

# Visualizing and Exploring Data

Sargur Srihari

University at Buffalo

The State University of New York

# Visual Methods for finding structures in data

- Power of human eye/brain to detect structures
  - Product of eons of evolution
- Display data in ways that capitalize on human pattern processing abilities
- Can find unexpected relationships
  - Limitation: very large data sets

# Exploratory Data Analysis

- Explore the data without any clear ideas of what we are looking for
- EDA techniques are
  - Interactive
  - Visual
- Many graphical methods for low-dimensional data
- For higher dimensions -- Principal Components Analysis

# Topics in Visualization

## 1. Summarizing Data

Mean, Variance, Standard Deviation, Skewness

## 2. Tools for Single Variables (histogram)

## 3. Tools for Pairs of Variables (scatterplot)

## 4. Tools for Multiple Variables

## 5. Principal Components Analysis

- Reduced number of dimensions

# 1. Summarizing the data

- Mean  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x(i)$ 
  - Centrality
    - Minimizes sum of squared errors to all samples
    - If there are  $n$  data values, mean is the value such that the sum of  $n$  copies of the mean equals the sum of data values
  - Measures of Location
    - Mean is a measure of location
    - Median (value with equal no of points above/ below)
    - First Quartile (value greater than a quarter of data points)
    - Third Quartile (value greater than three quarters)
- Mode
  - Most Common Value of Data
    - Multimodal
      - 10 data points take value 3, ten value 7 all other values less often than 10

# Measures of Dispersion, or Variability

## Variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n [x(i) - \mu]^2$$

Average squared error  
in mean representing data

## Sample Variance

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n [x(i) - \hat{\mu}]^2$$

Unbiased Estimate

## Standard Deviation

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n [x(i) - \mu]^2}$$

# Skewness

Measures how much the data is one-sided  
(single long tail)

$$\frac{\sum (x(i) - \hat{\mu})^3}{\left( \sum (x(i) - \hat{\mu})^2 \right)^{3/2}}$$

Symmetric distributions have zero skewness

Distribution of people's income is skewed with large majority having low and moderate income, with few having very large income

## 2. Tools for Displaying Single Variables

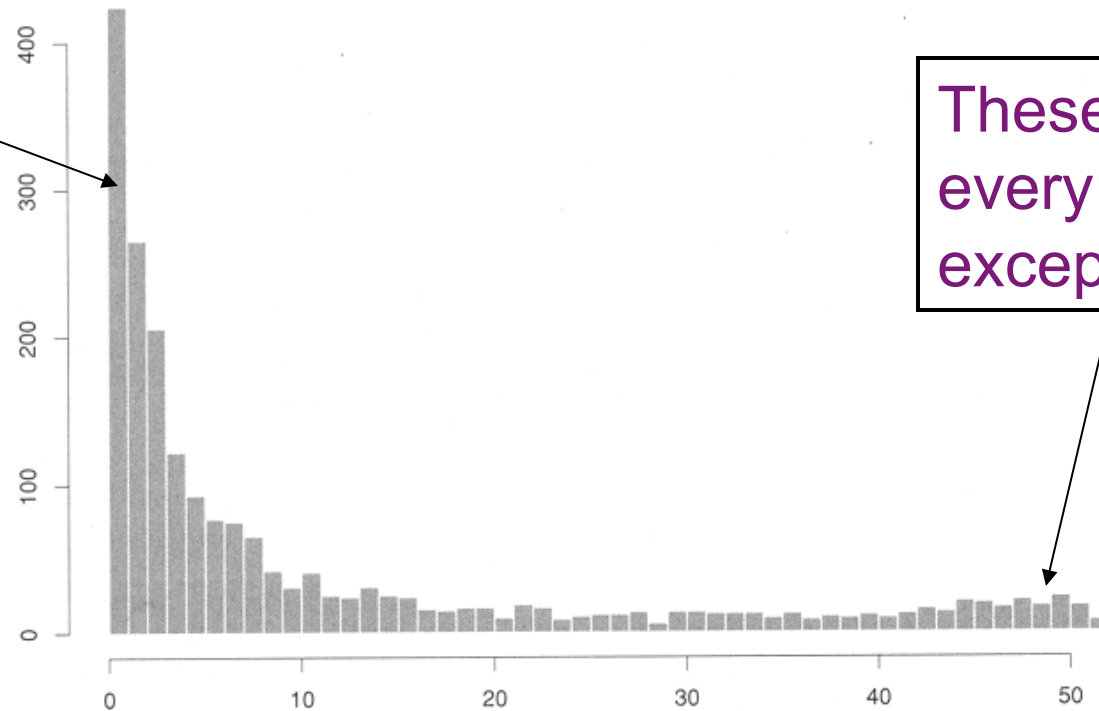
- Basic display for univariate data is the histogram
  - No of values of the variable that lie in consecutive intervals



# Histogram

(supermarket use of particular credit card)

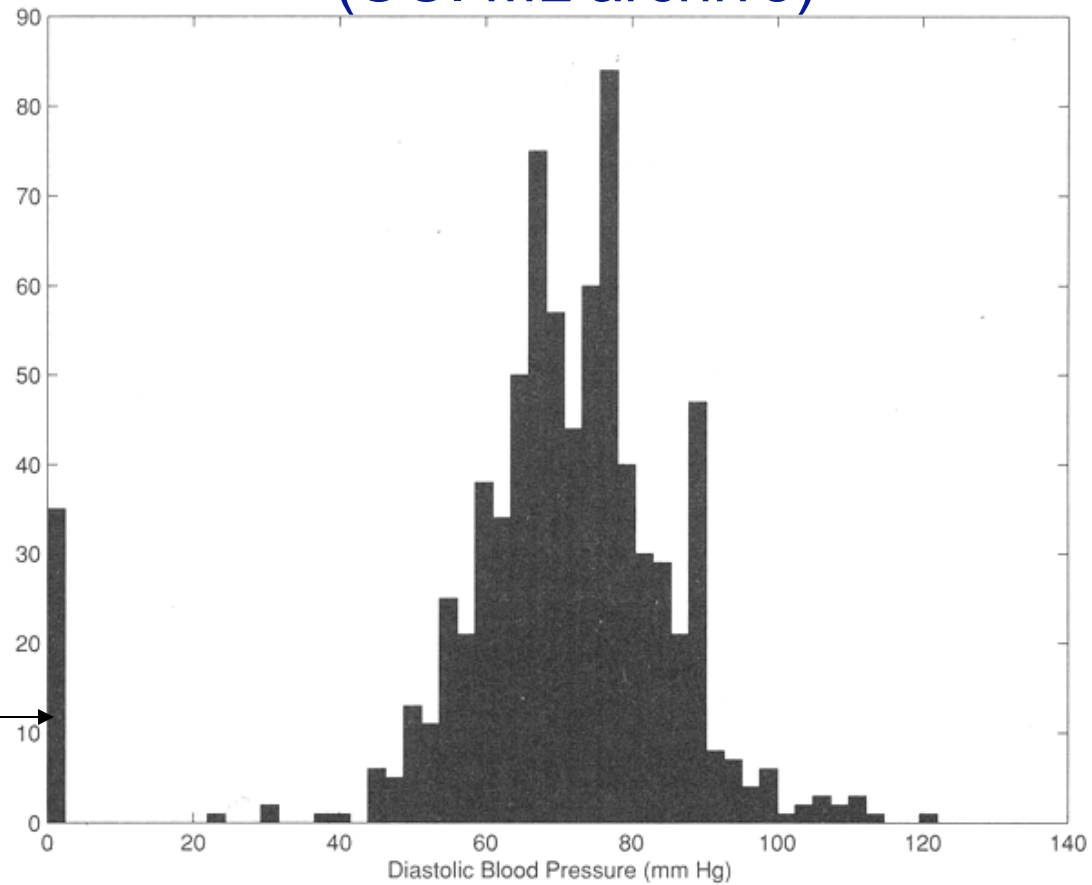
Many did not use it at all



These used it every week except holidays

Weeks (0-52)

# Histogram of Diastolic blood pressure of individuals (UCI ML archive)



Zero BP  
means  
data missing

# Disadvantages of Histograms

- Random Fluctuations in values
- Alternative choices for ends of intervals give very different diagrams
- Apparent multimodality can arise then vanish for different choices of intervals or for different small sample
- Effects diminish with increasing size of data set

# Smoothing Estimates

- Tackling disadvantages of histograms
- Kernel Function  $K$
- Estimated density at point  $x$  is

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x - x(i)}{h}\right)$$

- Gaussian Kernel with std dev  $h$

$$K(t, h) = C e^{-\frac{1}{2}\left(\frac{t}{h}\right)^2} \quad \text{where } t = x - x(i)$$

# Kernel Estimates with two values of $h$

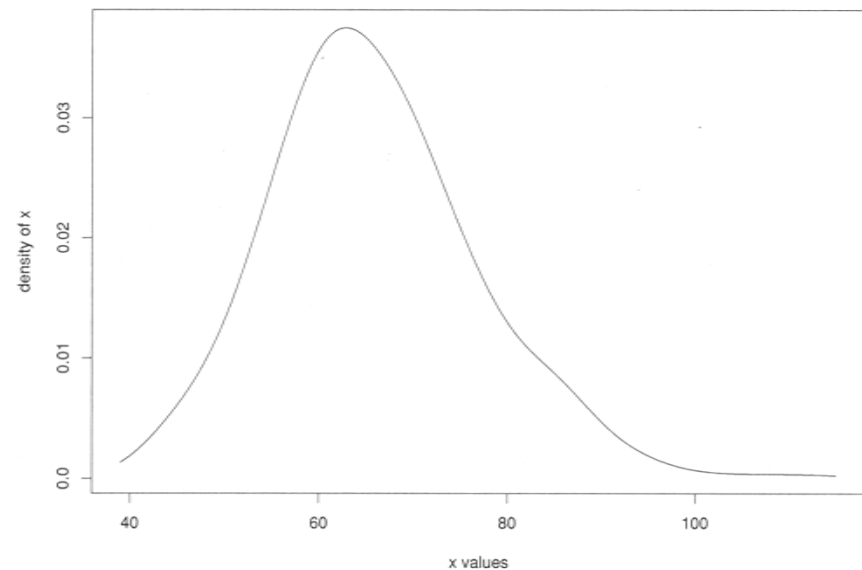
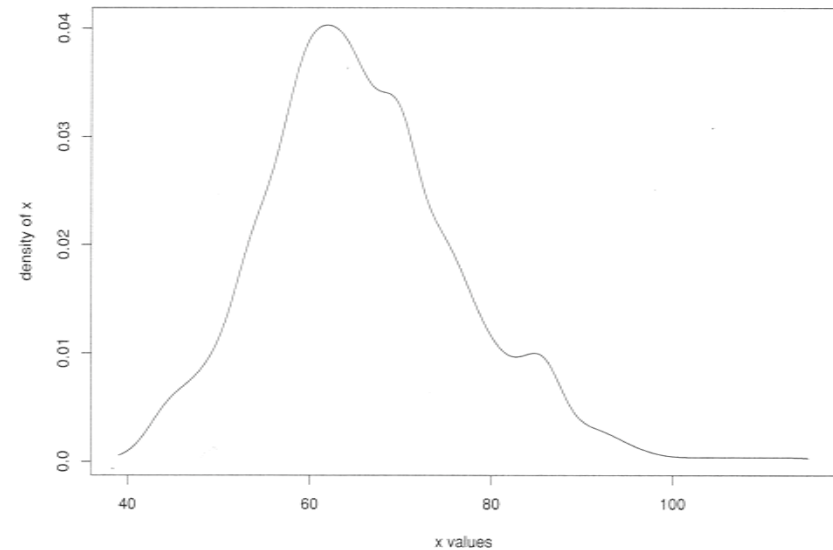
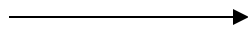
Small values lead to spiky estimates

Data is right skewed  
with hint of multimodality

$$K(t, h) = Ce^{-\frac{1}{2}\left(\frac{t}{h}\right)^2}$$

Higher  $h$

More smoothing



Srihari

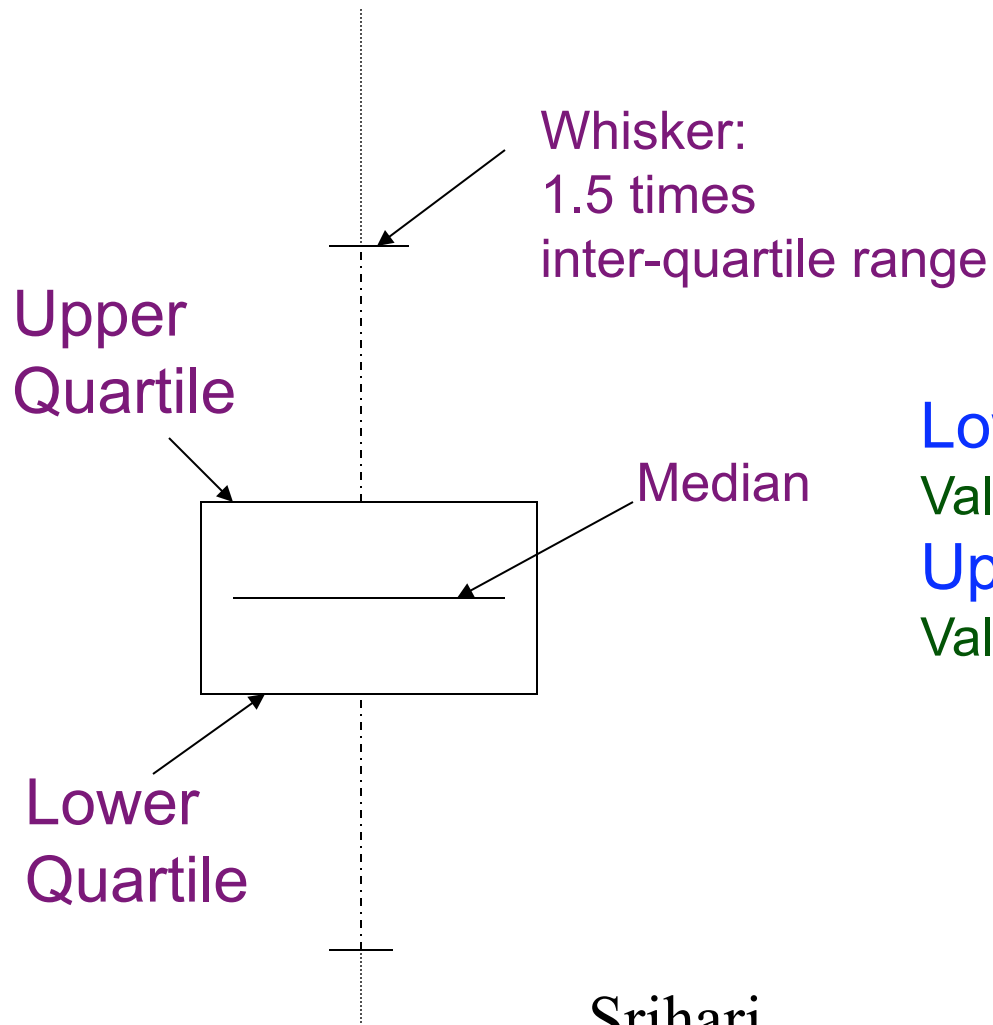
# 3. Tools for Displaying Relationship between two variables

- Box Plots
- Scatter Plots
- Contour Plots
- Time as one of the two variables

# Box Plot

Box contains bulk of data

E.g., interval between first and third quartiles



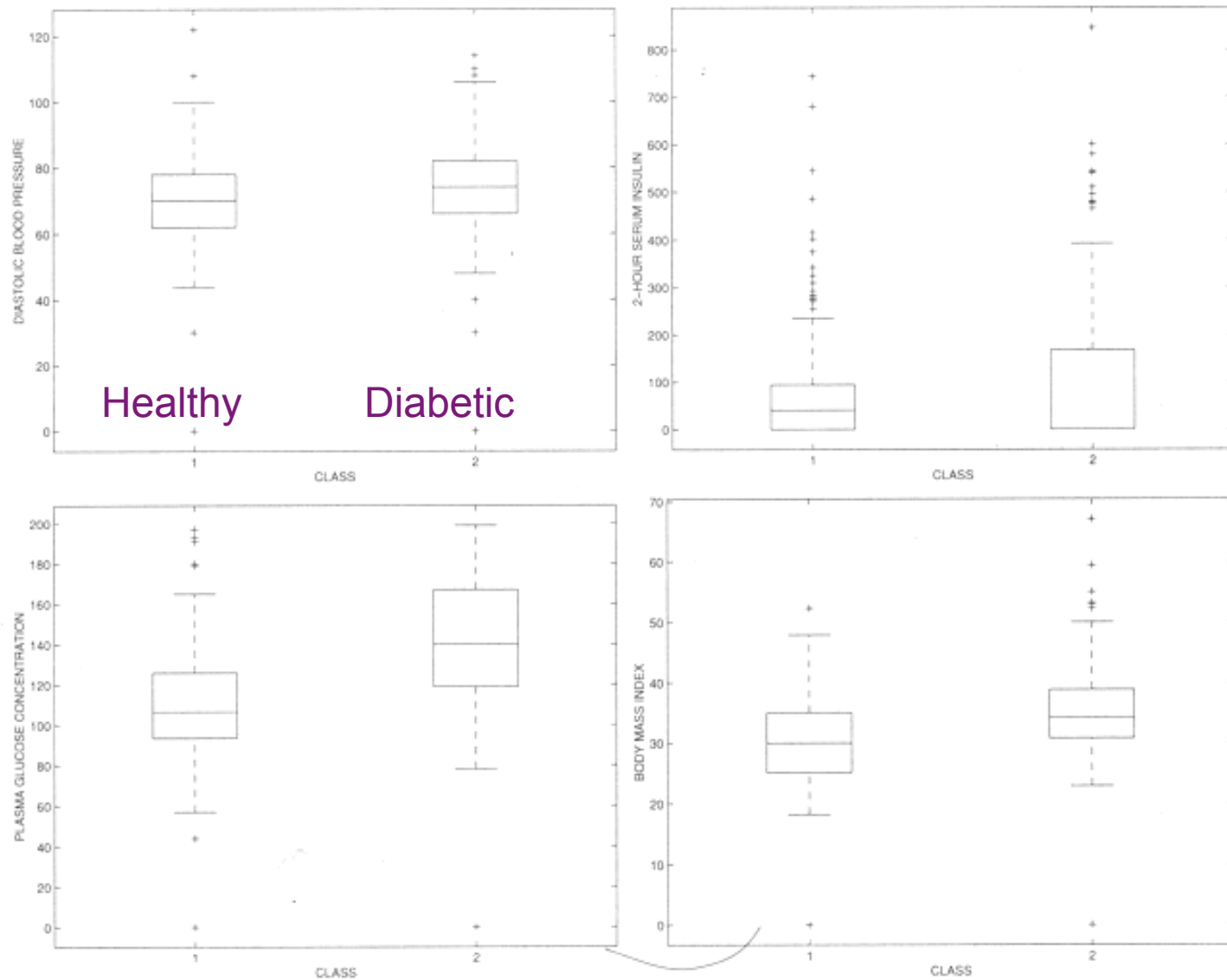
Lower Quartile:

Value greater than quarter of points

Upper Quartile:

Value less than a quarter of points

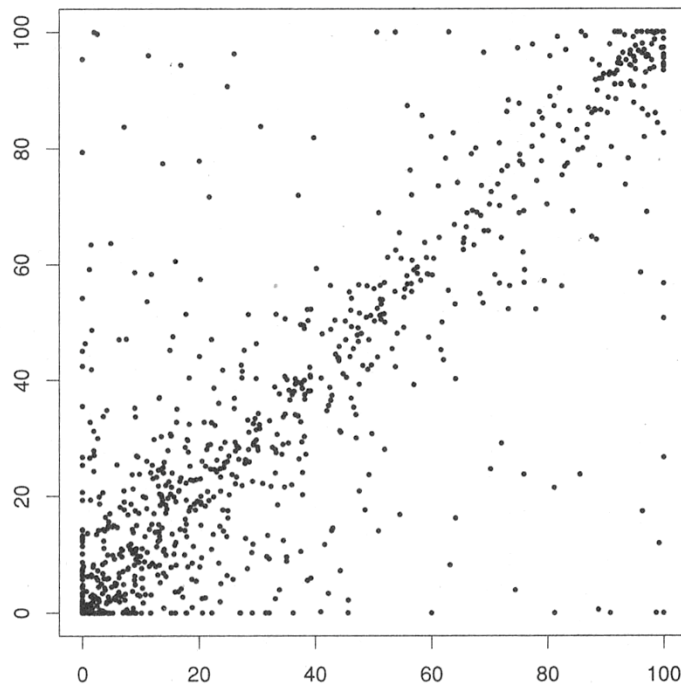
# Box Plots with Multiple Variables





# Scatterplot

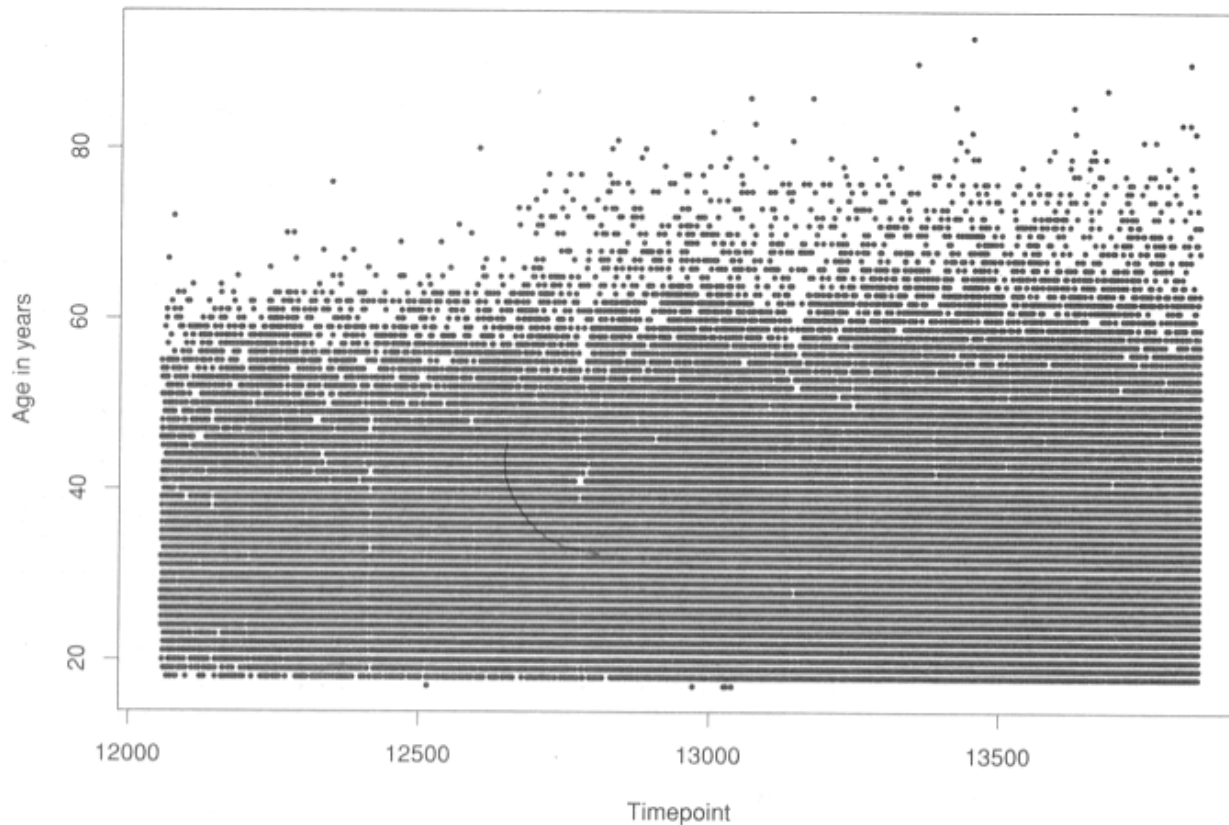
Credit card repayment data (Two banking variables)



Highly correlated data  
Significant number depart  
from pattern: worth investigating

# Scatterplot Disadvantages

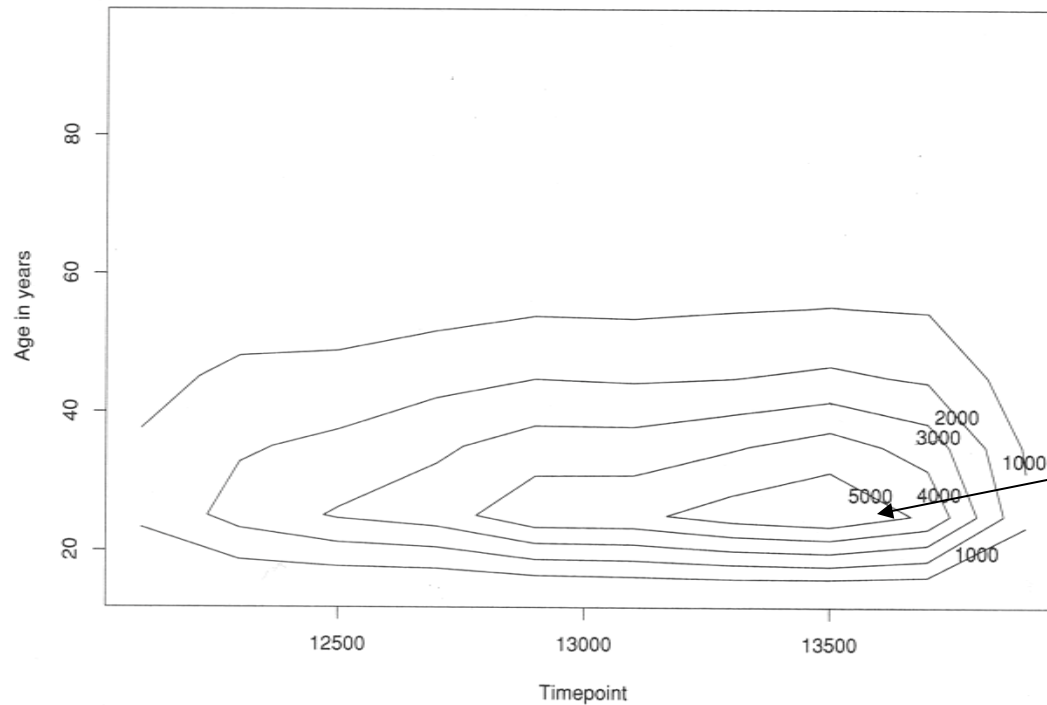
1. With large no of data points reveals little structure



2. Can conceal overprinting which can be significant for multimodal data

# Contourplot

## 1. Overcomes some scatterplot problems



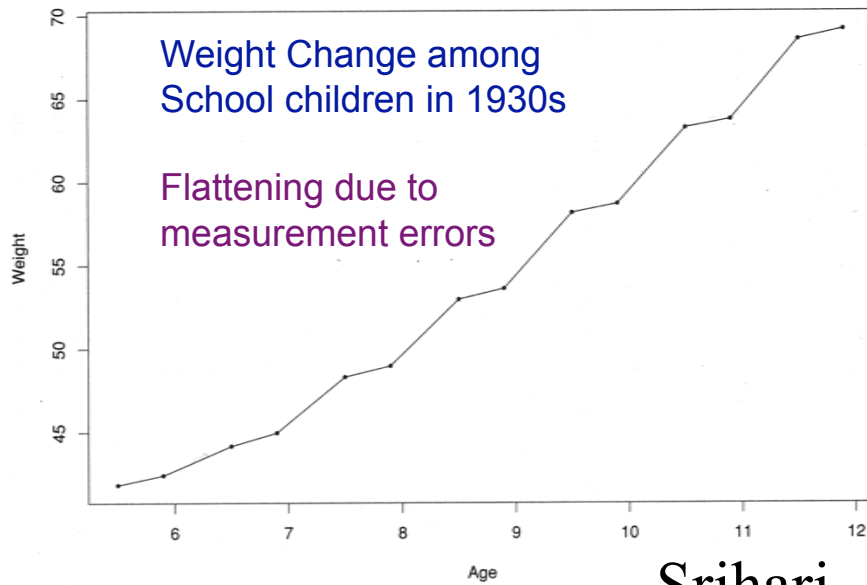
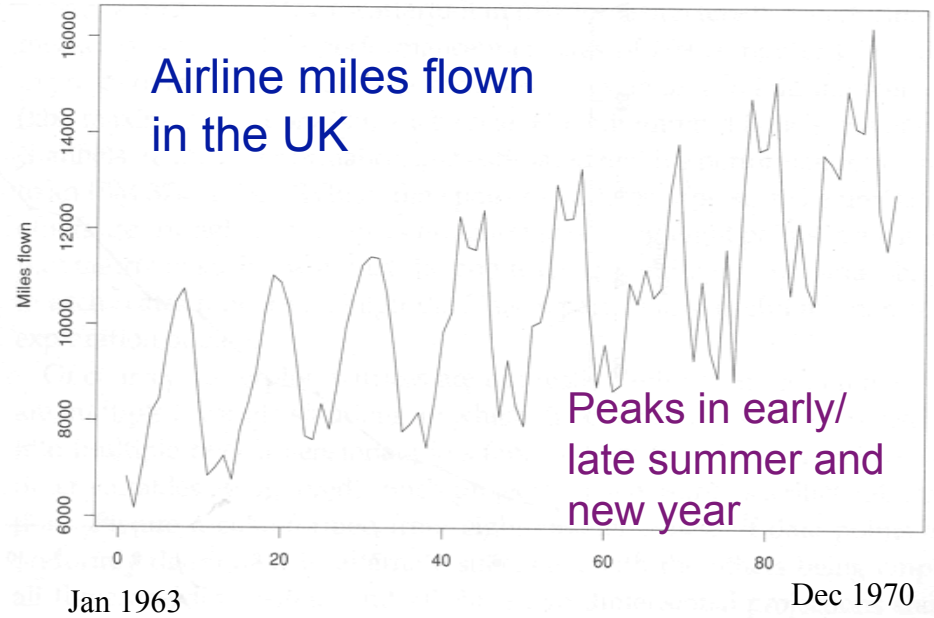
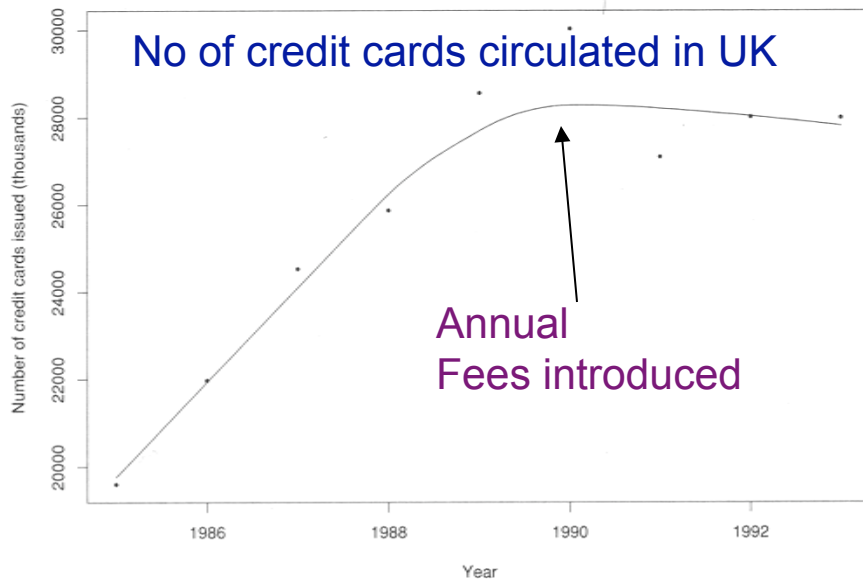
Unimodality  
can be seen:  
Not apparent  
in scatterplot

Same Data as previous

## 2. Requires a 2-D density estimate to be constructed with a 2-D kernel

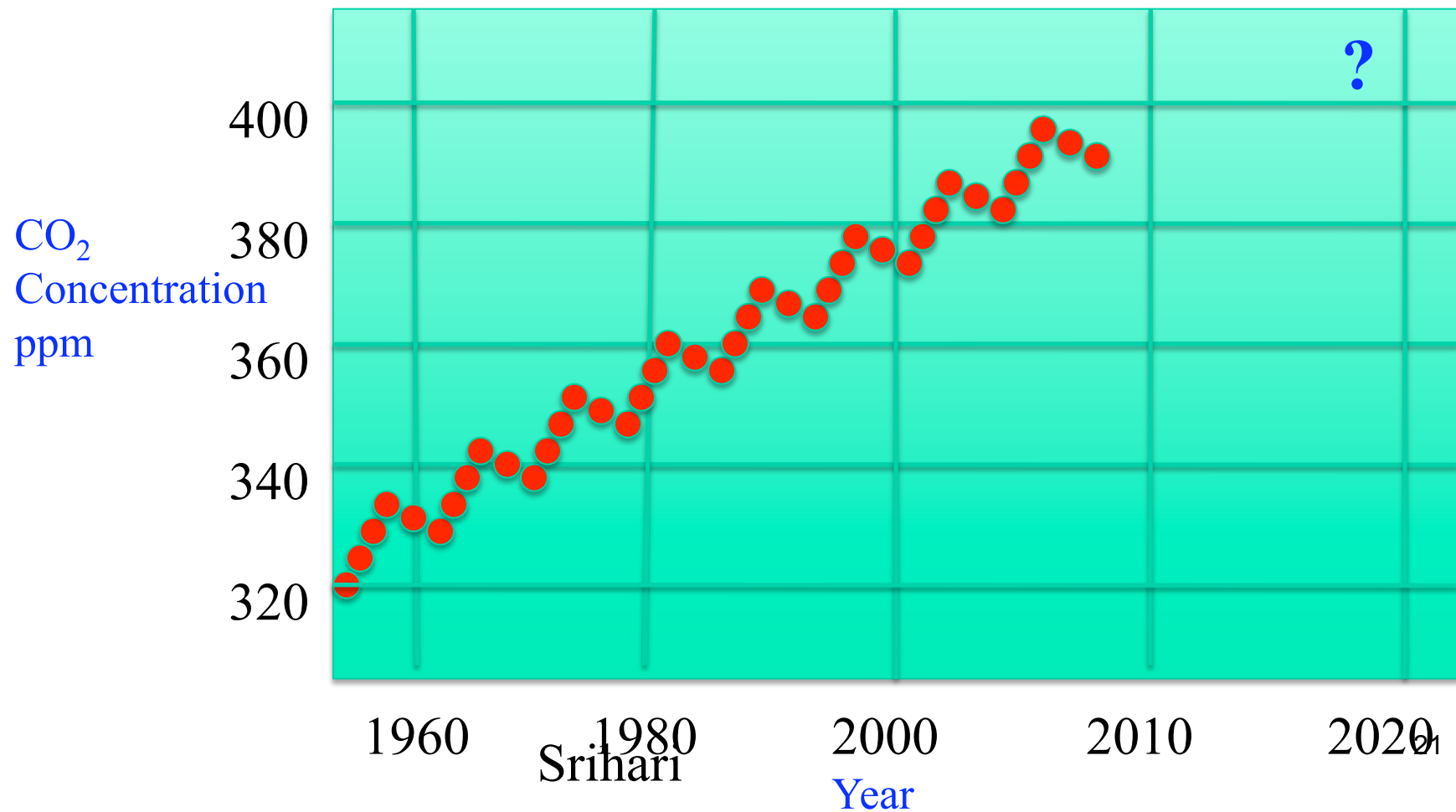
Srihari

# Display when one of the variables is time



Srihari

# Carbon Dioxide in Atmosphere



# Tools for Displaying More than Two Variables

- Scatter plots for all pairs of variables
- Trellis Plot
- Parallel Coordinates Plot

# More than two variables

- Sheets of Paper and Computer screens are fine for two variables
- Need projections from higher-dimensional data to 2-D plane
- Methods
  - Examine all pairs of variables
    - Scatterplot matrix
    - Trellis plot
    - Icons

# Scatter Plot Matrix

## CPU performance

209 CPU data:

Cycle Time

Minimum Memory

Maximum Memory

Cache Size (Kb)

Minimum Channels

Maximum Channels

Relative Performance

Estimated rel perf (wrt IBM)





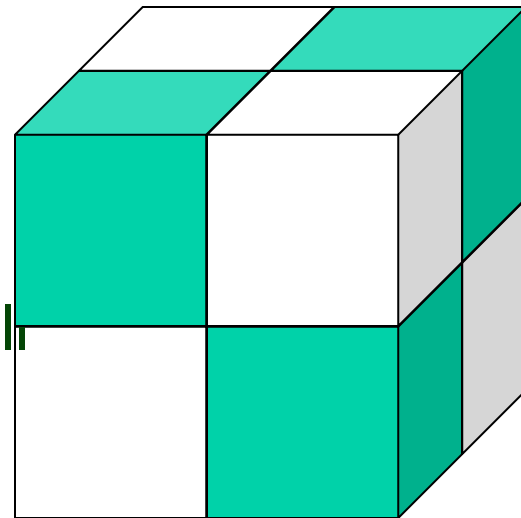
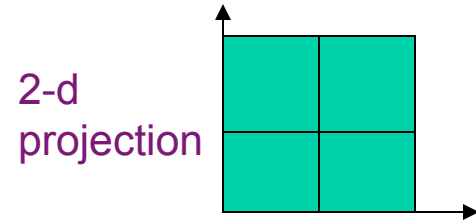
# Disadvantage of Scatter Plot Matrices

- Scatter Plot Matrices are multiple bivariate solutions
- Not a multivariate solution
- Such projections sacrifice information

3 variables

8 cubes: alternately empty and full

Each 1-D and 2-D projection is uniformly distributed!



# Trellis Plot

- Rather than displaying scatter plot for each pair of variables
- Fix a particular pair of variables and produce a series of scatter plots, histograms, time series plots, contour plots etc

# Trellis Plot

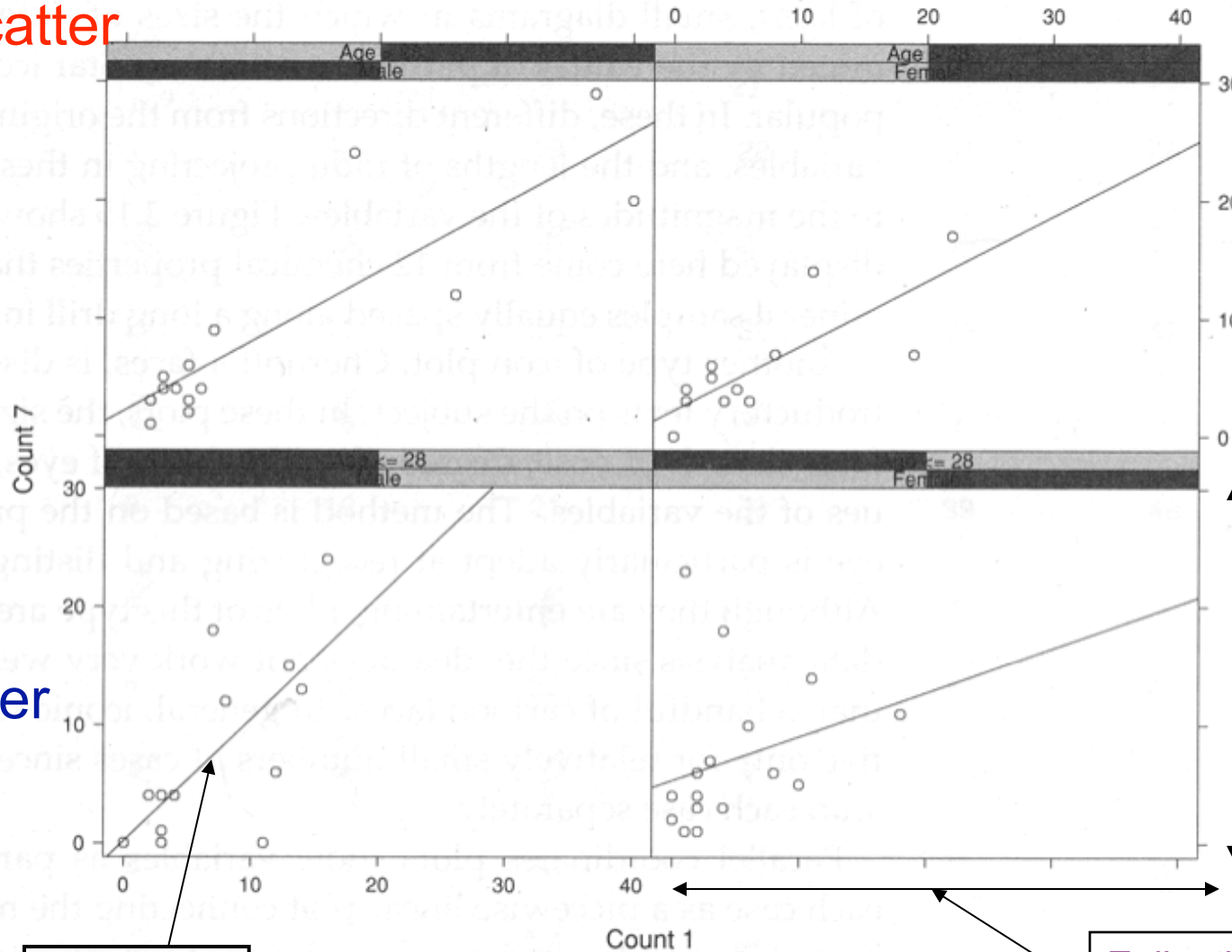
(with scatter plots)

Male

Female

Older

Younger



Best fit line

Epileptic Seizures in later 2 week period

Epileptic Seizures in 2 week period

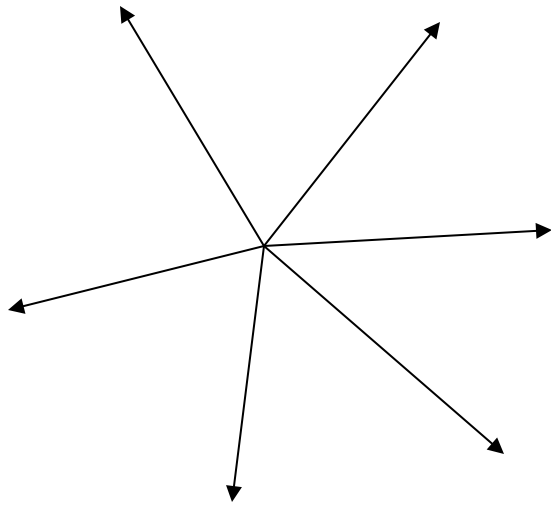
Srihari

# Icon Plot

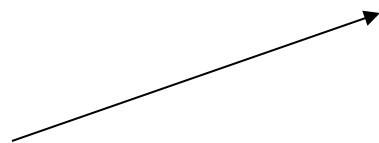
## Star Plot:

Each direction corresponds to a variable.

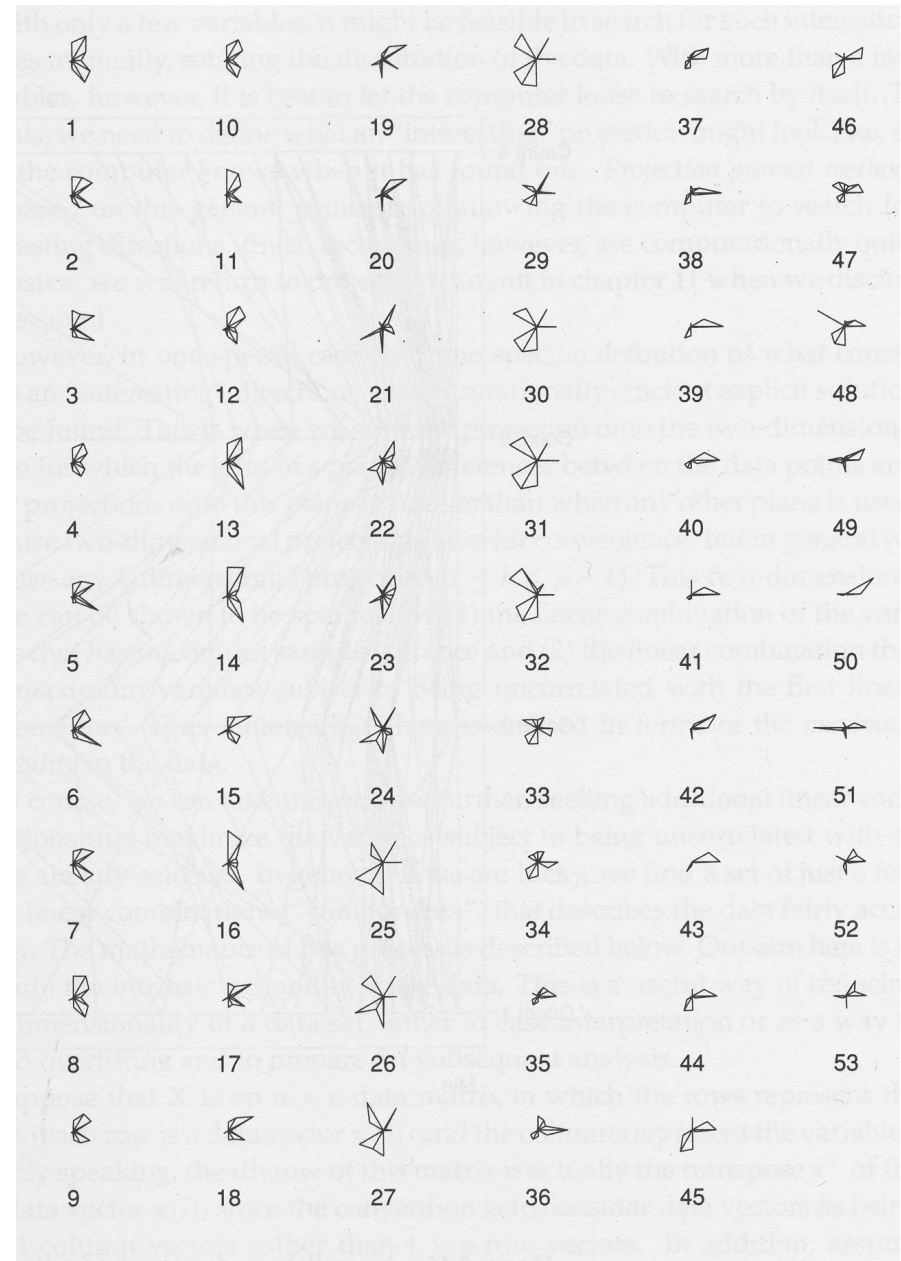
Length corresponds to a value



53 samples of minerals  
12 chemical properties



Srihari



# Parallel Coordinates Plot

Each path represents  
an individual

Each count  
Represents 2-week  
period

