

Visualizing and Exploring Data

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

$$\text{Mean, } \mu = \frac{1}{n} \sum_{i=1}^n x(i)$$

- **Centrality**
 - Minimizes the 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 that has equal no of points above and below)
 - Quartile (value greater than a quarter of the data points)

Measures of Dispersion, or Variability

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

Average squared error
in mean representing data

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

$$\text{Skewness} = \frac{\sum (x(i) - \hat{\mu})^3}{(\sum (x(i) - \hat{\mu})^2)^{3/2}}$$

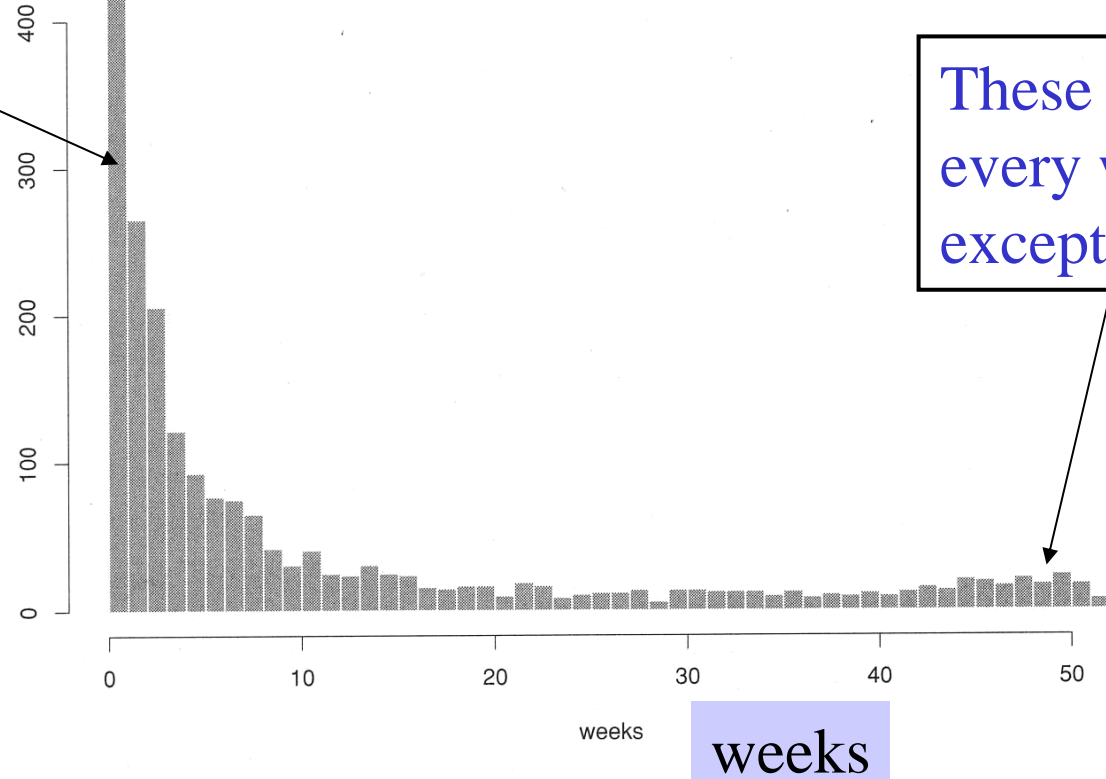
Measures how much the data
is one-sided (single long tail)

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

Many did not use it at all

Histogram of supermarket credit card usage



These used it every week except holidays

Figure 3.1 Histogram of the number of weeks of the year a particular brand of credit card was used.

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

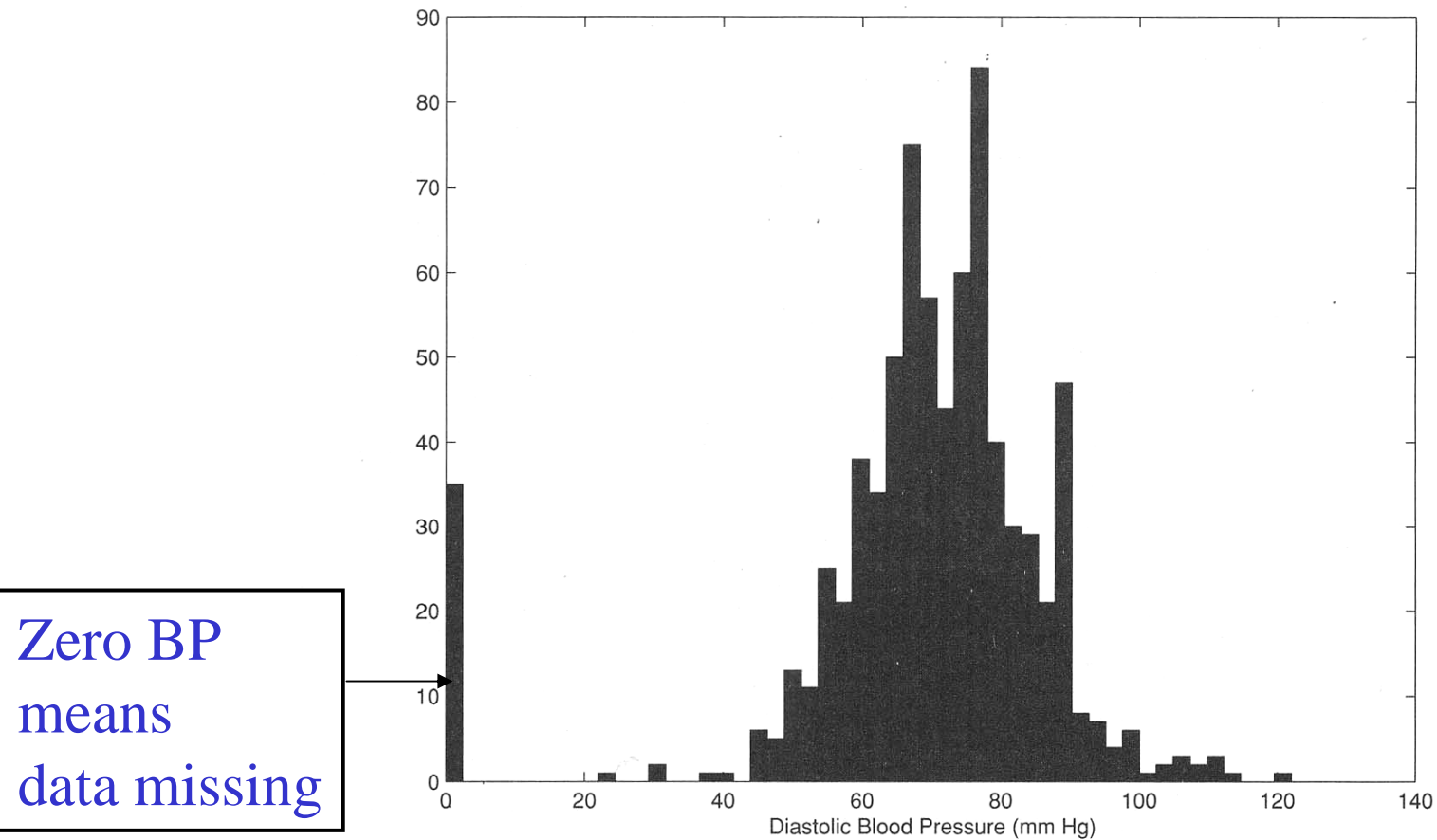


Figure 3.2 Histogram of diastolic blood pressure for 768 females of Pima Indian descent.

Smoothing estimates

- 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 different 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 smoothing \longrightarrow

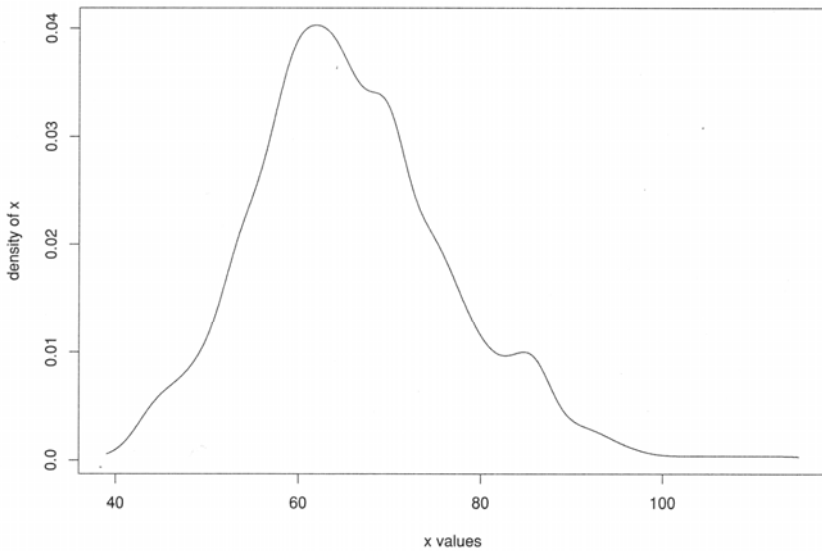


Figure 3.3 Kernel estimate of the weights (in kg) of 856 elderly women.

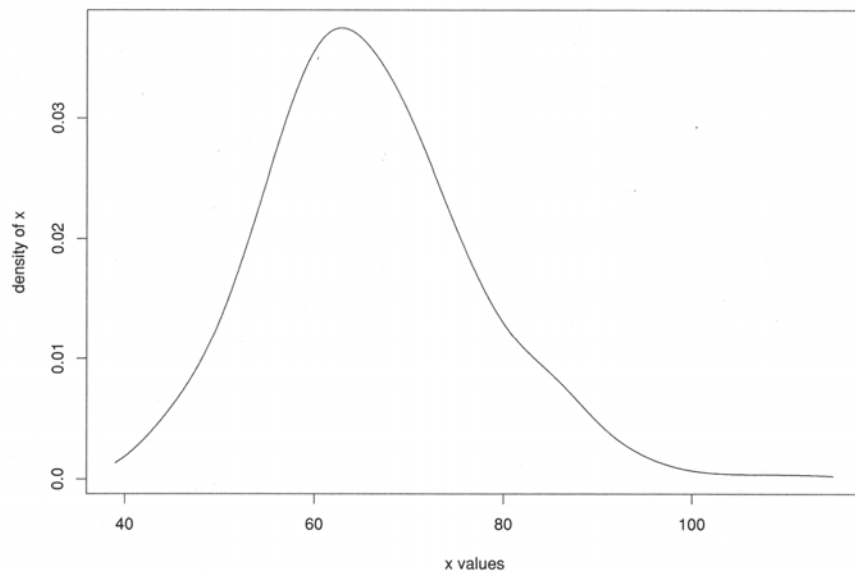
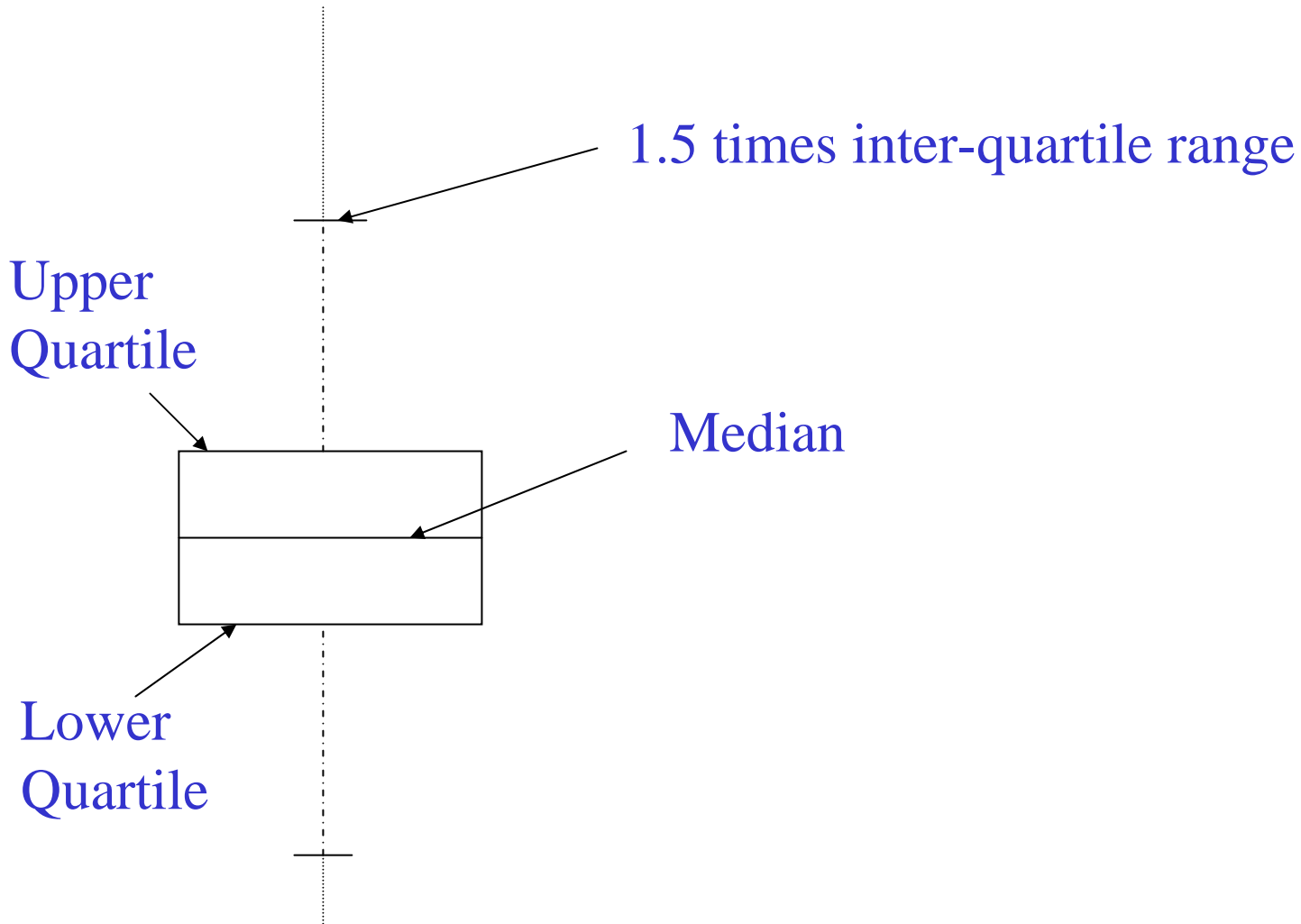


Figure 3.4 As figure 3.3, but with more smoothing.

3. Tools for Displaying Relationship between two variables

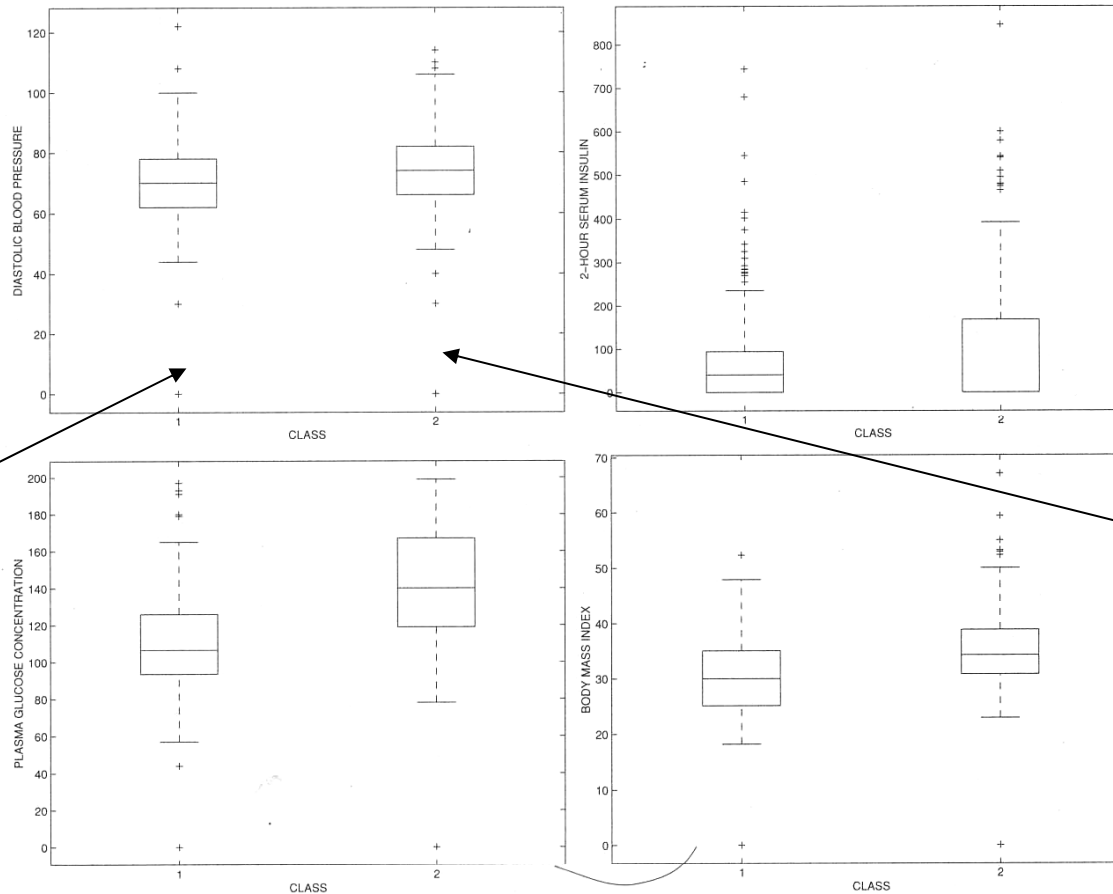
- Box Plots
- Scatter Plots
- Contour Plots
- Time as one of the two variables

Box Plot



Multiple Variables

Healthy

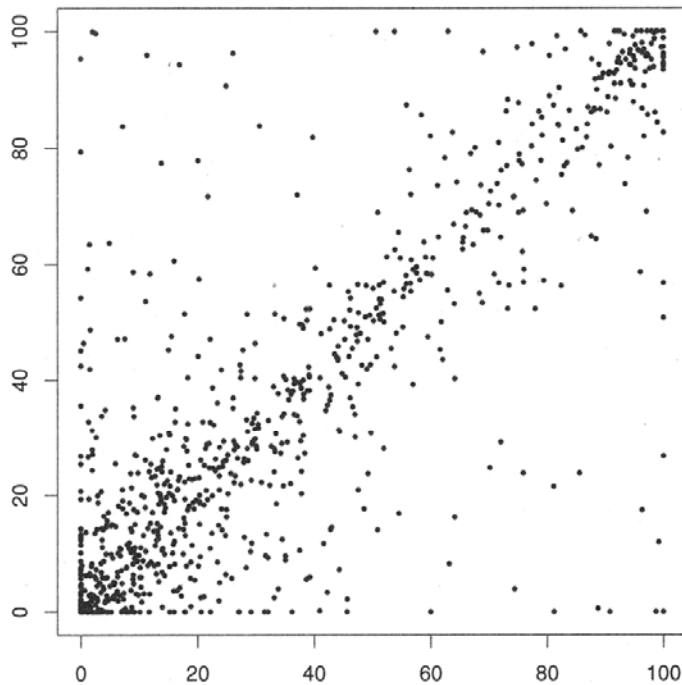


Diabetic

Figure 3.5 Boxplots on four different variables from the Pima Indians diabetes data set. For each variable, a separate boxplot is produced for the healthy subjects (labeled 1) and the diabetic subjects (labeled 2). The upper and lower boundaries of each box represent the upper and lower quartiles of the data respectively. The horizontal line within each box represents the median of the data. The whiskers extend 1.5 times the interquartile range from the end of each box. All data points outside the whiskers are plotted individually (although some overplotting is present, e.g., for values of 0).

Scatterplot

Credit card repayment data



Highly correlated data
Significant number depart
from pattern: worth investigating

Figure 3.6 A standard scatterplot for two banking variables.

Scatterplot Disadvantages

1. With large no of data points reveals little structure

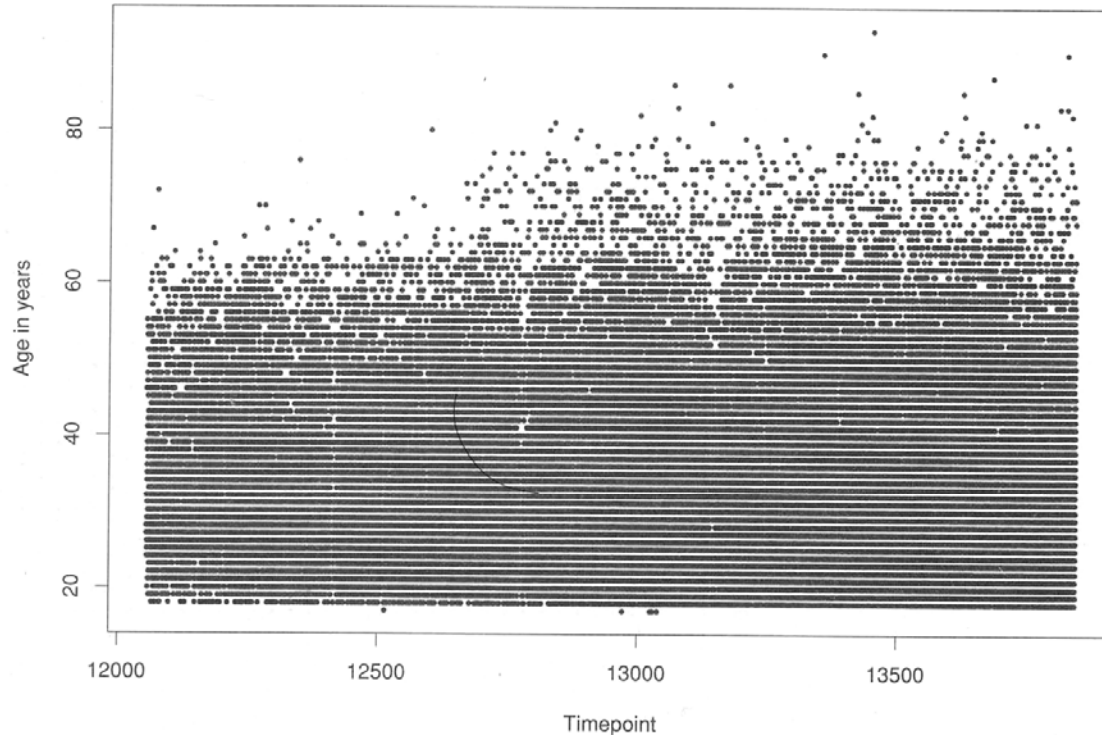
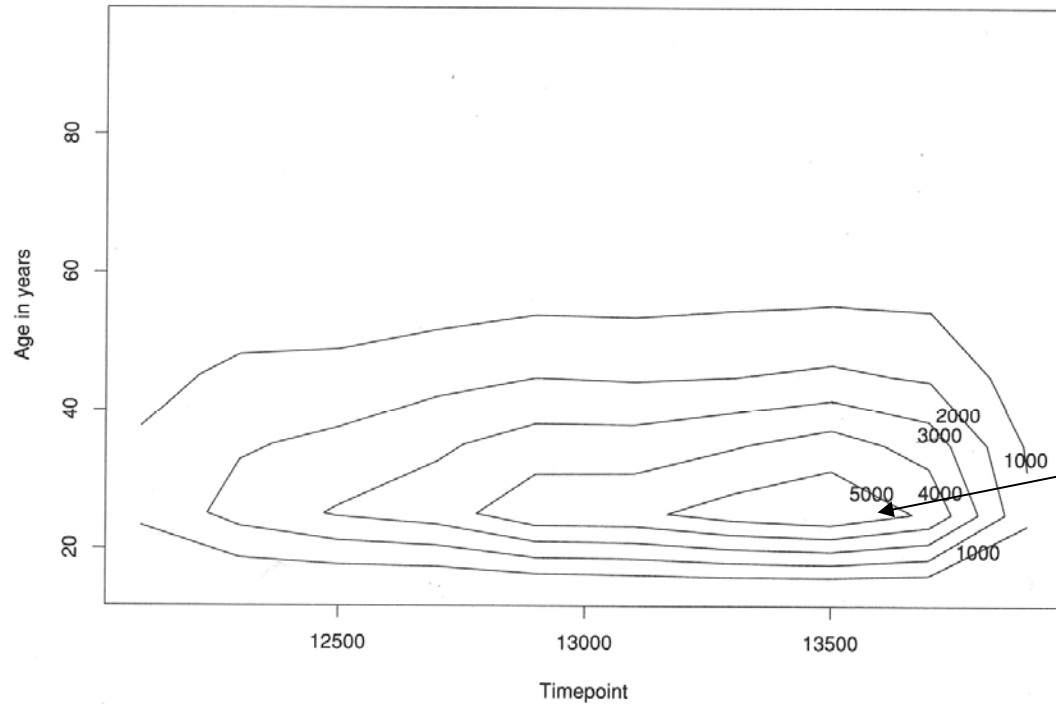


Figure 3.7 A scatterplot of 96,000 cases, with much overprinting. Each data point represents an individual applicant for a loan. The vertical axis shows the age of the applicant, and the horizontal axis indicates the day on which the application was made.

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

Figure 3.9 A contour plot of the data from figure 3.7.

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

Display when one of the variables is time

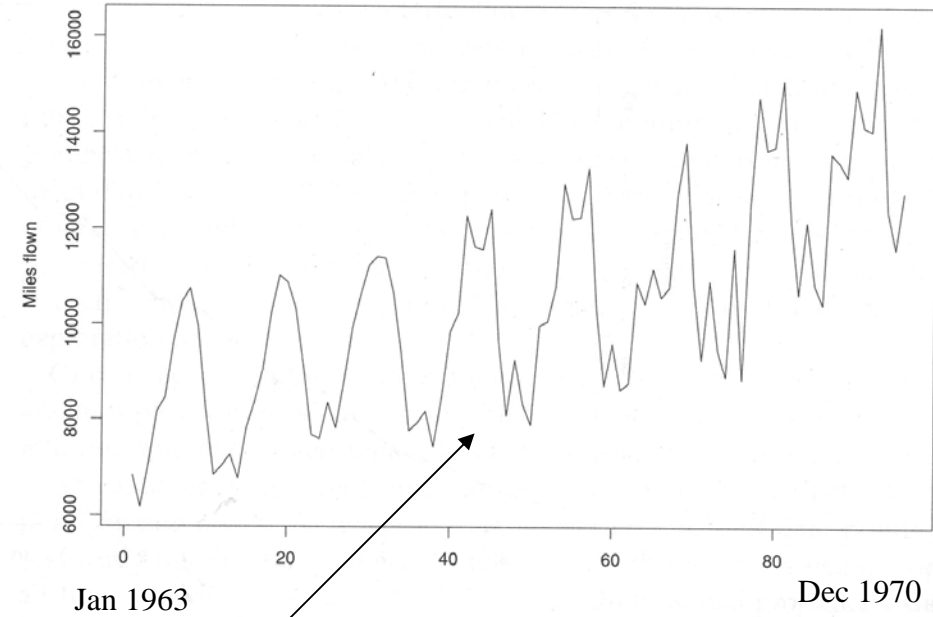
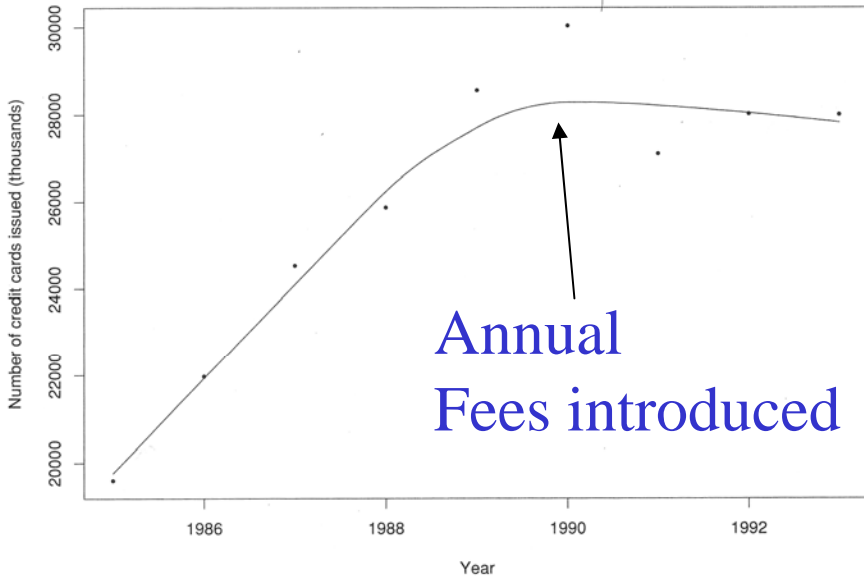


Figure 3.11 Patterns of change over time in the number of miles flown by UK airlines in the 1960s.

Peaks in early and late summer and around new year

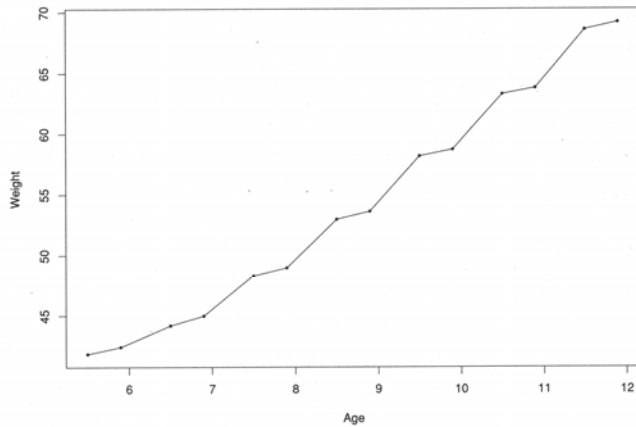


Figure 3.12 Weight changes over time in a group of 10,000 school children in the 1930s. The steplike pattern in the data highlights a problem with the measurement process.

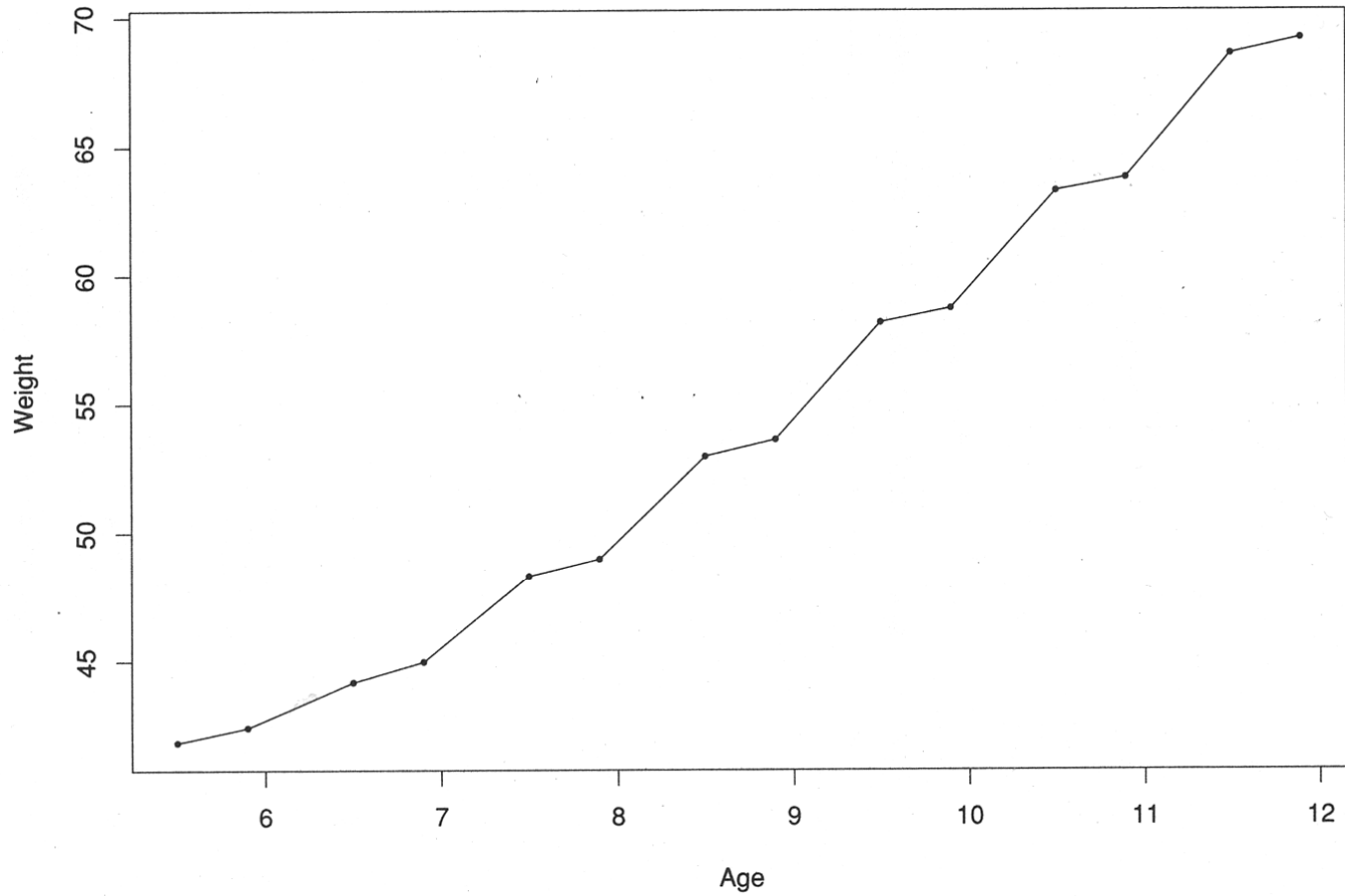


Figure 3.12 Weight changes over time in a group of 10,000 school children in the 1930s. The steplike pattern in the data highlights a problem with the measurement process.

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

CPU performance Scatter Plot Matrix

209 CPU data:

Cycle Time

Minimum Memory

Maximum Memory

Cache Size (Kb)

Minimum Channels

Maximum Channels

Relative Performance

Estimated rel perf (wrt IBM)

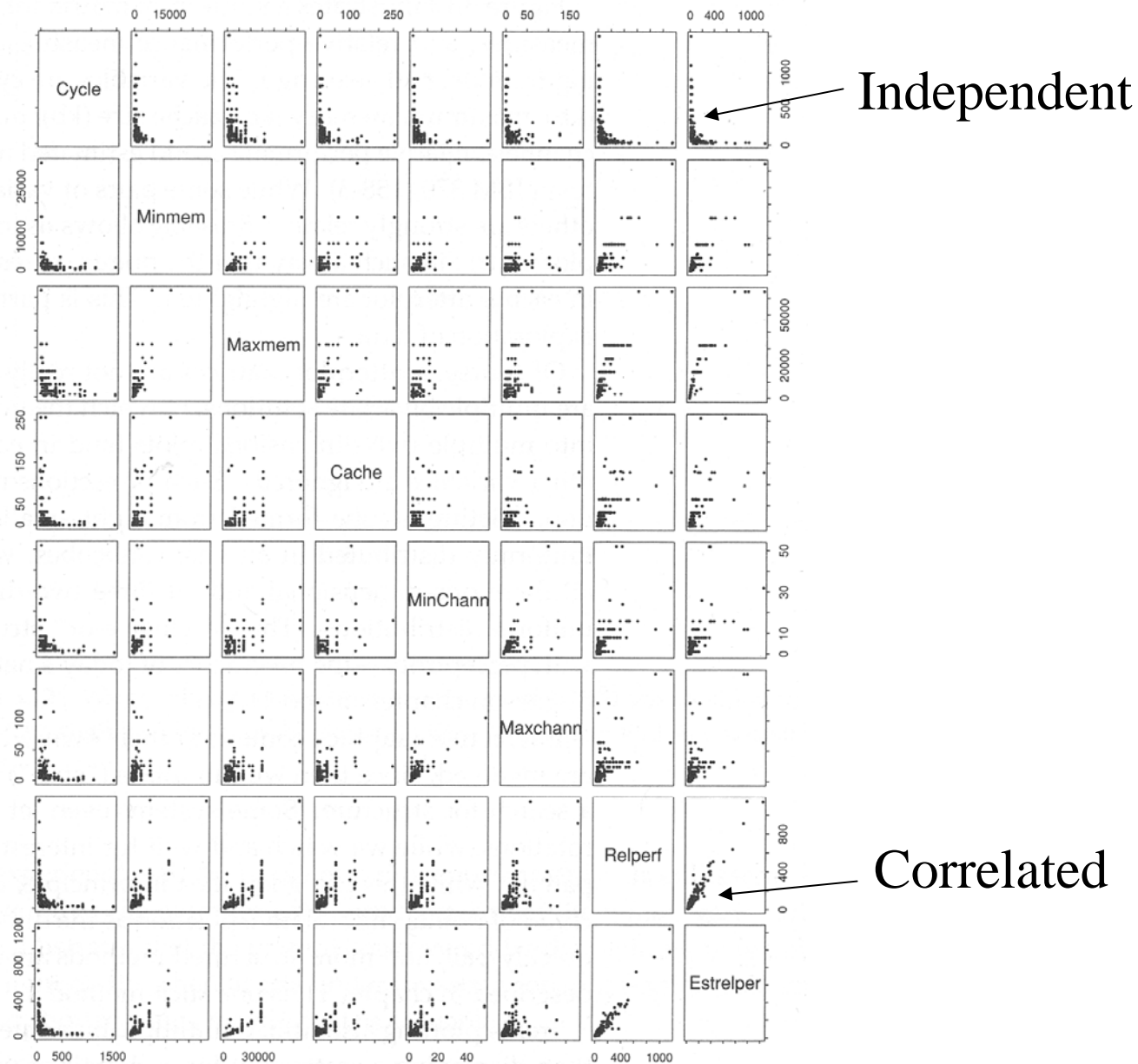


Figure 3.13 A scatterplot matrix for the computer CPU data.

Disadvantage of Scatter Plot Matrices

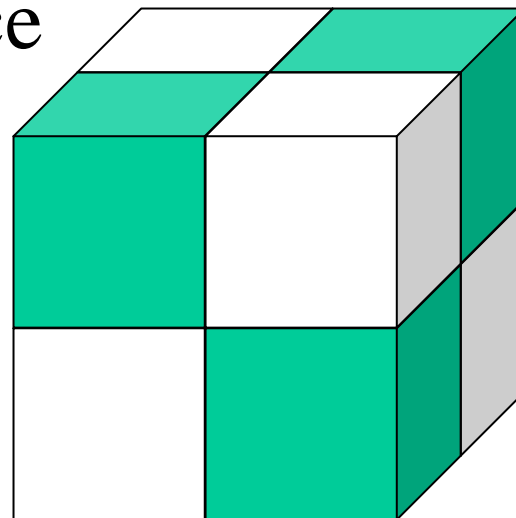
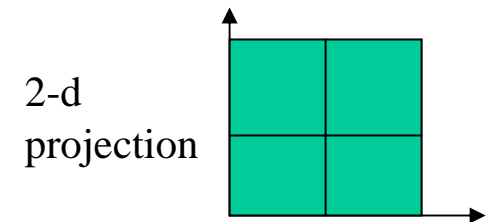
Matrices

- 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)

Older

Younger

Male

Female

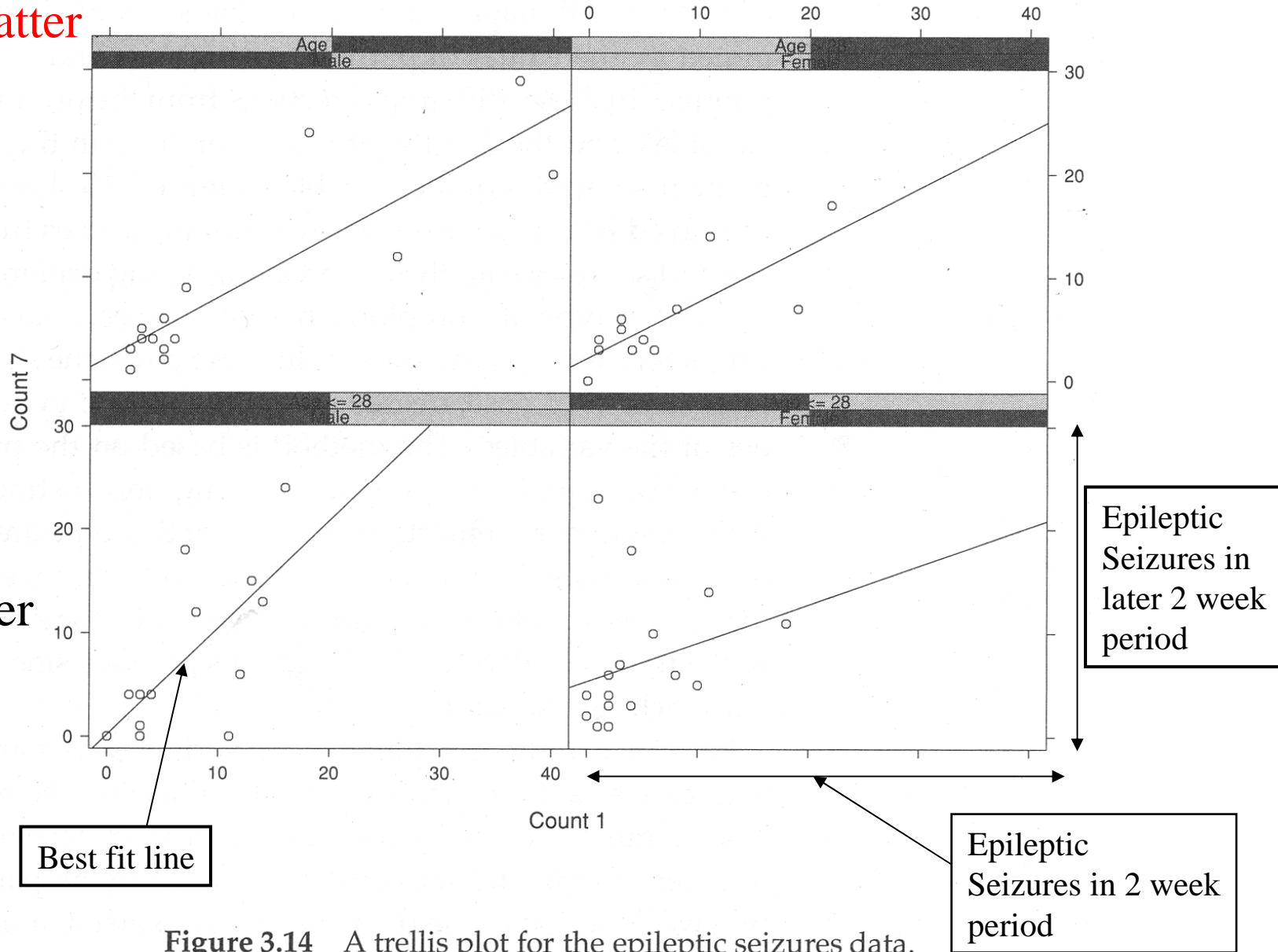
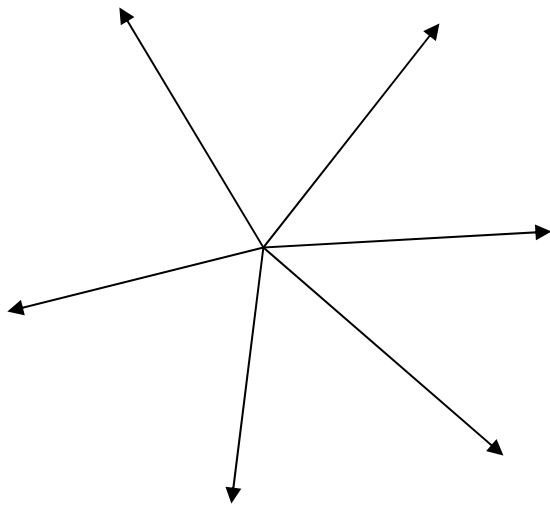


Figure 3.14 A trellis plot for the epileptic seizures data.

Icon Plot

Star Plot: Each direction corresponds to a variable. Length corresponds to a value



53 samples of minerals
12 chemical properties

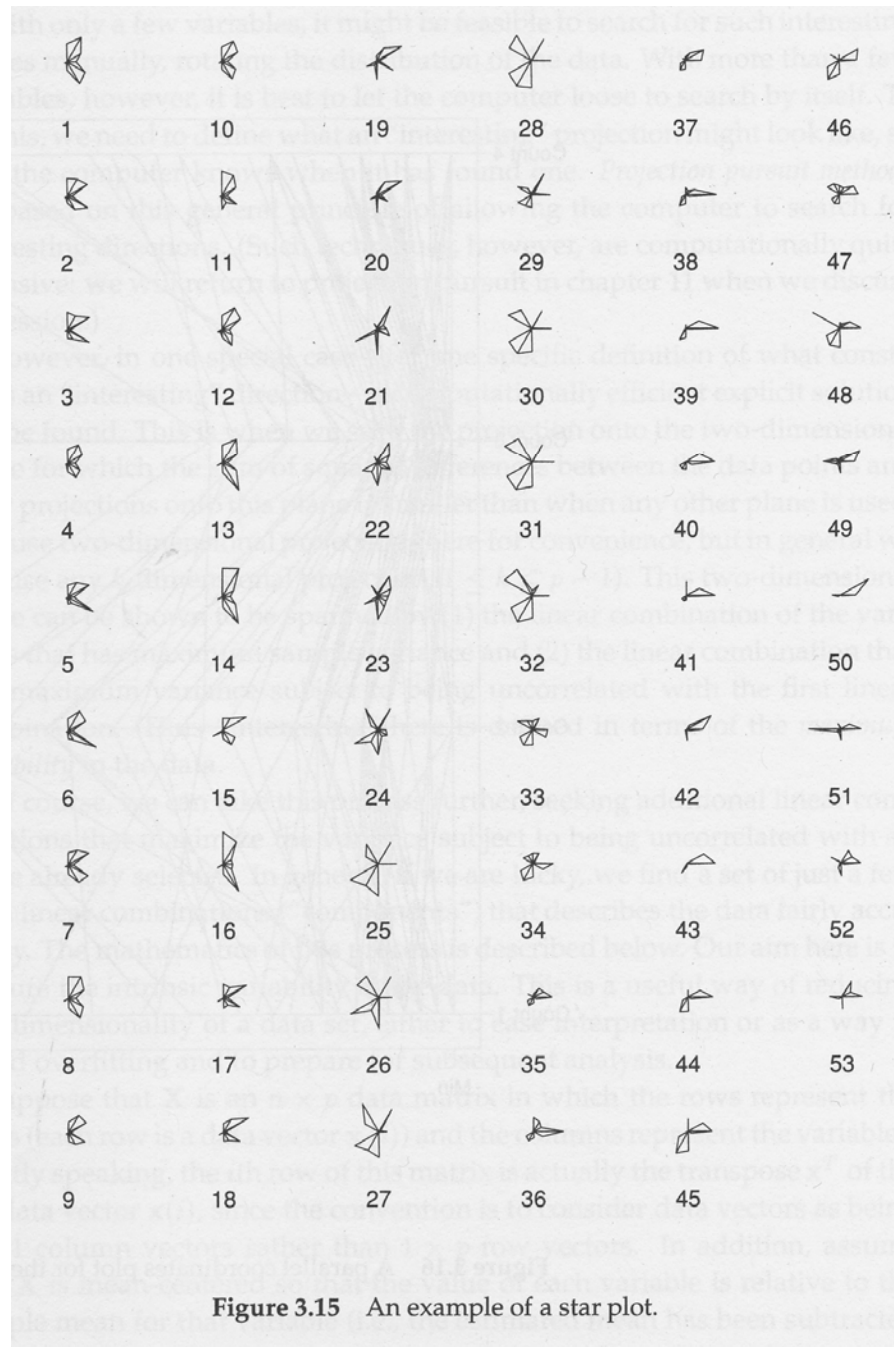


Figure 3.15 An example of a star plot.

Parallel Coordinates Plot

Each path represents
an individual

Each count
Represents 2-week
period

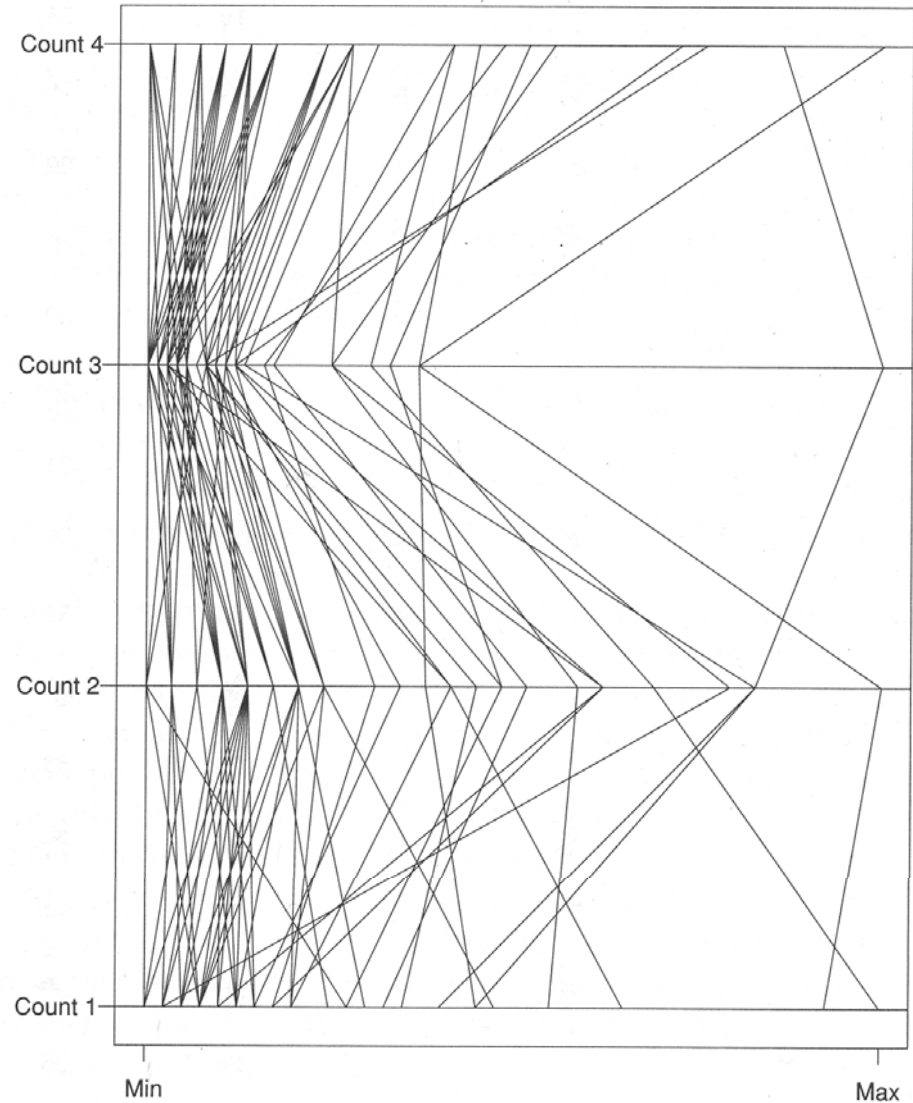


Figure 3.16 A parallel coordinates plot for the epileptic seizure data.