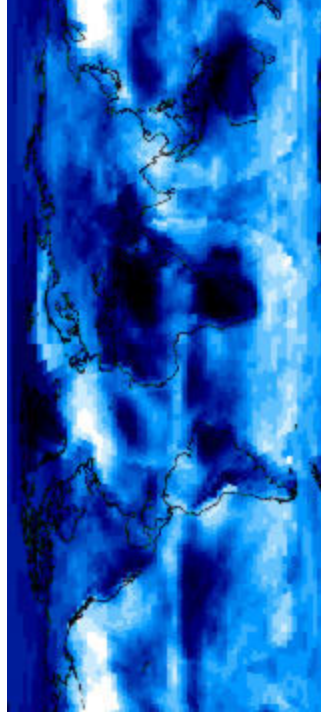
There are many other ideas ………………………………

# American Statistical Association

# Have challenges each year
# This is the 2006 one

**Pollen Data Set**

the **data set** from the 1986 **JSM** Exposition's **dataset** and was assembled by David Coleman of RCA Labs

**JSM = Joint Statistical Meeting**

# Data Visualisation
# Software
CrystalVision - E. Wegman
GGobi
XmdvTool
Orange

- **The data set:** The data are geographic and atmospheric measures on a very coarse 24 by 24 grid covering Central America. The variables are: elevation, temperature (surface and air), ozone, air pressure, and cloud cover (low, mid, and high). With the exception of elevation, all variables are monthly averages, with observations for Jan 1995 to Dec 2000. These data were obtained from the NASA Langley Research Center Atmospheric Sciences Data Center (with permission; see important copyright terms below).
- More details about the data, including descriptions of the variables, are available here.
- Download the data as a gzipped tar ball or as a zip file.
- There is also a flyer available.
- **The question:** The aim of the Data Expo is to provide a *graphical* summary of important features of the data set. This is intentionally vague in order to allow different entries to focus on different aspects of the data. For example, the focus can be on: the fact that the data are multivariate, or time-series, or spatial; or the fact that the data contain missing values; or the focus could even be on the *process* of exploring the data.
- Some obvious general questions that could be answered are: What are the important relationships between the variables? Are there any important trends in the data? Are there any important groupings or clusters in the data? Are there any unusual locations or time periods in the data set?
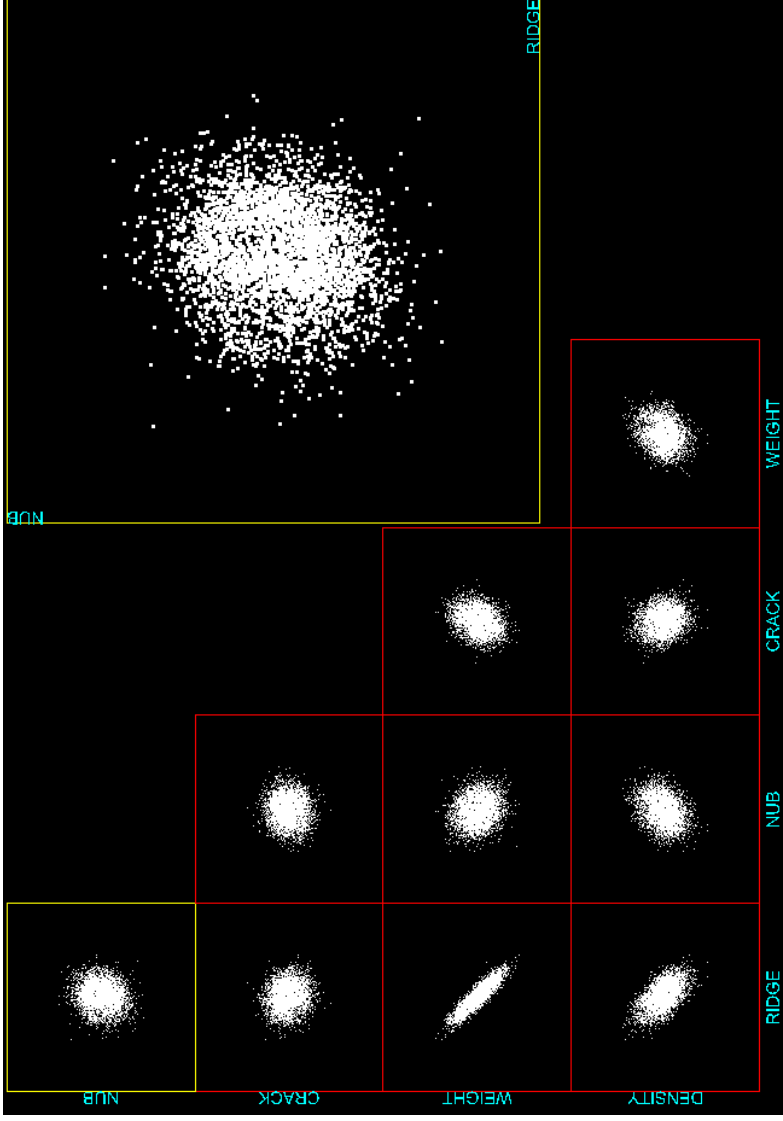
In graphics, a portion of each pixel's data that is reserved for transparency information. 32-bit graphics systems contain four channels -- three 8-bit channels for red, green, and blue (RGB) and one 8-bit alpha channel. The **alpha** channel is really a mask -- it specifies how the pixel's colors should be merged with another pixel when the two are overlaid, one on top of the other.

1) Try this on the pollen data set with CrystalVision

2) Now parallel coordinates.

Problem - how do you study an N-Dimensional space (N>2) when you only have a flat screen ?

This is one solution - with colour mixing (blending) and the alpha channel (transparency) - is very powerful
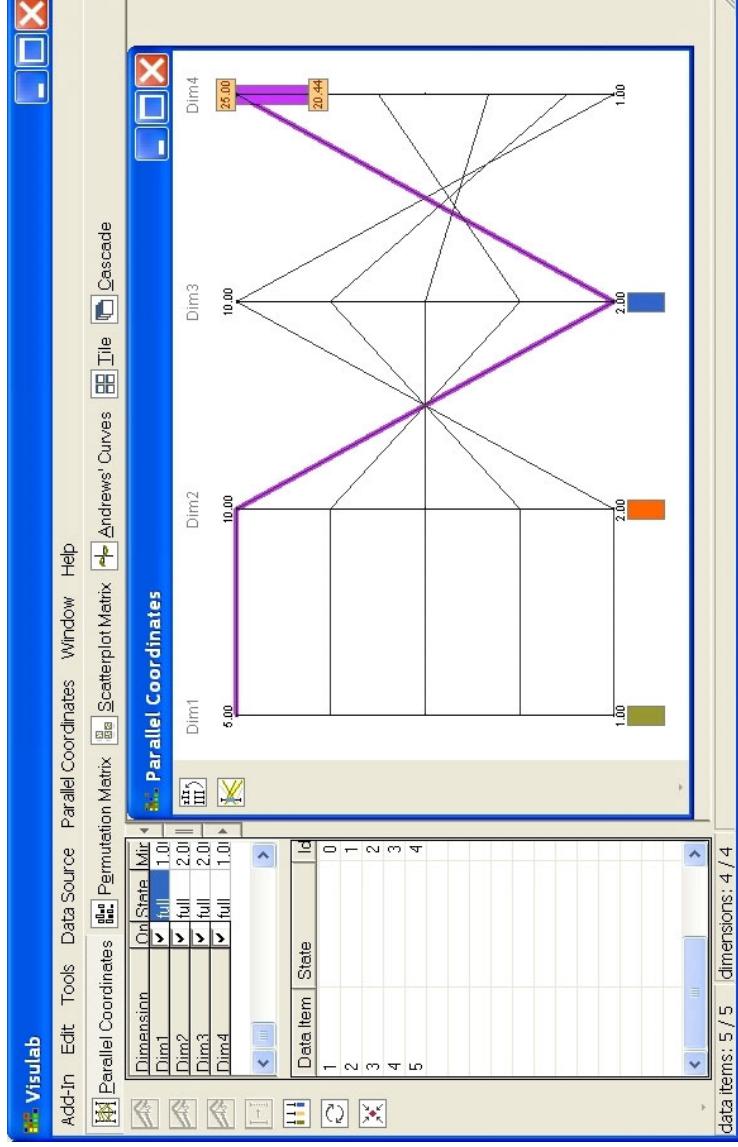
The pollen data - this is called a scatter matrix.
2D projections of this 5 variable space helps - but -

Greatly help matters using colour and the alpha channel

# Introduction to Parallel Coordinates

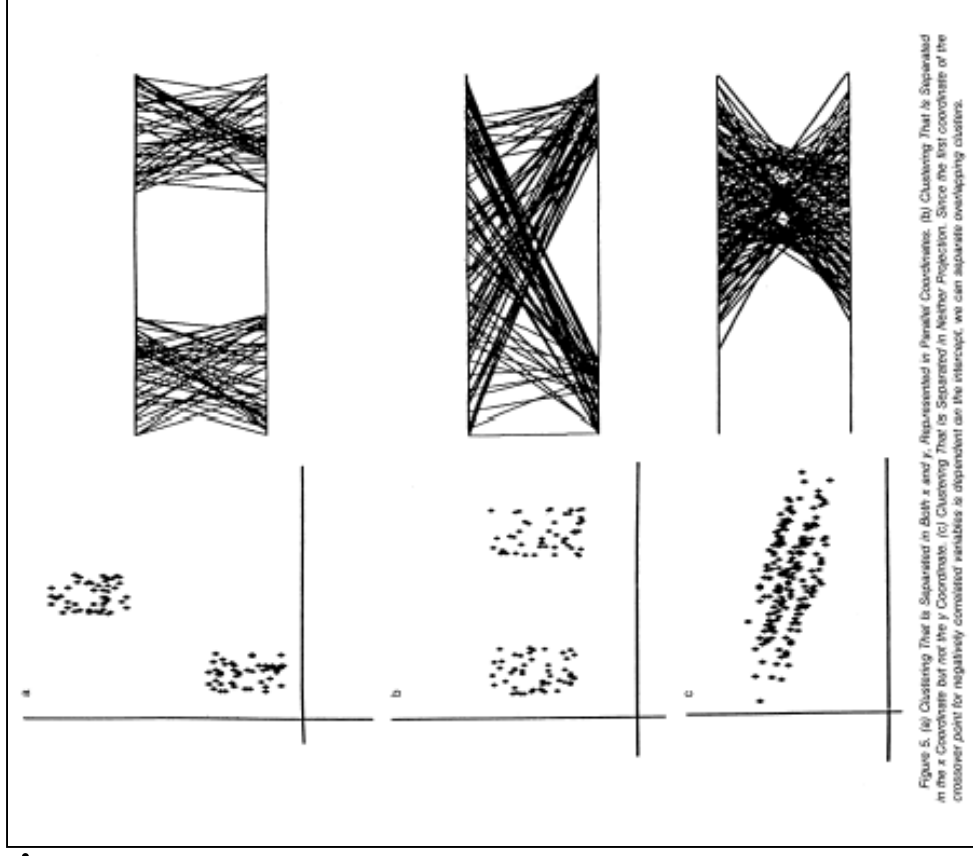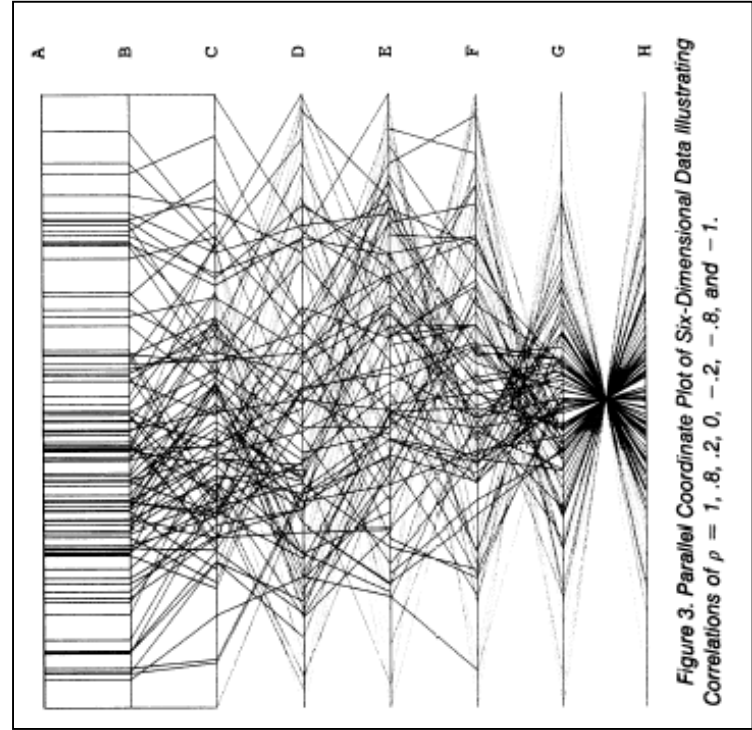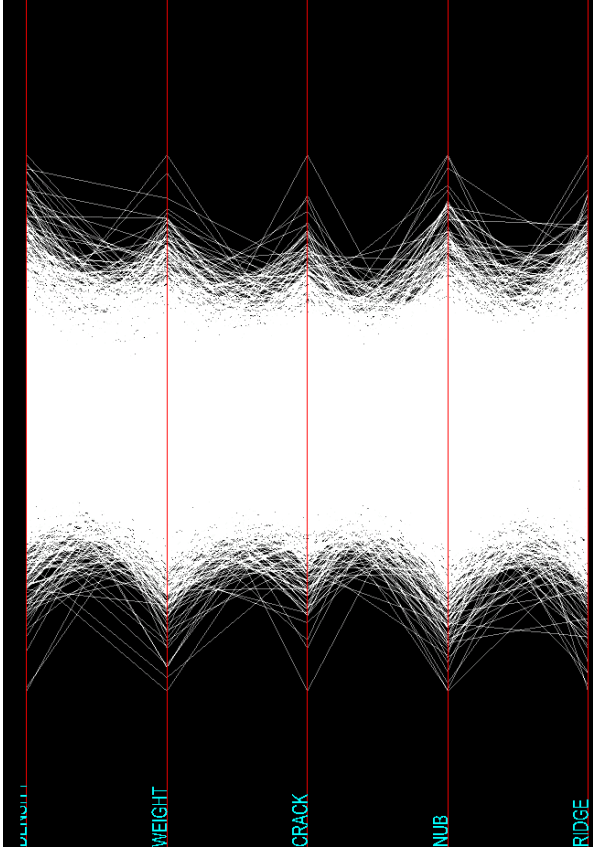| DataPoint | Dim1 | Dim2 | Dim3 | Dim4 |
|-----------|------|------|------|------|
| 1 | 1 | 2 | 10 | 1 |
| 2 | 2 | 4 | 8 | 4 |
| 3 | 3 | 6 | 6 | 9 |
| 4 | 4 | 8 | 4 | 16 |
| 5 | 5 | 10 | 2 | 25 |

Simple Implementation with EXCEL plugin
http://www.inf.ethz.ch/personal/hinterbe/Visulab/



This also shows the idea of brushing

**Hyperdimensional Data Analysis Using Parallel Coordinates**

Edward J. Wegman

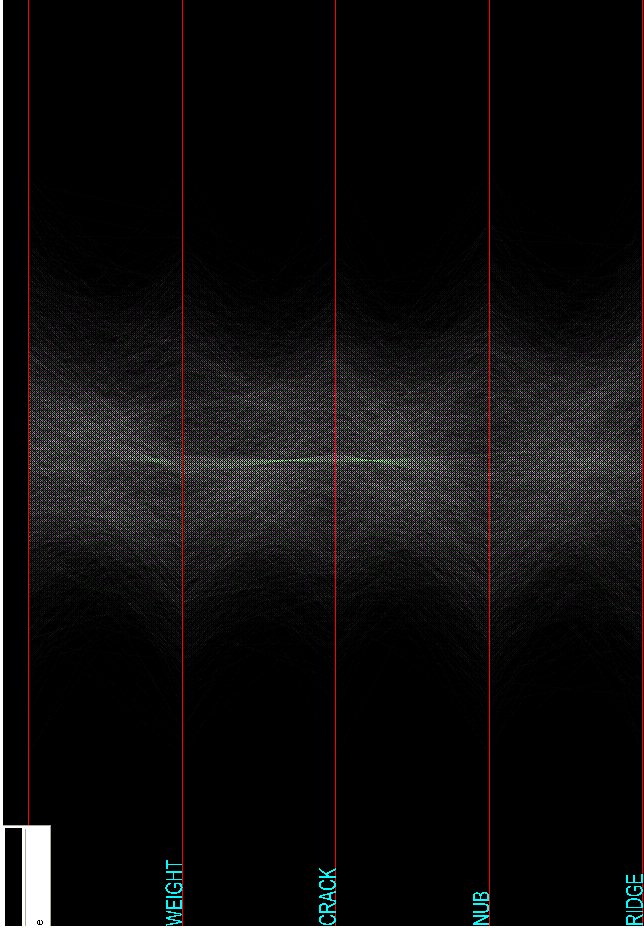*Journal of the American Statistical Association*, Vol. 85, No. 411 (Sep., 1990), 664-675.

Wegman has done much on the use of parallel coords.
Some useful things to note........



Figure 5. (a) Clustering That Is Separated in Both x and y, Represented in Parallel Coordinates. (b) Clustering That Is Separated in the x Coordinate but not the y Coordinate in Neither Projection. (c) Clustering That Is Separated in the x Coordinate but not the y Coordinate. Since the first coordinate of the crossover point for negatively correlated variables is dependent on the intercept, we can separate overlapping clusters.



Figure 3. Parallel Coordinate Plot of Six-Dimensional Data Illustrating Correlations of ρ = 1, .8, .2, 0, −.2, −.8, and −1.

Low alpha
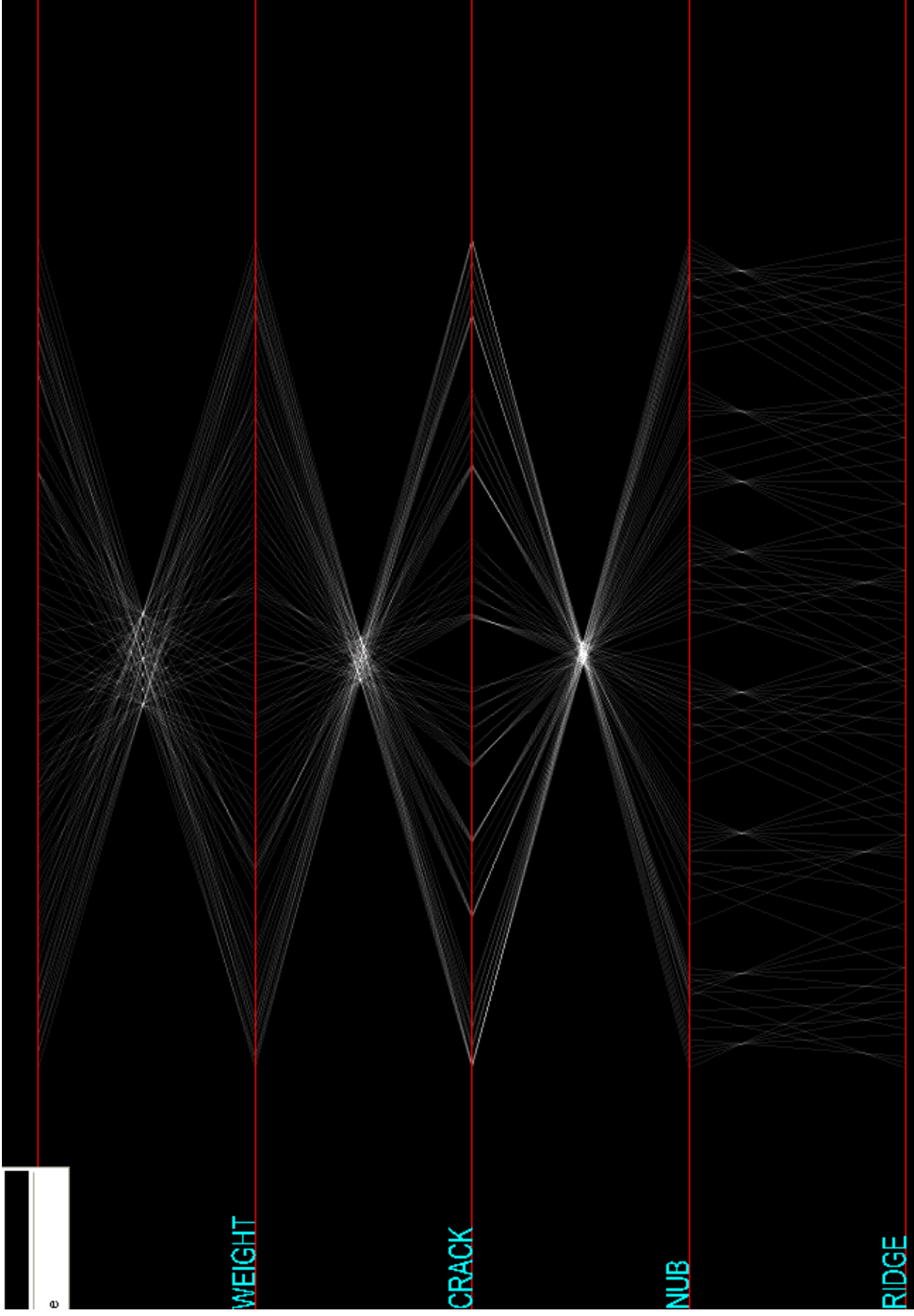
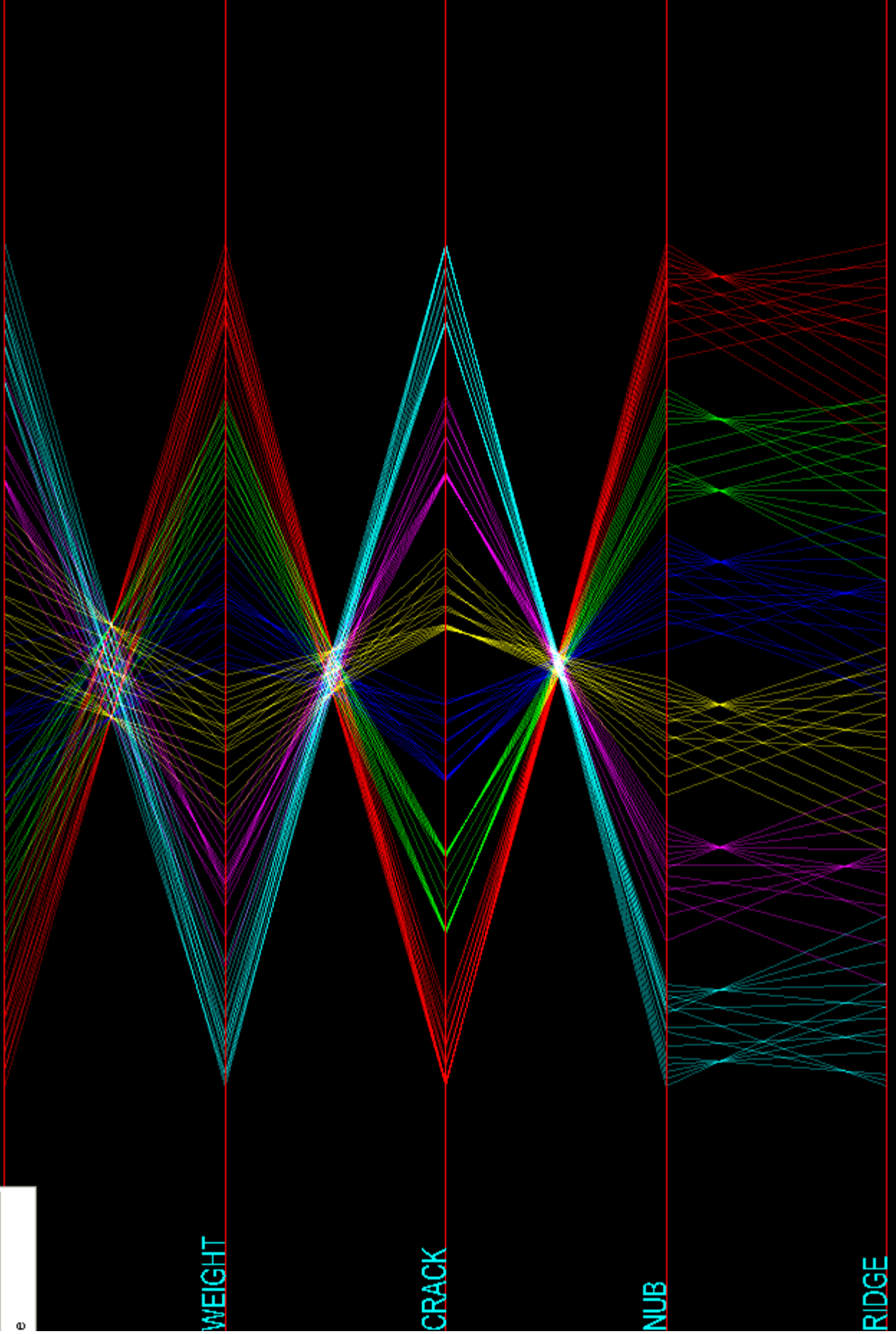High alpha

Now brushing - colour the data with chosen colours

and pruning - cut data you do not want

# First lets PRUNE



WEIGHT
CRACK
NUB
RIDGE

WEIGHT

CRACK

NUB

RIDGE

e
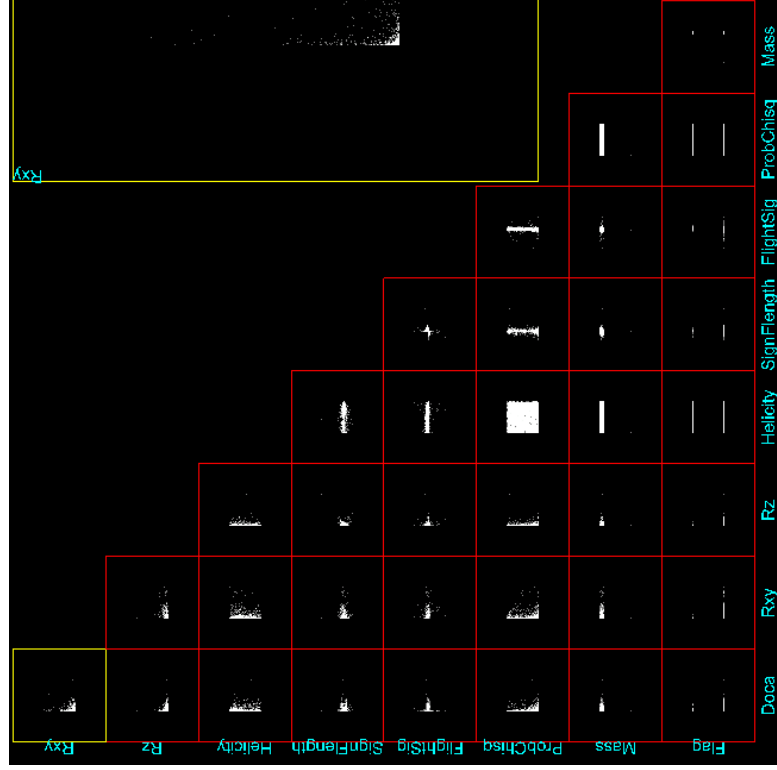
98/3848 points. S/B = 2.6%
There are other features in the data. See E. Wegman
Contrived example, but helps a newcomer to use this
type of graph.

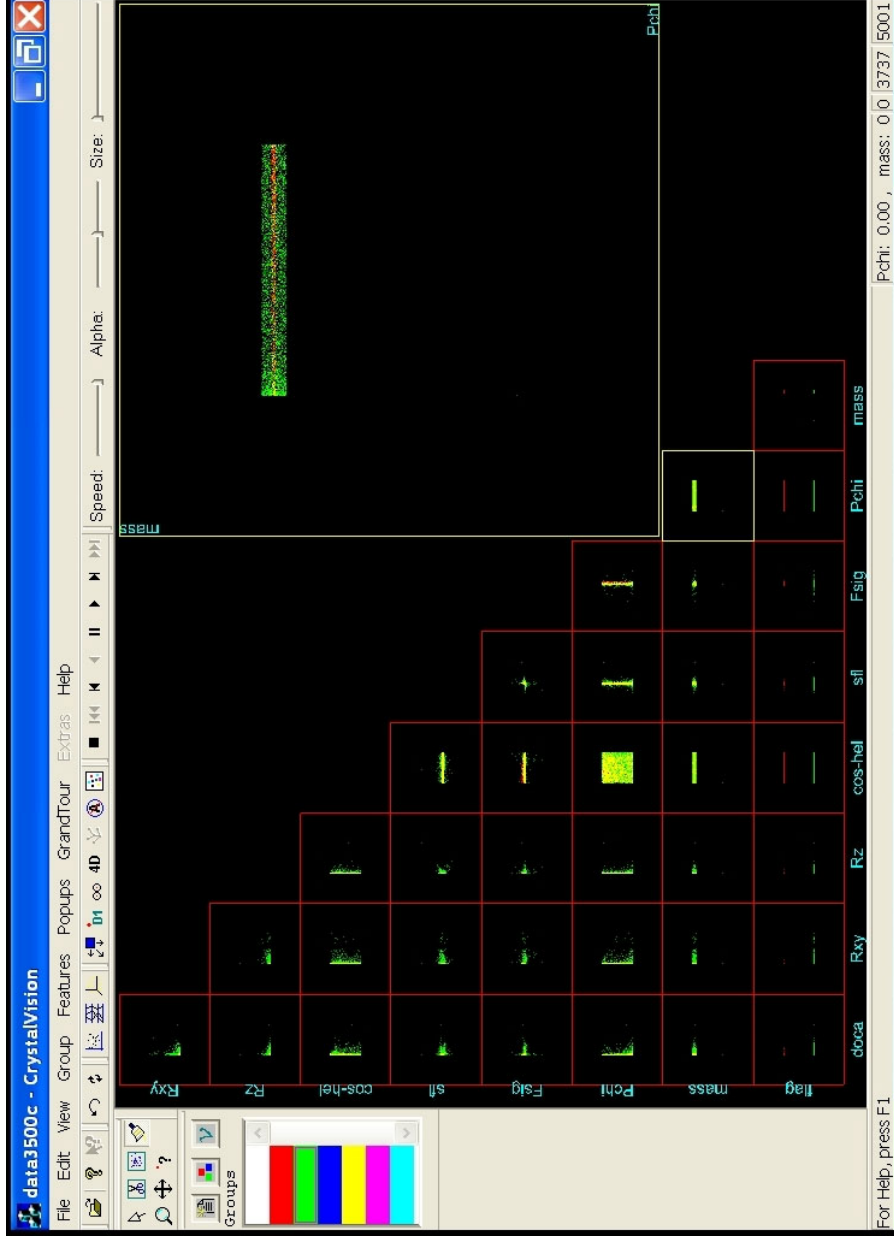Now lets try some particle physics monte carlo data

From Liliana Teodorescu - 1264 Kzero + 3734 background
(and a flag to tell us which is which !  Flag =1 S Flag=0 B
LT has shown how to use GEP on this dataset in another talk.

$K_s \rightarrow \pi^+ \pi^-$

Doca = distance of closest approach
Rxy radius of cylinder for interaction region
Rz abs. half length of cylinder defining the IR
Cos_hel abs. Value of cosine of Ks helicity angle
SFL – signed flight length
Fsig stat. Sig. Of Ks flight length
Pchi chisq prob of Ks vertex
Mass – reconstructed mass of the Ks



CrystalVision – E. Wegman
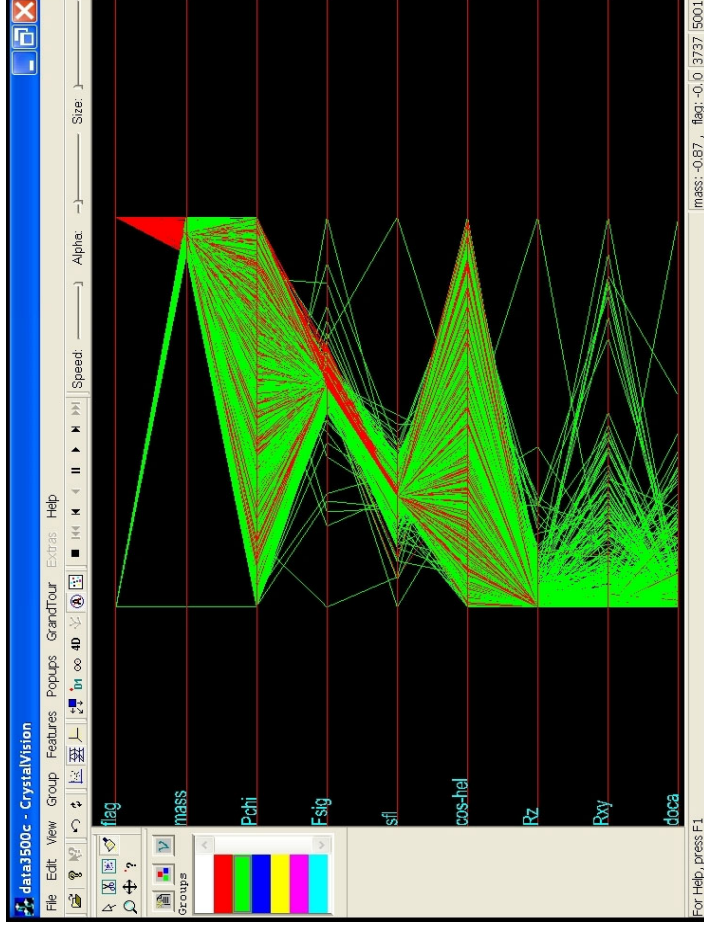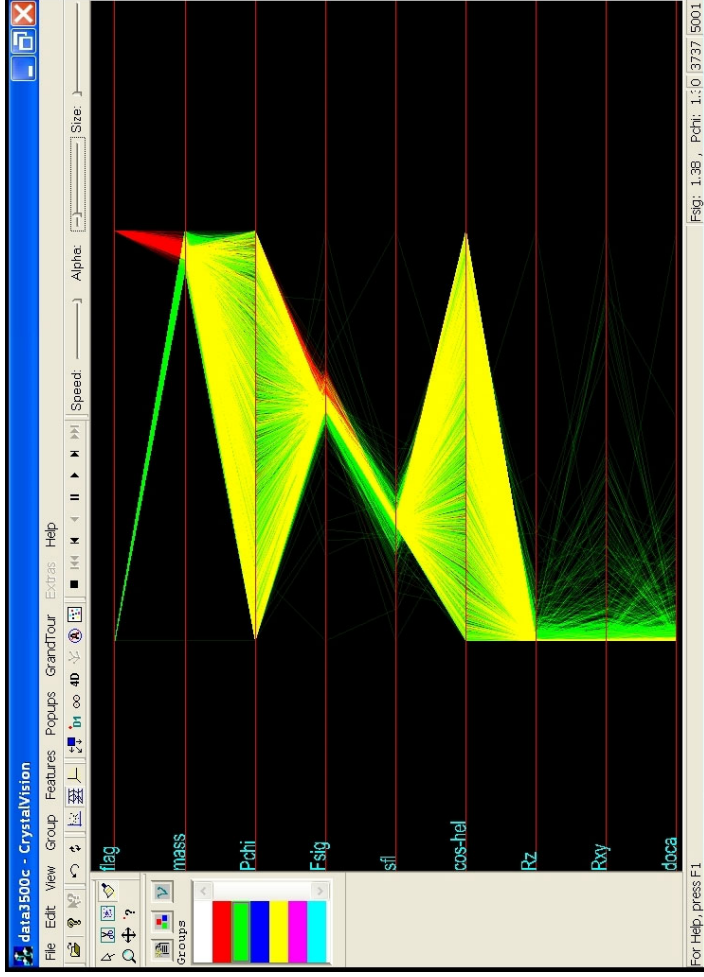
CrystalVision
(E. Wegman)

Has blending
and control of
intensity

VERY Powerful



Brush signal RED and background GREEN

If they overlap RED + GREEN = YELLOW (yellow)

Now go to parallel coordinates - adjust alpha

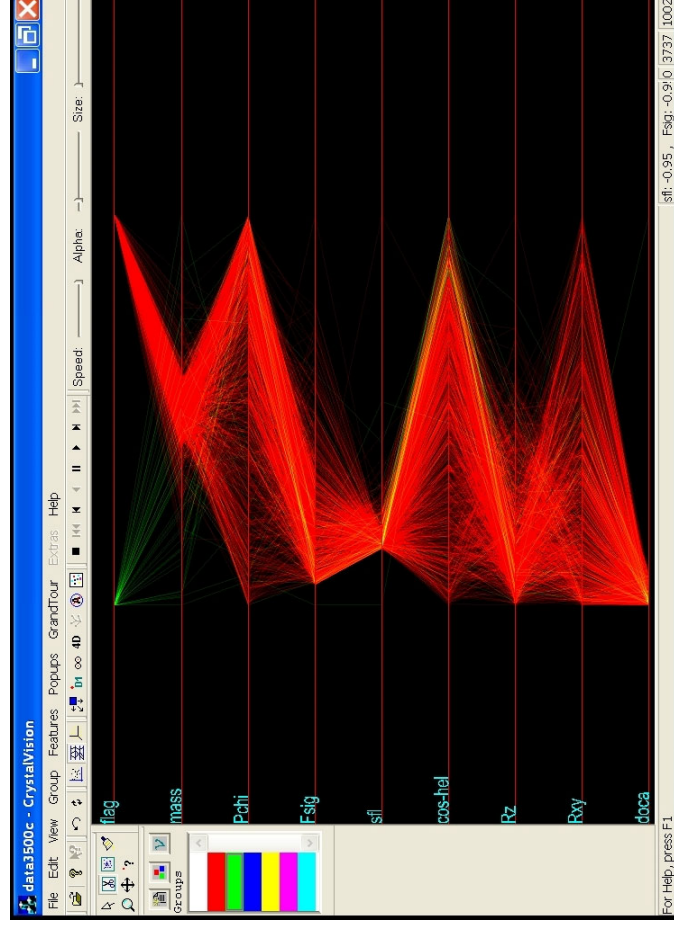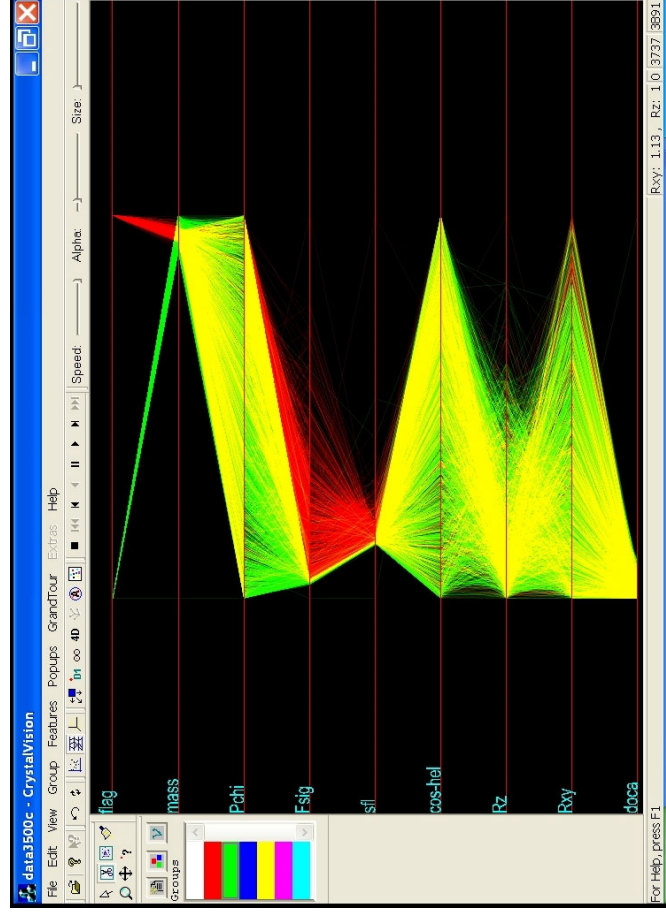Note: See affect of turning alpha channel on and off

Note:Parallel Coords Vertical. Sclaes data between min. and max.

Immediately see that $R_{xy}$, Doca (and sfl less so) discriminate the background
Only variable where signal can be seen is Fsig.
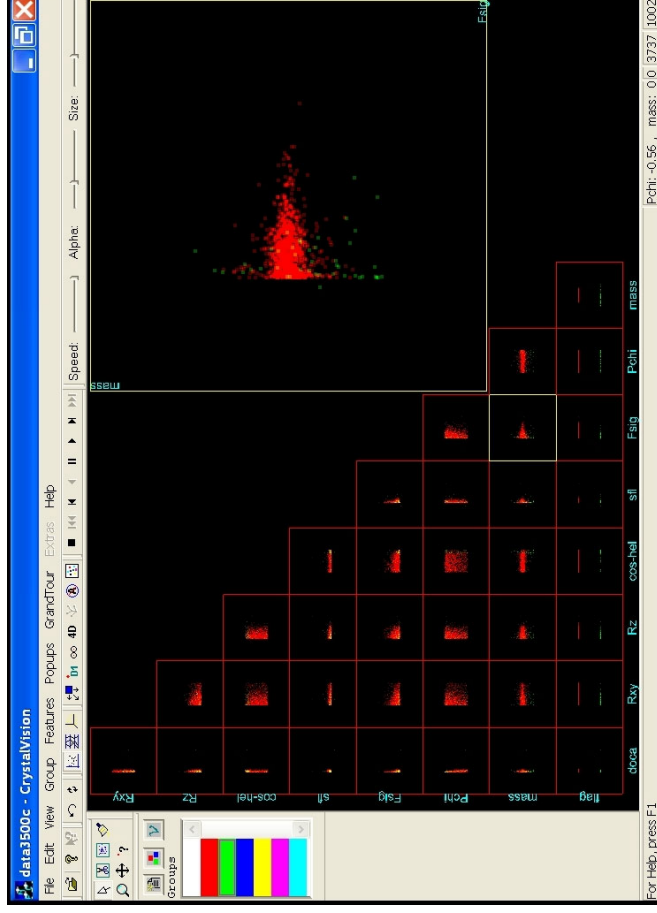
How to clean up this data - " what is the order of cuts ?"
Remove obvious background (Prune Doca and Rxy)
Then select signal (FSig)



Takes just a couple of minutes to do this…

# Back in scatterplot space



958 S 44 B   95% Purity 80% Efficiency

Did not spend long on this – Exploratory Visual Data Analysis

**Powerful way to decide which variables matter and the order in which cuts should be applied.**

**Precursor to machine learning approach ….to be continued**

**Multivariate classifiers**

An algorithm/tool that assigns a point in a multidimensional space to one of several possible categories e.g. signal or background.

- Cut    (binary split or stump)
- Decision Trees
- Support Vector Machine (SVM)
  (NOTE: radial basis function SVM equiv. to a type of Neural Net)
- Neural Net
- k nearest neighbours (kNN)
- VizRank
- Genetic Algorithm,  Genetic Programming,
  Gene Expression Programming (GEP) – can do other things also.

- etc..………………

# ViZRank – finding informative data projections in Functional genomics by machine learning
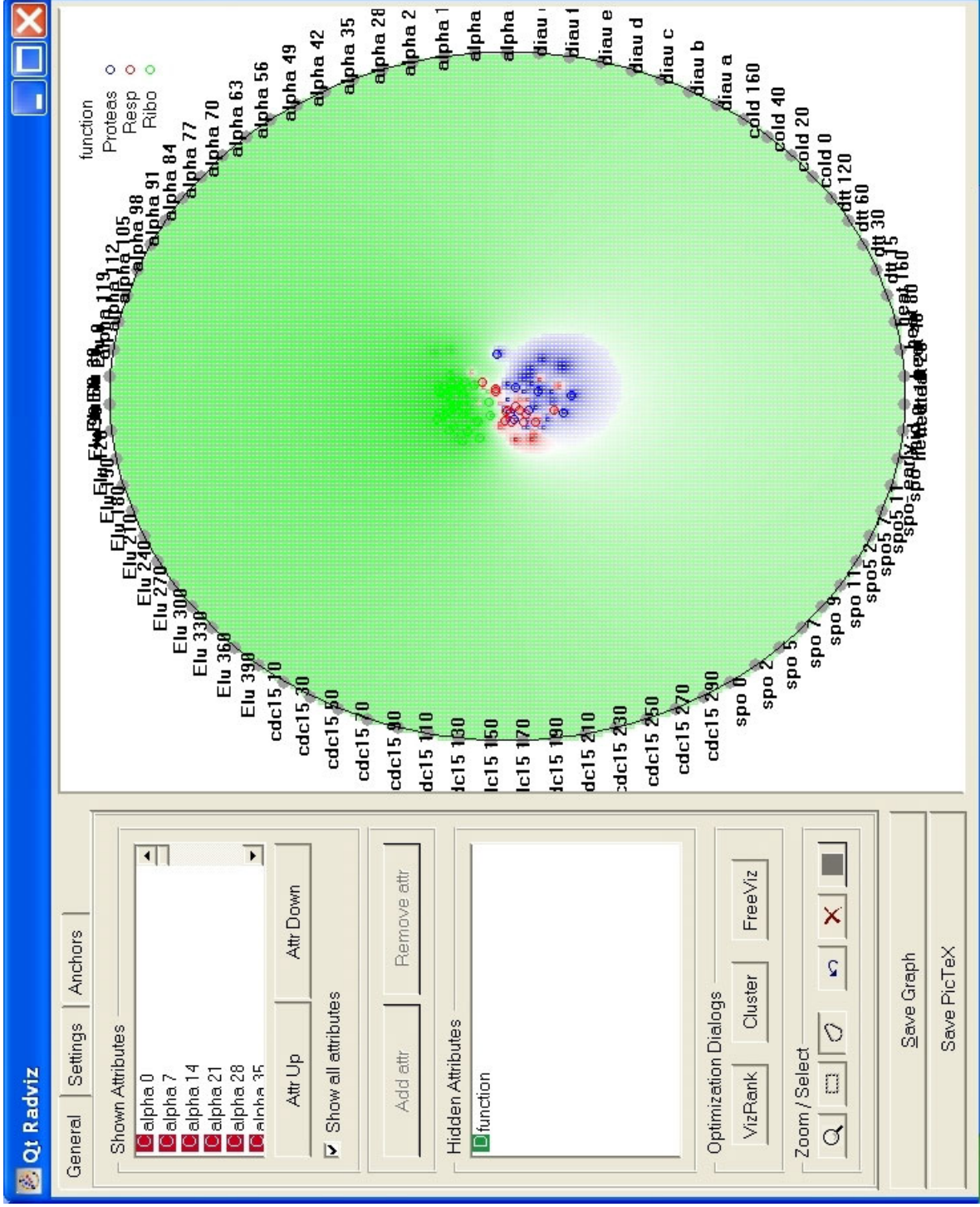
Gregor Leban et al.......Use ORANGE

Project N-dimension space onto a reduced number of dimensions (use radviz or polyviz visualisation)
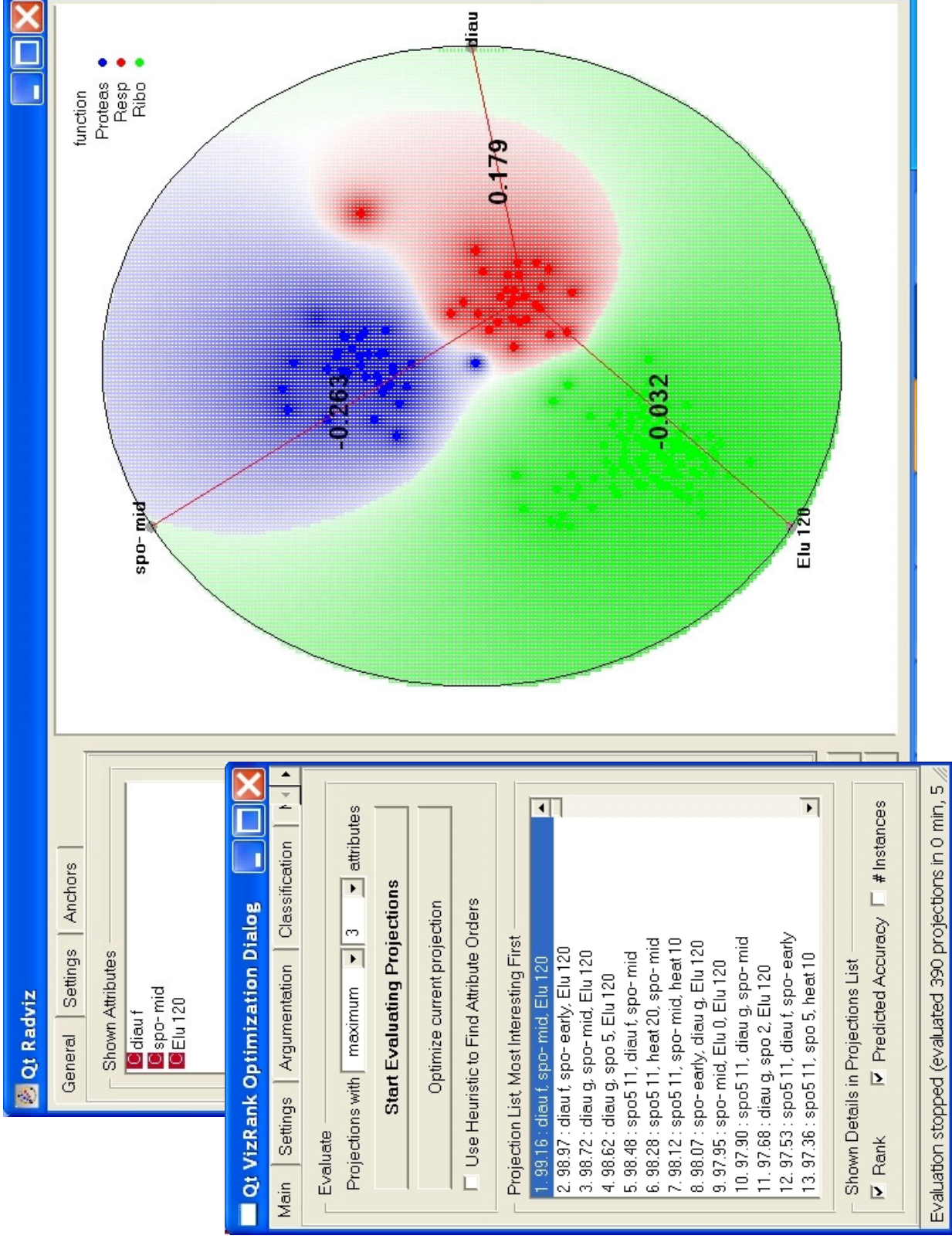
Rank projections using kNN classifier

Yeast *Saccharomyces cerevisiae*
79 different DNA microassay hybridization measurements
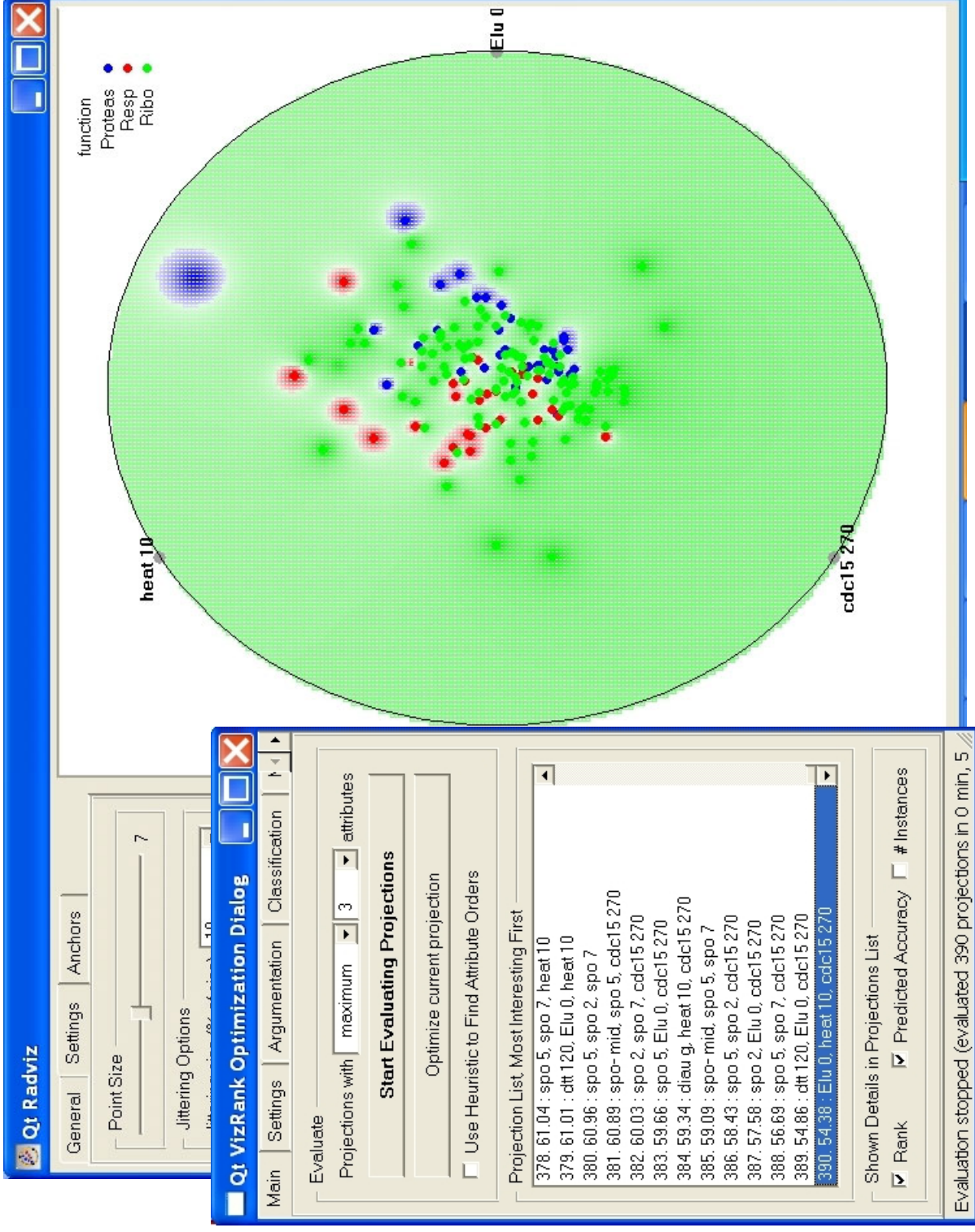Describe each gene.......which ones matter ?

# RadViz Visualisation – invented by P. Hoffman

# Rank 3 attribute projections 99 %

# Poorly ranked projection – 54%

Illustrate some other classifiers and visualisation using the wine dataset.....

# THE WINE DATA SET

4. Relevant Information: -- These data are the results of a chemical analysis of wines grown in the same region in Italy but derived from three different cultivars. The analysis determined the quantities of 13 constituents found in each of the three types of wines.

-- I think that the initial data set had around 30 variables, but for some reason I only have the 13 dimensional version. I had a list of what the 30 or so variables were, but a.) I lost it, and b.), I would not know which 13 variables are included in the set. !!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!

5. Number of Instances class 1 - 59 class 2 - 71 class 3 - 48

6. Number of Attributes 13

Use WEKA  - histogram of all the wine dataset

Once you get used to it, this plot contains a lot of information

Wine dataset.....

Flavanoids, colour, hue, alcohol, proline matter…..

VIZRank algorithm applied…..

50% selection........wines overlap....

Polyviz – radviz + value of variable indicated

# The GRAND tour

2D projections of an N-D space - choose suitable axes of rotation and an algorithm that ensures you explore all the space. (The maths is complicated – See E. Wegman or Asimov

**The Grand Tour via Geodesic Interpolation of 2-frames**[*]

Daniel Asimov and Andreas Buja[†]
Report RNR-94-004, February 1994

Facinating idea – useful for looking for clusters in data

# Grand Tour of the wine dataset 2

# Grand Tour of the wine dataset 3

# Find hyperplane that separates clusters

# Use RBF when plane will not work.



class 1, $y = +1$ ( *buys_computer = yes* )

class 2, $y = -1$ ( *buys_computer = no* )

$A_1$

$A_2$

large margin

Support vectors. The SVM finds the maximum separating hyperplane, that is, the one v maximum distance between the nearest training tuples. The support vectors are shown v a thicker border.

# Support Vector Machine – rather confusing name.

# For math. details see books referenced earlier.

http://www.csie.ntu.edu.tw/~cjlin/libsvm/

SVM-toy is a good introduction

# SVM works well on this dataset  – Grand Tour agrees

J. Platt (1998). "Fast Training of Support Vector Machines using Sequential Minimal Optimization". Advances in Kernel Methods - Support Vector Learning, B. Schoelkopf, C. Burges, and A. Smola, eds., MIT Press.

**Weka Explorer**

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Classifier

Choose  **SMO** -C 1.0 -E 1.0 -G 0.01 -A 250007 -L 0.0010 -P 1.0E-12 -N 0 -V -1 -W 1

Test options
- Use training set
- Supplied test set    Set...
- Cross-validation    Folds  10
- Percentage split    %  66

More options...

(Nom) wine

Start    Stop

Result list (right-click for options)
23:45:17 - functions.SMO
23:46:04 - functions.SMO

Status
OK

Classifier output

```
Time taken to build model: 0.11 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances         175               98.3146 %
Incorrectly Classified Instances         3                1.6854 %
Kappa statistic                          0.9745
Mean absolute error                      0.226
Root mean squared error                  0.279
Relative absolute error                 51.4678 %
Root relative squared error             59.5404 %
Total Number of Instances              178

=== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
  1        0.008      0.983        1        0.992      W1
  0.958    0          1            0.958    0.978      W2
  1        0.015      0.96         1        0.98       W3

=== Confusion Matrix ===

  a  b  c   <-- classified as
 59  0  0 |  a = W1
  1 68  2 |  b = W2
  0  0 48 |  c = W3
```

Now lets try some other machine learning algorithms on the particle physics dataset…

•Decision Tree   C4.5
• VizRank
• SVM
• Neural Net with back propagation

**SVM DOES NOT WORK WELL GRAND TOUR EXPLAINS WHY.**

# DECISION TREES

# BIG SUBJECT

# FOR WINE DATASET

---

Flavanoids
Colour
Proline

94% selection

---

**Weka Classifier Tree Visualizer: 00:01:03 - trees.J48 (wine-weka.filters....**

Tree View

Flavanoids
- <= 1.57 → Colour Intensity
  - <= 3.8 → W2 (13.0)
  - > 3.8 → W3 (49.0/1.0)
- > 1.57 → Proline
  - <= 720 → W2 (54.0/1.0)
  - > 720 → Colour Intensity
    - <= 3.4 → W2 (4.0)
    - > 3.4 → W1 (58.0)

---

**Weka Explorer**

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Classifier

Choose  J48 -C 0.25 -M 2

Test options
- ○ Use training set
- ○ Supplied test set    Set...
- ● Cross-validation    Folds  10
- ○ Percentage split    %  66

More options...

(Nom) wine

Start    Stop

Result list (right-click for options)
23:45:17 - functions.SMO
23:46:04 - functions.SMO
00:01:03 - trees.J48

Classifier output

Time taken to build model: 0.02 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances        167                93.8202 %
Incorrectly Classified Instances       11                 6.1798 %
Kappa statistic                         0.9058
Mean absolute error                     0.0486
Root mean squared error                 0.2019
Relative absolute error                11.0723 %
Root relative squared error            43.0865 %
Total Number of Instances             178

=== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
0.983     0.034     0.935       0.983    0.959       W1
0.944     0.056     0.918       0.944    0.931       W2
0.875     0.008     0.977       0.875    0.923       W3

=== Confusion Matrix ===

 a  b  c   <-- classified as
58  1  0 |  a = W1
 3 67  1 |  b = W2
 1  5 42 |  c = W3

Status
OK

# Decision Trees

Decision trees emerged in mid 80's: CART (Breiman, Friedman etc), C4.5 (Quinlan) etc

## Criteria used for commercial trees

(p = fraction of correctly classified events)

$$Q(p) = p$$
$$Q(p) = -2p(1-p) \qquad \text{Gini index}$$
$$Q(p) = p\log p + (1-p)\log(1-p) \quad \text{cross} - \text{entropy}$$

**Split nodes recursively until a stopping criterion is satisfied.**



S/B
52/48

< 100    PMT Hits?    ≥ 100

B
4/37

S/B
48/11

< 0.2 GeV    Energy?    ≥ 0.2 GeV

S/B
9/10

S
39/1

< 500 cm    Radius?    ≥ 500 cm

S
7/1

B
2/9

**Parent node with *W* events and correctly classified *p*\*W events is split into two daughters nodes iff**

$$WQ(p) < W_1\underline{Q}(p_1) + W_2\underline{Q}(p_2)$$

**Stopping criteria:**

- unable to find a split that satisfies the split criterion
- maximal number of terminal nodes in the tree
- minimal number of events per node

Output of a decision tree is discrete: 1 if an event falls into a signal node, 0 otherwise.

Ilya Narsky Caltech Seminar Oct 2005.

Key variables are: Fsig, Rxy, SFL, cos-hel, doca  95% selection
Remember the parallel coords analysis….

---

**Weka Explorer**

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Classifier

Choose   J48 -C 0.25 -M 2

Test options
- ○ Use training set
- ○ Supplied test set   Set...
- ● Cross-validation   Folds  10
- ○ Percentage split   %   66

More options...

(Nom) LFLAG

Start      Stop

Result list (right-click for options)
19:02:04 - functions.SMO
19:07:34 - functions.SMO
07:32:02 - trees.J48
07:33:13 - trees.J48

Status
OK

Classifier output

Time taken to build model: 0.78 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      4767       95.34   %
Incorrectly Classified Instances     233        4.66   %
Kappa statistic                     0.8754
Mean absolute error                 0.072
Root mean squared error             0.2011
Relative absolute error            19.0589 %
Root relative squared error        46.2633 %
Total Number of Instances           5000

=== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
0.893     0.026     0.92        0.893    0.906       TRUE
0.974     0.107     0.964       0.974    0.969       FALSE

=== Confusion Matrix ===

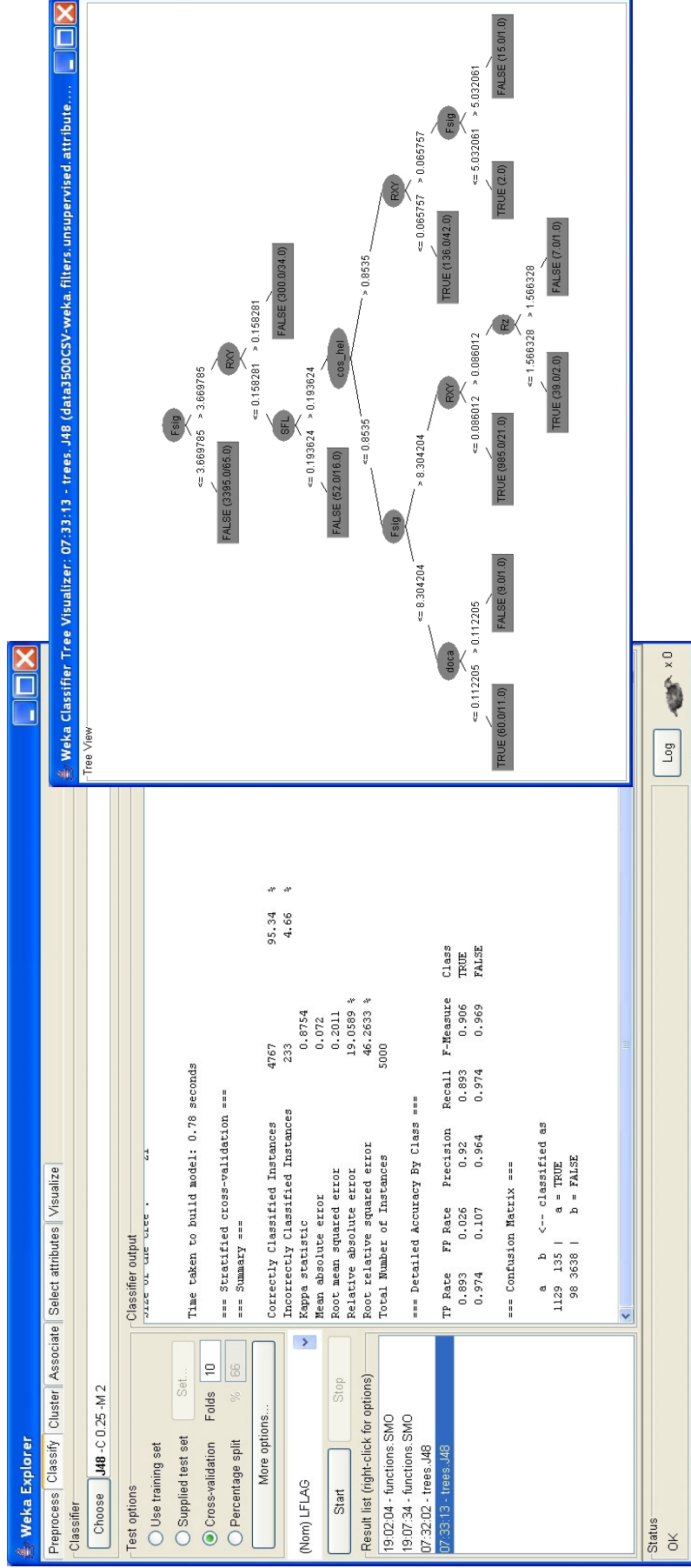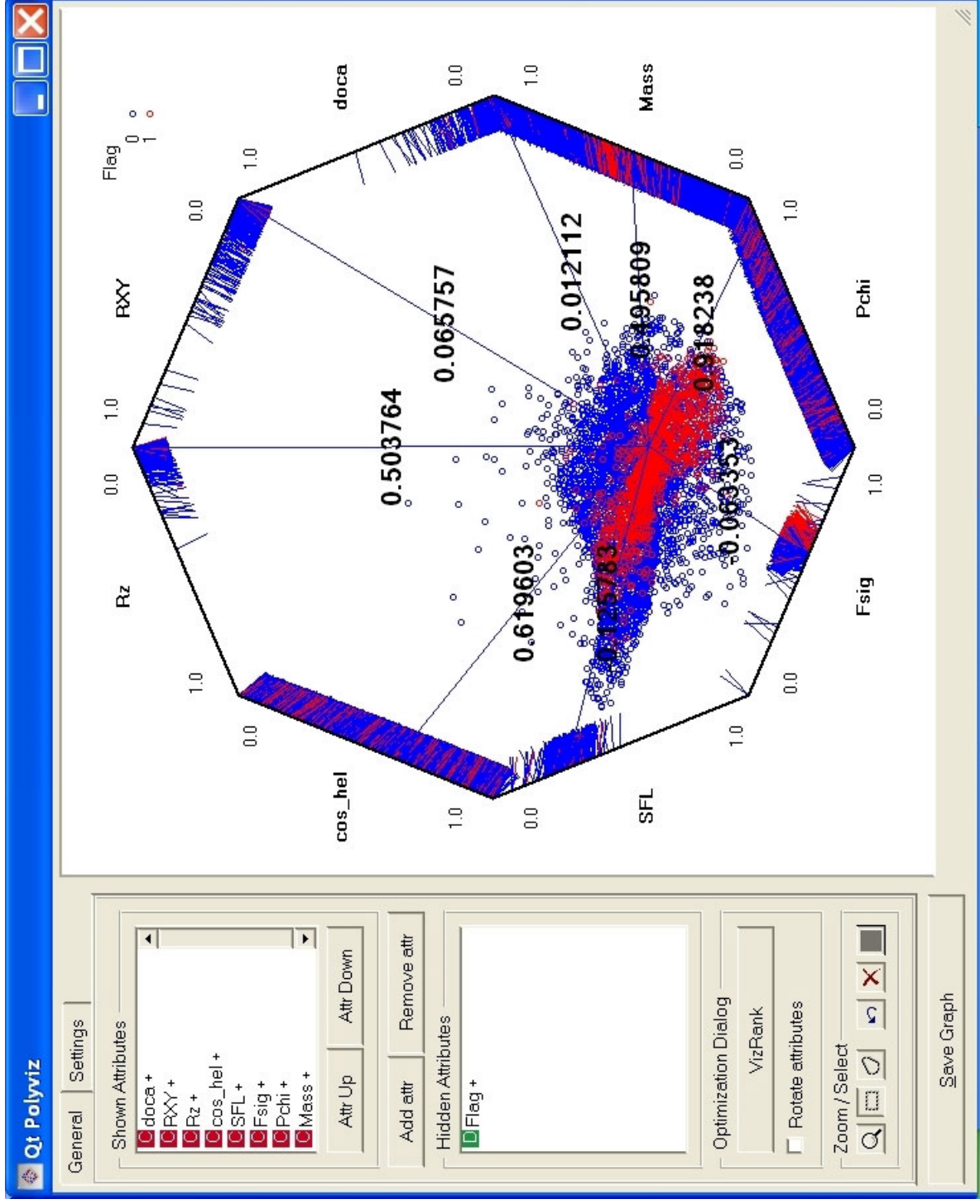   a    b   <-- classified as
1129  135 |   a = TRUE
  98 3638 |   b = FALSE

---

**Weka Classifier Tree Visualizer: 07:33:13 - trees.J48 (data3500CSV-weka.filters.unsupervised.attribute....**

Tree View

Fsig
<= 3.669785   > 3.669785

FALSE (3395.0/65.0)

Rxy
<= 0.158281   > 0.158281

SFL
<= 0.193624   > 0.193624

FALSE (52.0/16.0)

FALSE (300.0/34.0)

cos_hel
<= 0.8535   > 0.8535

Fsig
<= 8.304204   > 8.304204

doca
<= 0.112205   > 0.112205

TRUE (60.0/1.0)     FALSE (9.0/1.0)

Rxy
<= 0.086012   > 0.086012

TRUE (985.0/21.0)

Rz
<= 1.566328   > 1.566328

TRUE (39.0/2.0)

Rxy
<= 0.066757   > 0.066757

TRUE (136.0/42.0)

Fsig
<= 5.032061   > 5.032061

TRUE (2.0)     FALSE (15.0/1.0)

FALSE (7.0/1.0)

Log     x 0

Useful way to see what variables matter........

# VizRank (kNN)   Selects Fsig, Rxy, SFL, doca  88%

SVM (RBF not used) works very badly on this data
Grand Tour backs this conclusion up........

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Classifier

Choose  **SMO** -C 1.0 -E 1.0 -G 0.01 -A 250007 -L 0.0010 -P 1.0E-12 -N 0 -V -1 -W 1

Test options
- Use training set
- Supplied test set
- Cross-validation   Folds 10
- Percentage split   % 66

More options...

(Nom) LFLAG

Start    Stop

Result list (right-click for options)
19:02:04 - functions.SMO
19:07:34 - functions.SMO
07:32:02 - trees.J48
07:33:13 - trees.J48
07:43:28 - functions.SMO

Classifier output

Time taken to build model: 31.01 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances        3813              76.26   %
Incorrectly Classified Instances      1187              23.74   %
Kappa statistic                          0.0966
Mean absolute error                      0.2374
Root mean squared error                  0.4872
Relative absolute error                 62.831  %
Root relative squared error            112.1071 %
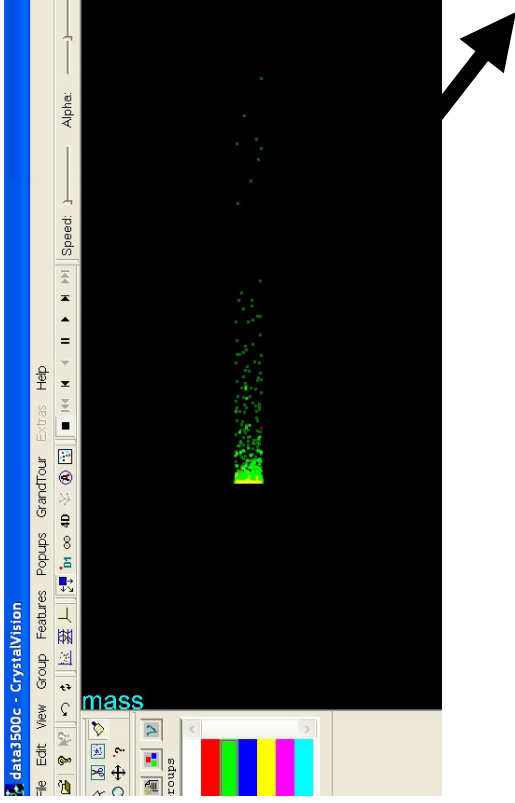Total Number of Instances             5000

=== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
0.07      0.003     0.881       0.07     0.13        TRUE
0.997     0.93      0.76        0.997    0.863       FALSE

=== Confusion Matrix ===

   a    b   <-- classified as
  89 1175 |   a = TRUE
  12 3724 |   b = FALSE

Status
OK                                                                    Log    x 0

Mass v Rxy
Standard Projection

CrystalVision
GrandTour

Cannot separate signal

# SUMMARY

|     | SF | BF |
|-----|-----|-----|
| **S** | | |
| **B** | | |

C4.5  1129  135  **95%**
      98  3638

VizRANK  **88%**

SVM  89  1175  **76%**
     12  3724

NN with backpropagation  927  337  **90%**
                         157  3579

Bagging  1145  119  **96%**
         96  3640

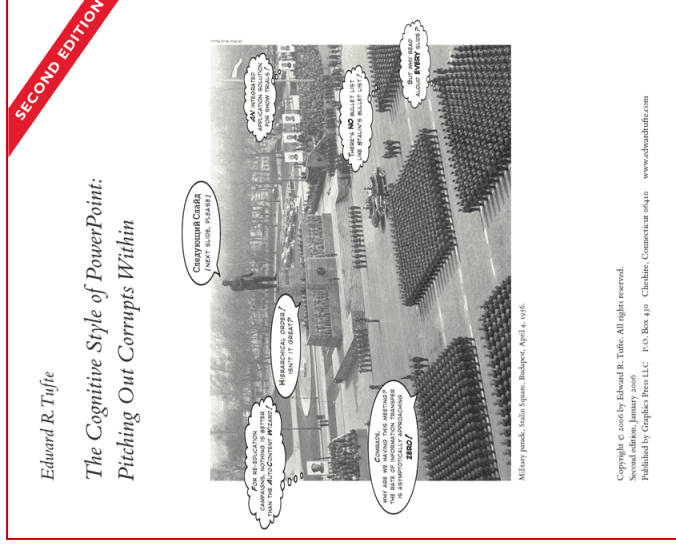Exploratory data analysis with crystalvision  968 S 44 B !!

VISUAL DATA ANALYSIS HELPS ONE TO UNDERSTAND THESE RESULTS

# SOME OTHER USEFUL DATA VISUALISATIONS

How can you display the data in a way that it can be easily understood ?

https://www.edwardtufte.com/tufte/





Should be read by all graduate students !

For last 100 years mathematicians have been suspicious of pictures and visual proofs.……perhaps this is changing.……

Synthese Lib Vol 327, 2005



Visualisations of some mathematical surfaces

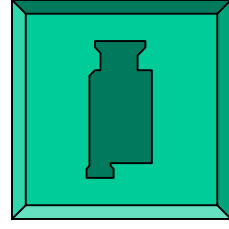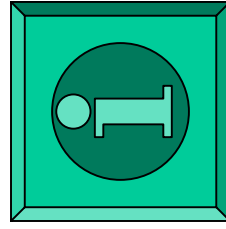http://vmm.math.uci.edu/3D-XplorMath/Surface/gallery.html

# Sphere Does Elegant Gymnastics in New Video,

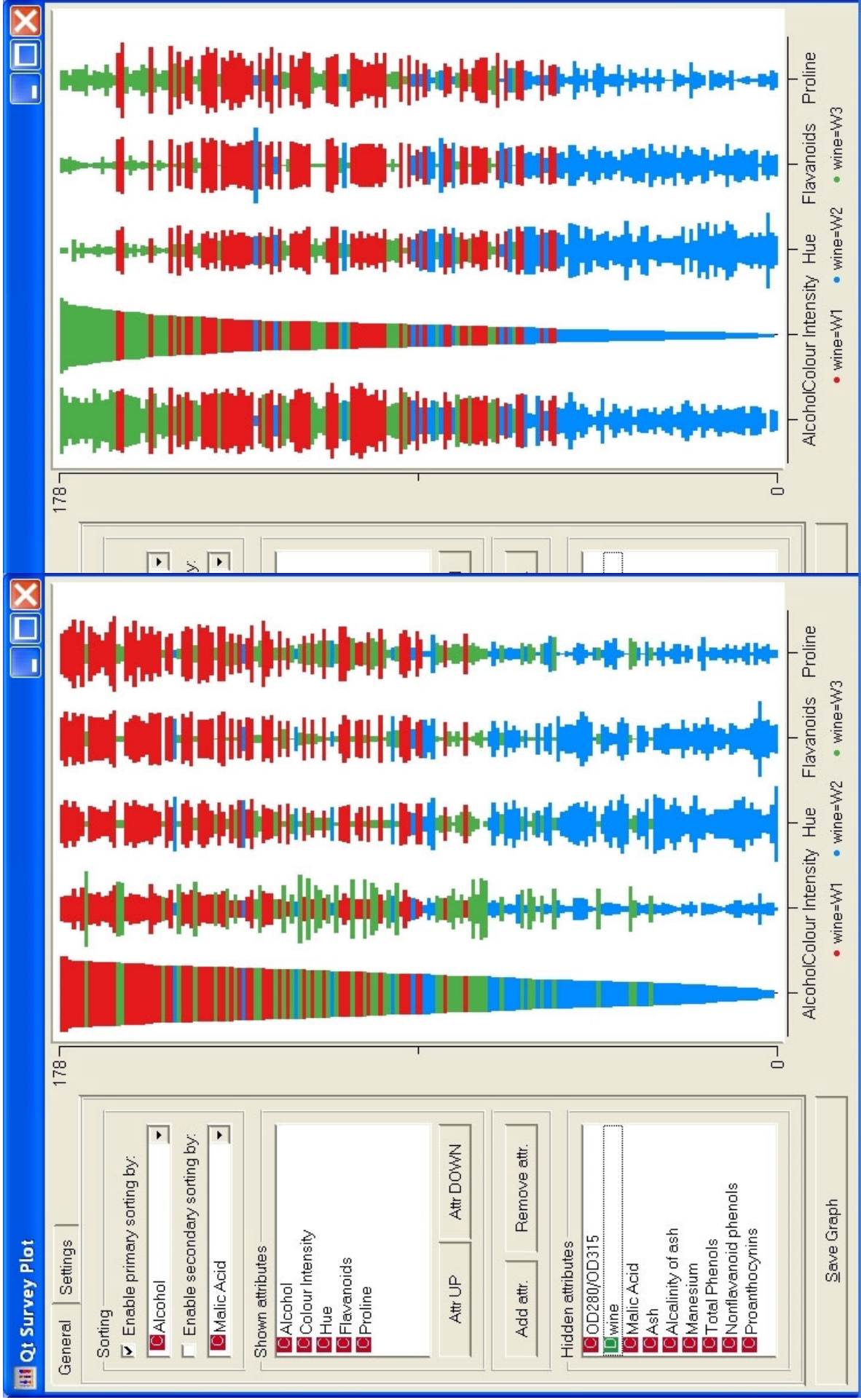Dana Mackenzie. Science 281:5377 (July 1998) 634-635.

## Is it possible to turn a sphere inside out without tearing or creasing it? The answer is yes.
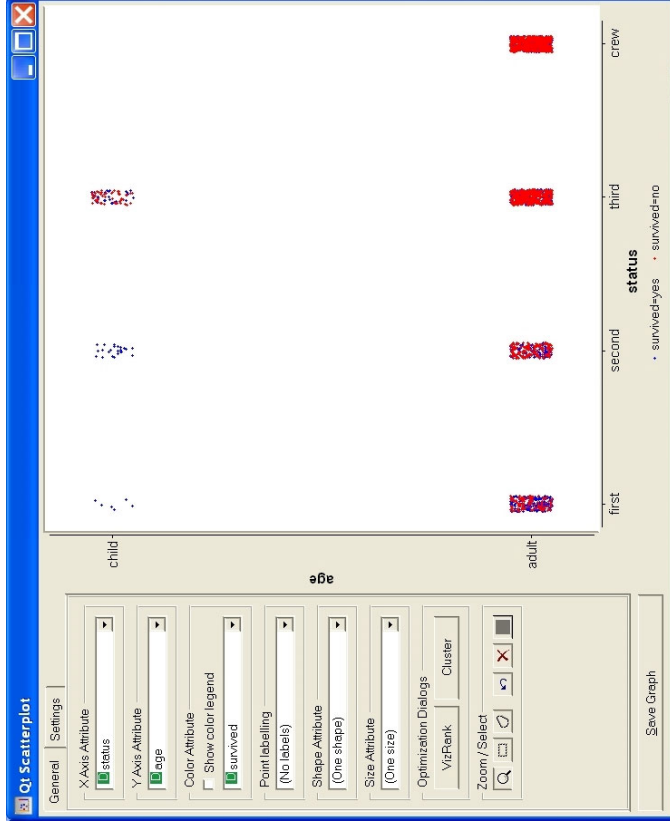
The proof, by Stephen Smale, is over forty years old but until now there has not been a satisfying demonstration. In this interesting article, the author explains how George Francis and John Sullivan at the University of Illinois were able to use a surface created by Robert Kusner in 1983 to create a 6 /2 minute computer animation of the eversion process. During the process, it is necessary to keep the energy level (defined so that the energy increases as more bending takes place) at a minimum at each stage. It was noted that the surface Kusner created had an energy level that would make it a candidate for the halfway point of the eversion process. Once the halfway point was known, it seemed possible to go backward to the sphere and forward to the sphere turned inside out. Francis and Sullivan showed with their animation that this is indeed the case. Their video has been shown recently at the Siggraph 98 convention in Orlando and the International Congress of Mathematicians in Berlin.

http://new.math.uiuc.edu/optiverse/
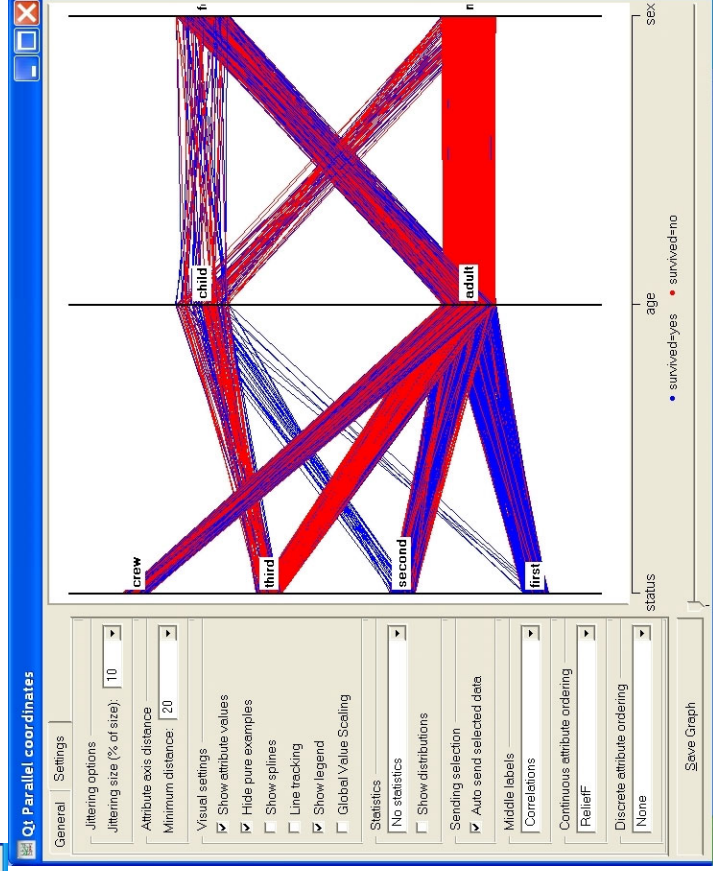
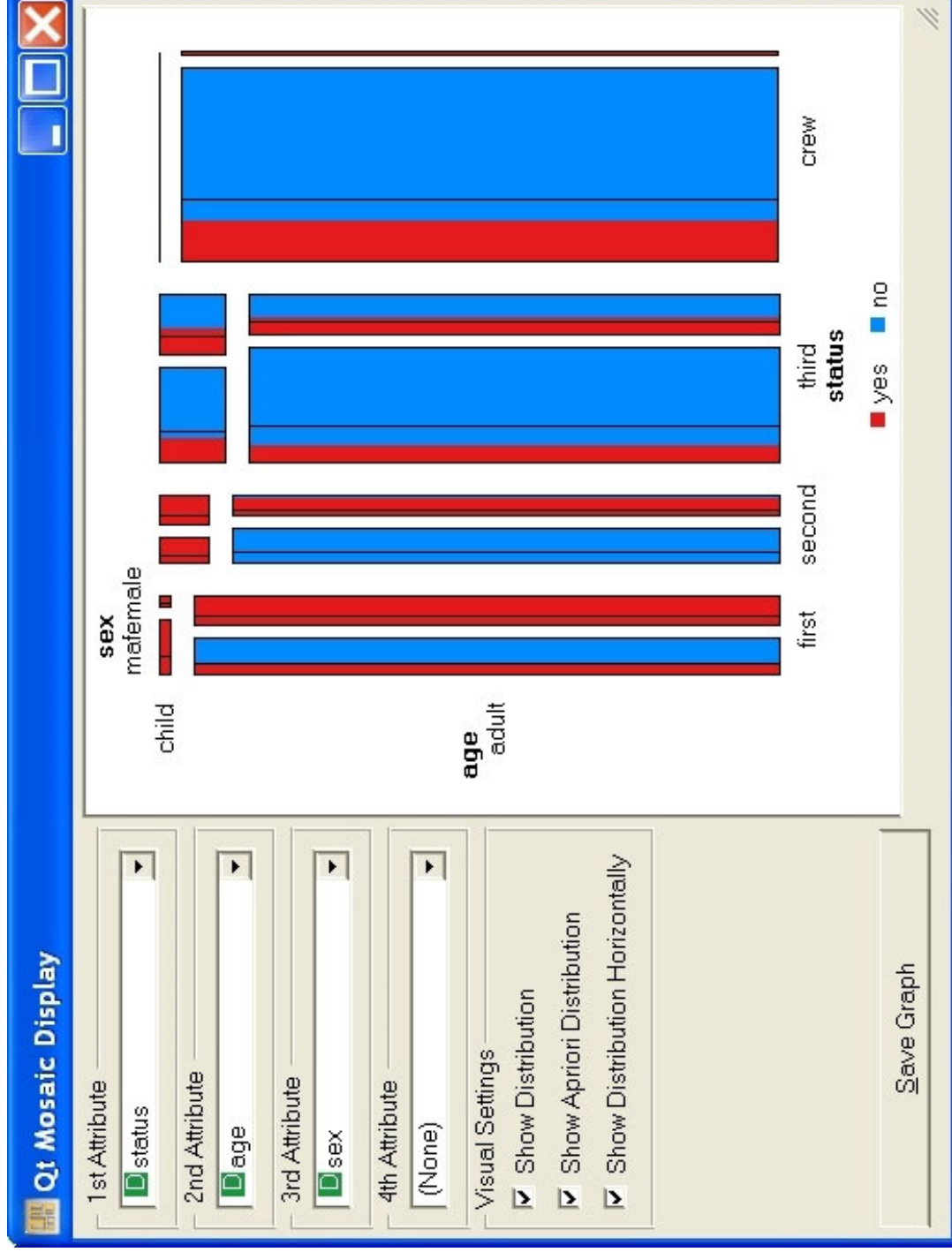# SURVEY PLOT  - WINE DATA

SURVIVED

RED = NO
BLUE = YES

FEMALE

MALE

Data on Titanic Disaster

Scatterplot and parallel coords do not immediately tell you who and who did not survive
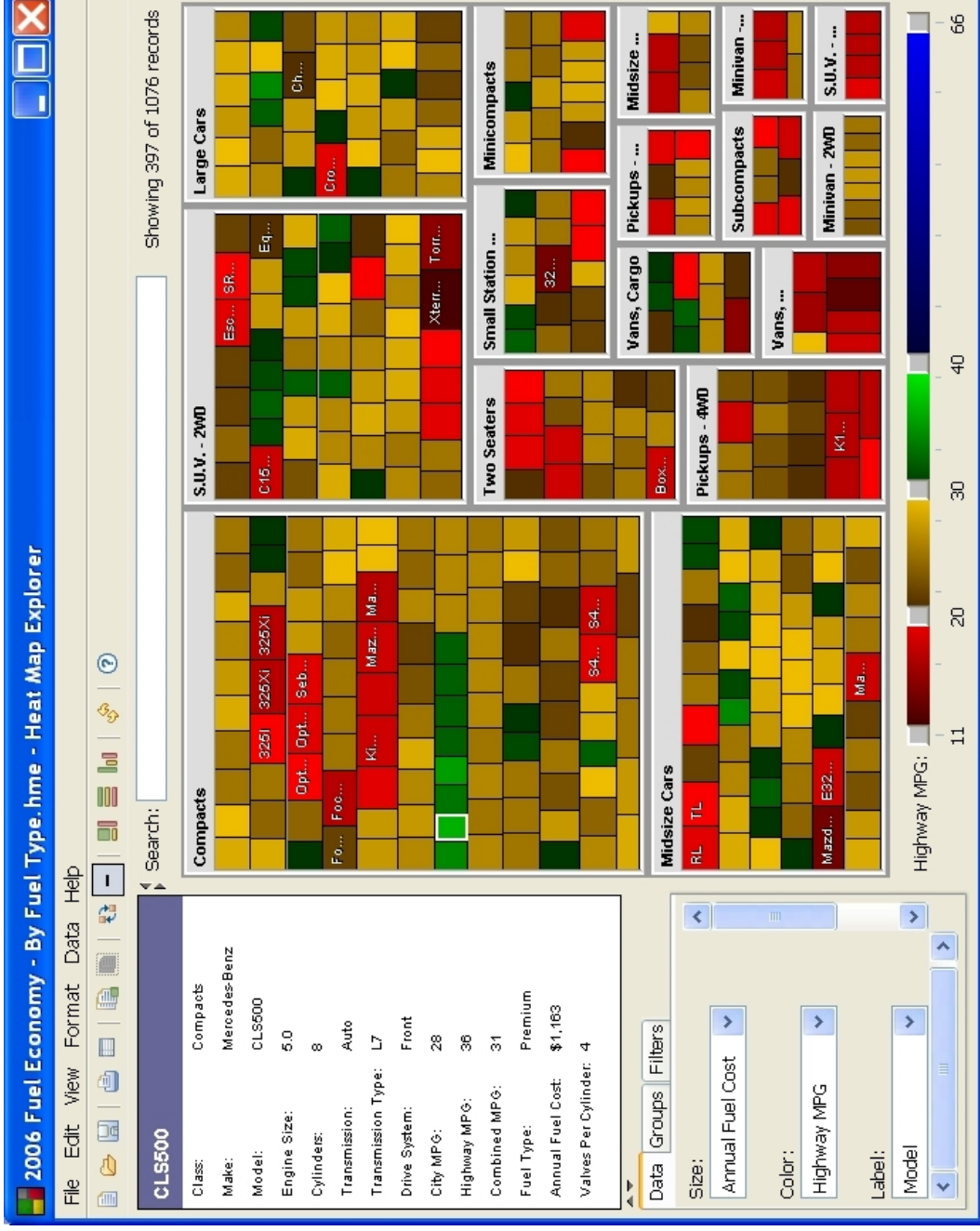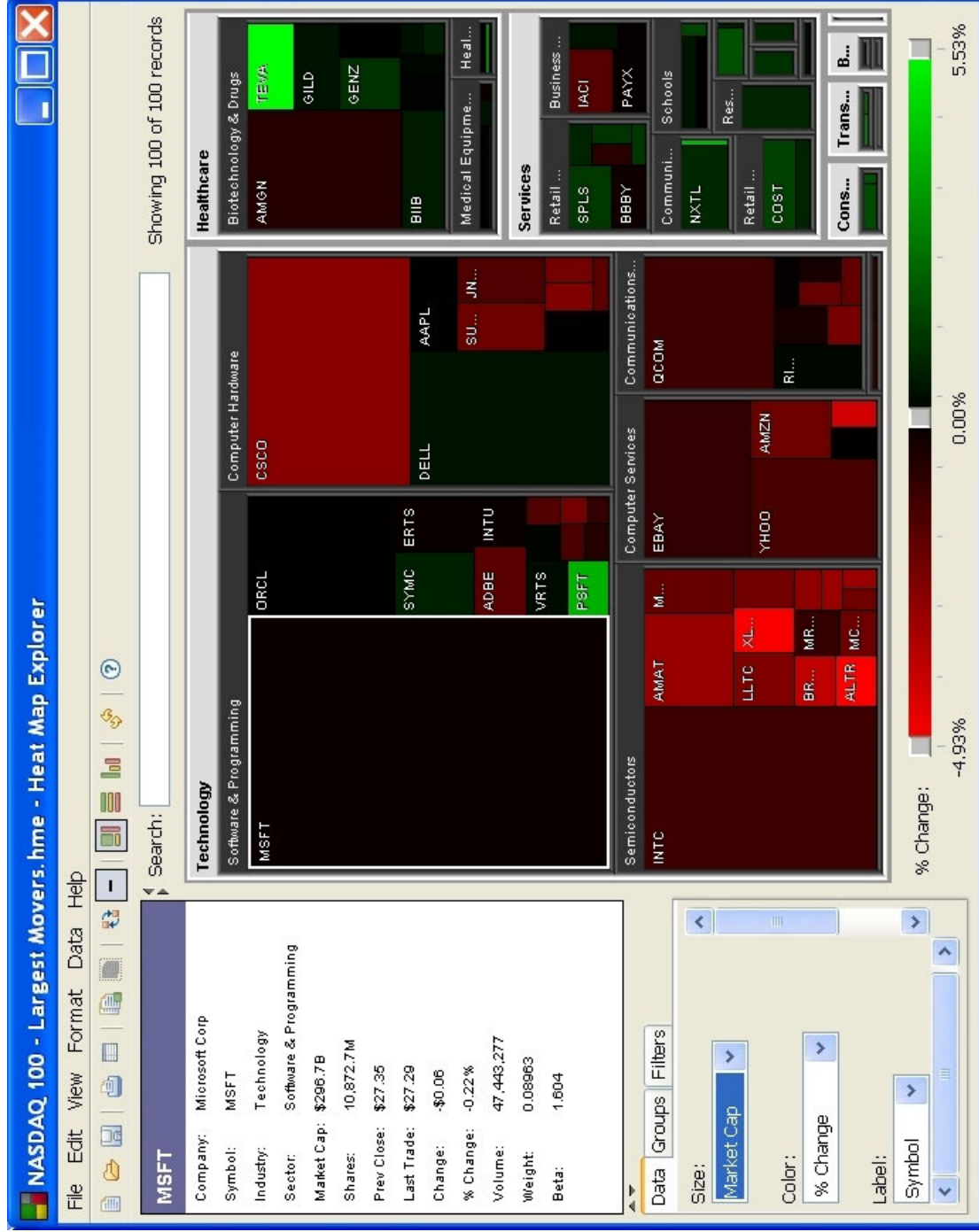
# MOSAIC DISPLAY WORKS VERY WELL WITH THIS DATA

**HEAT MAPS** - invented by USA professor who wanted to visualise the data on his hard disk

Has drill down procedure



**Excellent for looking at a large amount of data quickly – online histograms ?**

Useful for stock brokers....................

# NOW ITS YOUR TURN....

- ## WHERE TO START

- ## PRACTICAL ADVICE

| Software | Site | Comment |
|---|---|---|
| **CrystalVision** | ftp://www.galaxy.gmu.edu/pub/ | Windows. ExplorN Unix α-channel. GT, PC Needs development. |
| GGobi | www.ggobi.org | No α-channel.GT, PC All Platforms. Access to R. |
| Mondrian | http://stats.math.uni-augsburg.de/Mondrian/ | Java. α-channel. |
| Visulab | http://www.inf.ethz.ch/personal/hinterbe/Visulab/ | Excel plugin. PC only |
| **Orange** | http://www.ailab.si/orange | Component based data mining. C++ and python scripting. PC. |
| **WEKA** | http://www.cs.waikato.ac.nz/ml/weka/ | Java based data mining package Large no of algorithms included. |
| Datadesk | http://www.datadesk.com/ | Commercial. Linked plots. Stats. |
| Statistica | http://www.statsoft.com/ | Commercial. Very powerful. Not evaluated yet. Graphics + Stats. |
| VisualExplorer | www.curvaceous.com | Commercial. PC for process control Excel PlugIn. |

# When to use a particular visualisation ???

## Benchmark Development for the Evaluation of Visualization for Data Mining

Georges G. Grinstein[1], Patrick Hoffman[1], Sharon J. Laskowski[2], Ronald M. Pickett[1]

[1]Institute for Visualization and Perception Research
University of Massachusetts at Lowell, Lowell, MA 01854
{grinstein, phoffman, pickett}@cs.uml.edu

[2]The National Institute of Standards and Technology, Gaithersburg, MD 20899
sharon.laskowski@nist.gov

**Table 8 Scatter Plot Matrix**

| TASK / DATA SET | See Outliers | See Clusters | Find Class Clusters | See All Important Features | See Some Important Features | See Possible Rule/Model | See Exact Rule/Model |
|---|---|---|---|---|---|---|---|
| Balloons | | | | | | | |
| Balloons-flattened | | | | | | | |
| Lenses | | | | | | | |
| Lenses-flattened | | | | | | | |
| Orings | Y | Y | Y | Y | Y | Y | |
| Monks1-training | | | | | | | |
| Iris | Y | Y | Y | | Y | Y | |
| Congress | | | | | | | |
| Liver | Y | Y | | | | | |
| Cars | Y | Y | Y | | Y | Y | |
| Wine | Y | Y | Y | | Y | Y | |

**Table 6 Survey Plot**

| TASK / DATA SET | See Outliers | See Clusters | Find Class Clusters | See All Important Features | See Some Important Features | See Possible Rule/Model | See Exact Rule/Model |
|---|---|---|---|---|---|---|---|
| Balloons | | | | | | | |
| Balloons-flattened | | | | Y | Y | | Y |
| Lenses | | | | | Y | Y | |
| Lenses-flattened | | | | Y | Y | Y | |
| Orings | Y | | | Y | Y | Y | |
| Monks1-training | | | | Y | Y | Y | Y |
| Iris | Y | Y | Y | Y | Y | Y | |
| Congress | | | | Y | Y | Y | |
| Liver | | | | | Y | | |
| Cars | | | | Y | Y | Y | |
| Wine | | | | Y | Y | Y | |

Table 4  Parallel Coordinates

| TASK / DATA SET | See Outliers | See Clusters | Find Class Clusters | See All Important Features | See Some Important Features | See Possible Rule/Model | See Exact Rule/Model |
|---|---|---|---|---|---|---|---|
| Balloons | | | | | | | |
| Balloons-flattened | | | | | | | |
| Lenses | | | | | | | |
| Lenses-flattened | | | | | | | |
| Oranges | Y | Y | Y | Y | Y | Y | |
| Monks1-training | | Y | | Y | Y | Y | Y |
| Iris | Y | Y | Y | Y | Y | Y | |
| Congress | | | | | | | |
| Liver | Y | Y | | | | | |
| Cars | Y | Y | Y | | Y | Y | |
| Wine | Y | Y | Y | | Y | Y | |

Table 5  Radviz

| TASK / DATA SET | See Outliers | See Clusters | Find Class Clusters | See All Important Features | See Some Important Features | See Possible Rule/Model | See Exact Rule/Model |
|---|---|---|---|---|---|---|---|
| Balloons | | Y | Y | | | Y | |
| Balloons-flattened | | Y | Y | | Y | Y | |
| Lenses | | Y | Y | | Y | | |
| Lenses-flattened | | Y | Y | | Y | | |
| Oranges | Y | Y | Y | Y | Y | Y | |
| Monks1-training | | Y | | | | | |
| Iris | Y | Y | Y | | | Y | |
| Congress | Y | Y | Y | | Y | | |
| Liver | Y | Y | | | | | |
| Cars | Y | Y | Y | | Y | Y | |
| Wine | Y | Y | Y | | Y | Y | |

WITH THE RIGHT TOOLKIT

(WEKA/ORANGE/CRYSTAL VISION
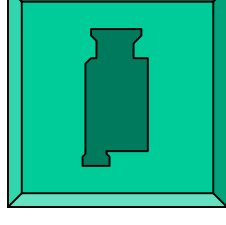ARE VERY POWERFUL STARTING
POINTS)

Visualise the data
Apply various data mining techniques to your data

Compare the results and then
use your own judgement on how to proceed

OCCAMS RAZOR is a good principle

WARNING ..…VIZ-O-Matic
The dangers of Glitziness

# CONCLUSIONS

How was this seminar researched ????

Google – data mined the Internet

Books and papers !!    The book is not dead yet.

There are powerful visualisation and machine learning techniques available. This field is rapidly growing.

USE THEM IN YOUR RESEARCH !!!!!!!

WE SHOULD BUILD THEM INTO OUR
DATA ANALYSIS SOFTWARE

STANDARD HEP REPOSITORY OF DATA SAMPLES
TO EVALUATE NEW METHODS

Can we data mine particle physics data without a priori ideas of the physics we wish to see ?????????????

FINAL THOUGHT.......

IDEAS GUIDE YOU WHEN YOU CANNOT SEE.