# Data Mining:

## Concepts and Techniques

## — Chapter 2 —

## Ali Shakiba

## Vali-e-Asr University of Rafsanjan

based on slides by

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign
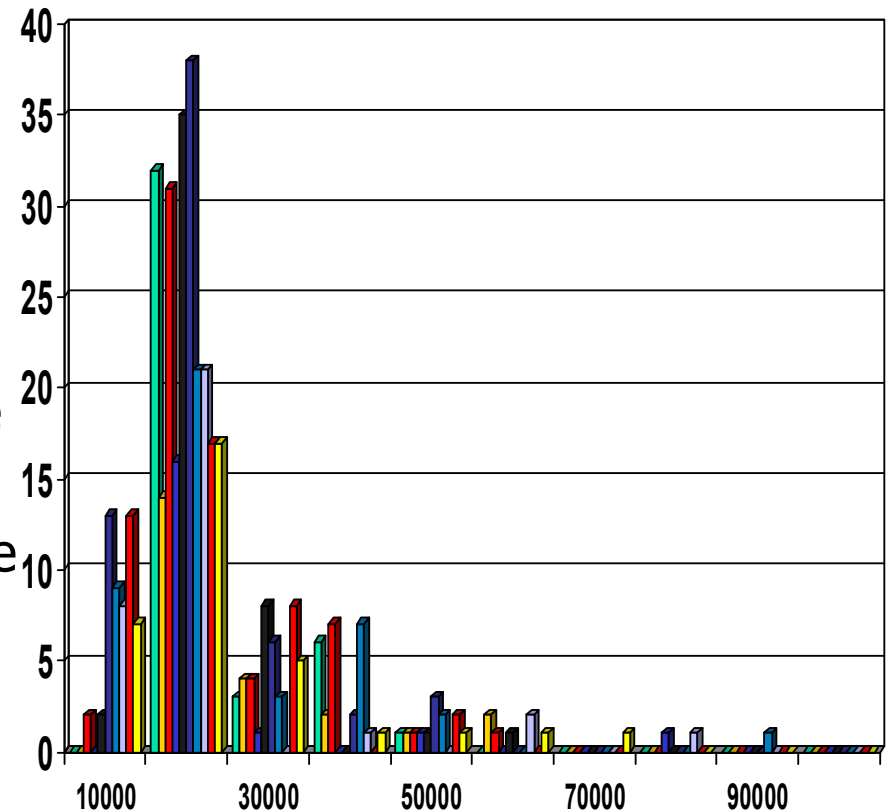
Simon Fraser University

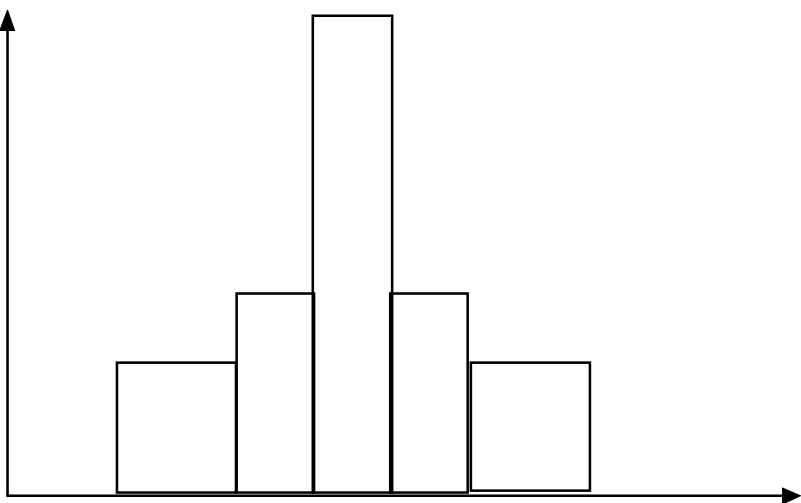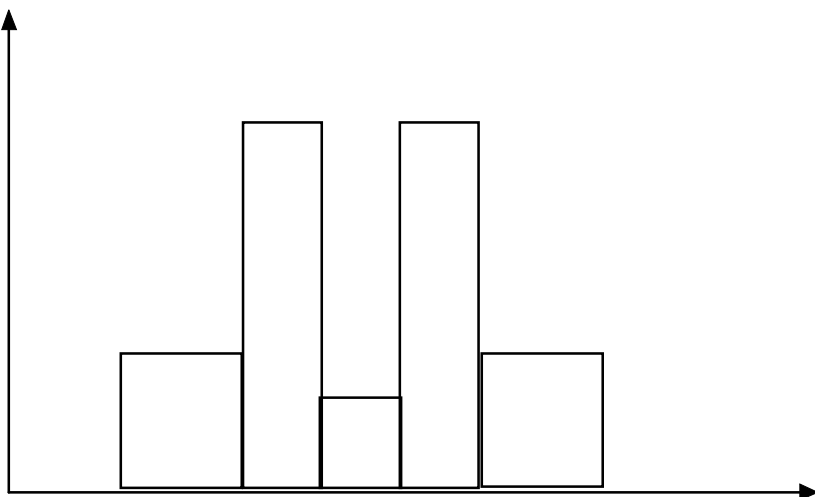# Graphic Displays of Basic Statistical Descriptions

- **Boxplot**: graphic display of five-number summary

- **Histogram**: x-axis are values, y-axis repres. frequencies

- **Quantile plot**: each value $x_i$ is paired with $f_i$ indicating that approximately 100 $f_i$% of data are $\leq x_i$

- **Quantile-quantile (q-q) plot**: graphs the quantiles of one univariant distribution against the corresponding quantiles of another

- **Scatter plot**: each pair of values is a pair of coordinates and plotted as points in the plane

# Histogram Analysis

- Histogram: Graph display of tabulated frequencies, shown as bars

- It shows what proportion of cases fall into each of several categories

- Differs from a bar chart in that it is the *area* of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width

- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent
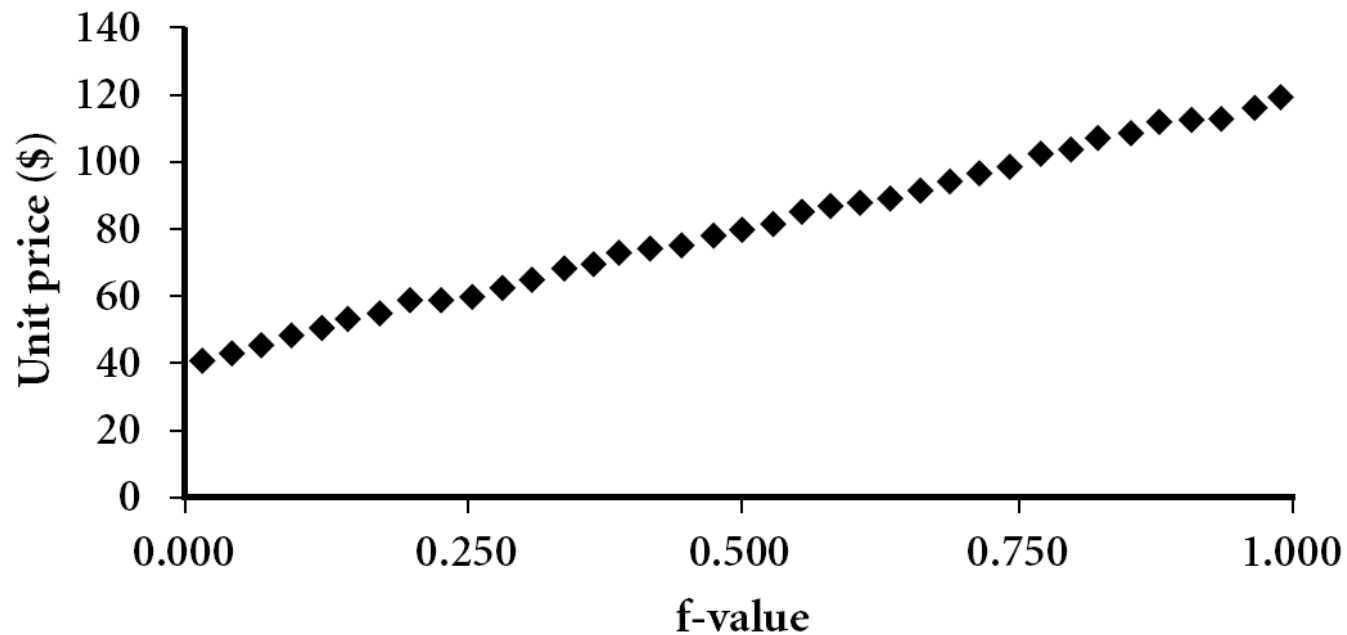
# Histograms Often Tell More than Boxplots



- The two histograms shown in the left may have the same boxplot representation
  - The same values for: min, Q1, median, Q3, max
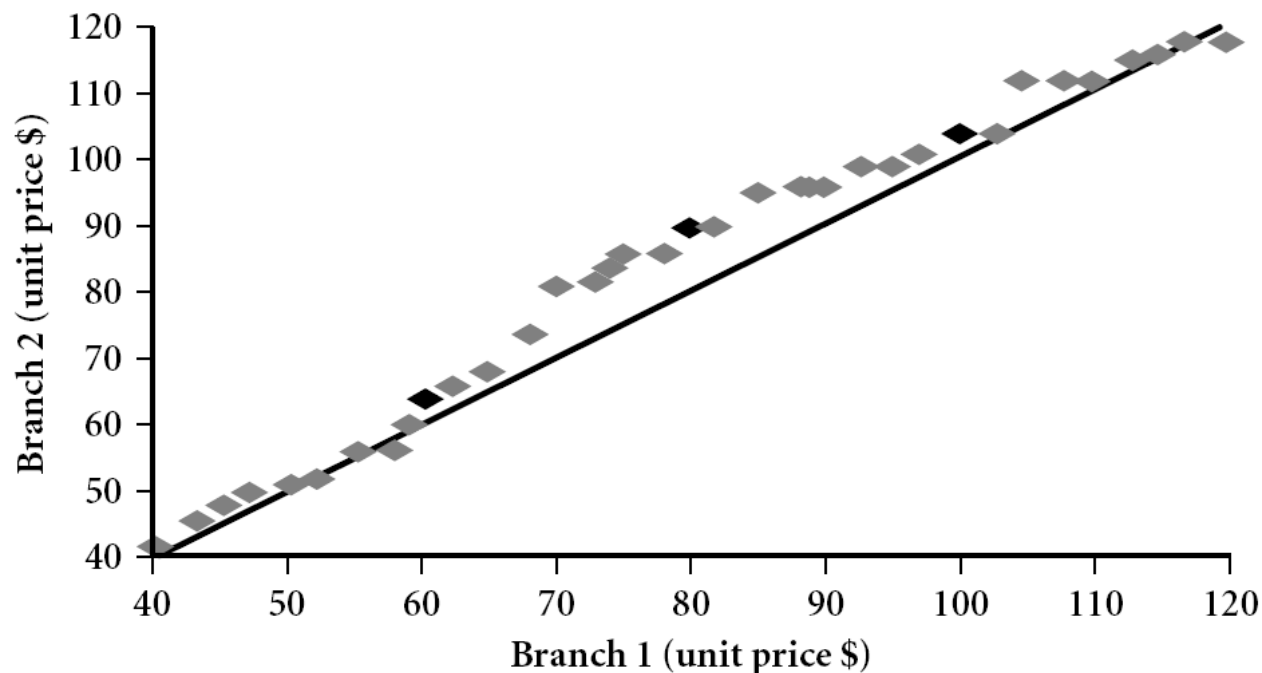- But they have rather different data distributions

# Quantile Plot

- Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)
- Plots **quantile** information
  - For a data $x_i$ data sorted in increasing order, $f_i$ indicates that approximately 100 $f_i$% of the data are below or equal to the value $x_i$
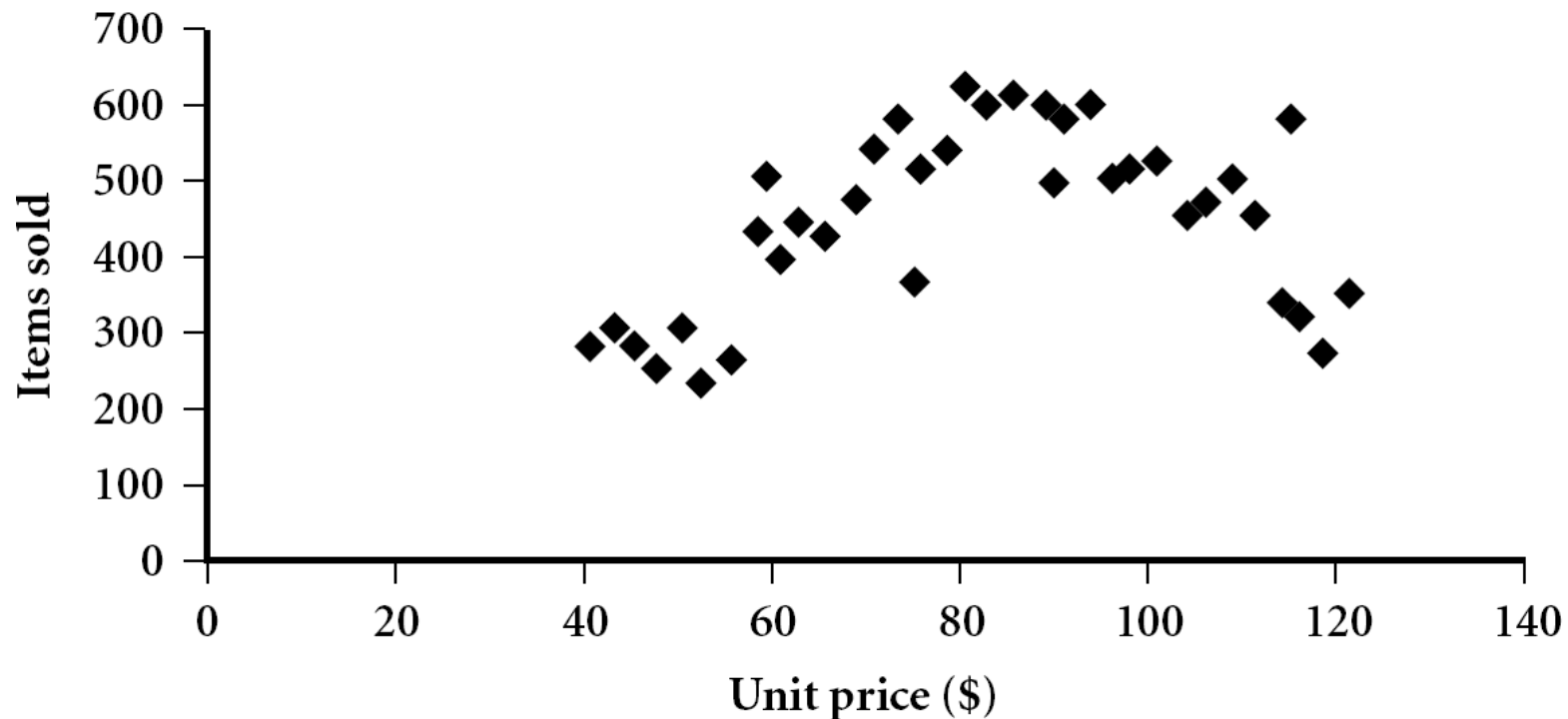
# Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?
- Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.
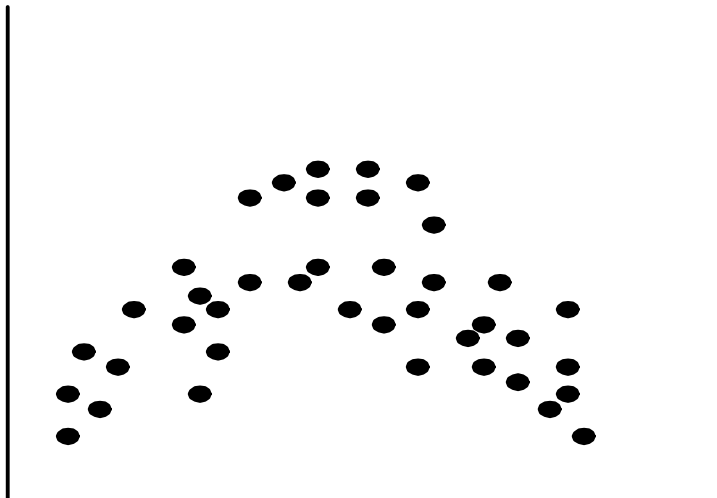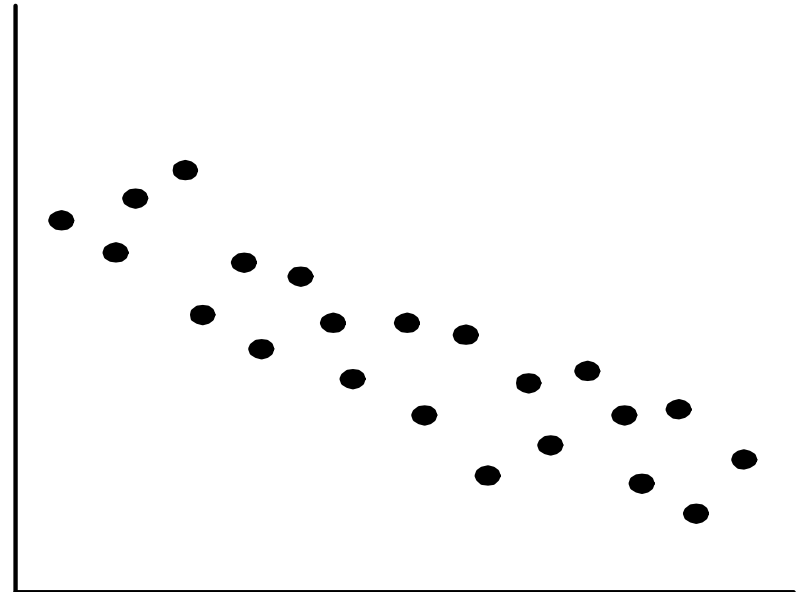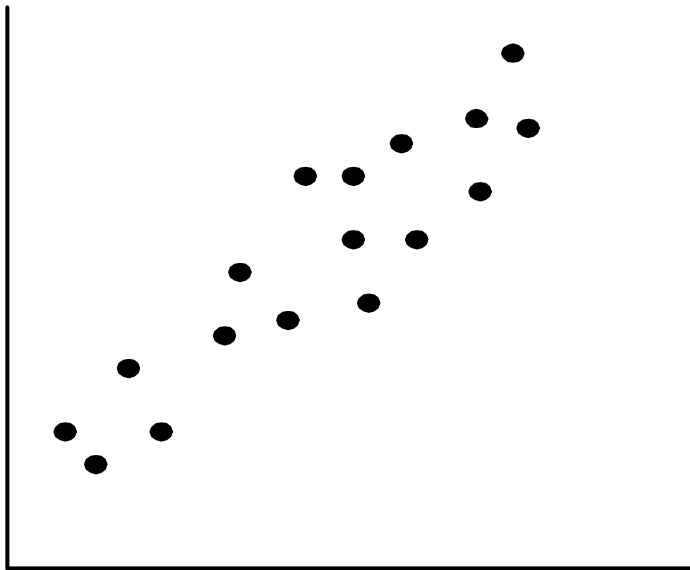
# Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc

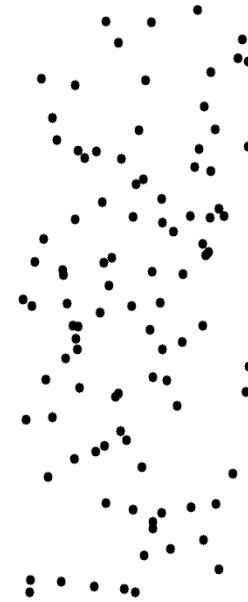- Each pair of values is treated as a pair of coordinates and plotted as points in the plane
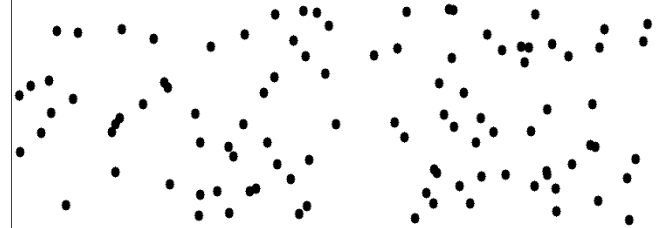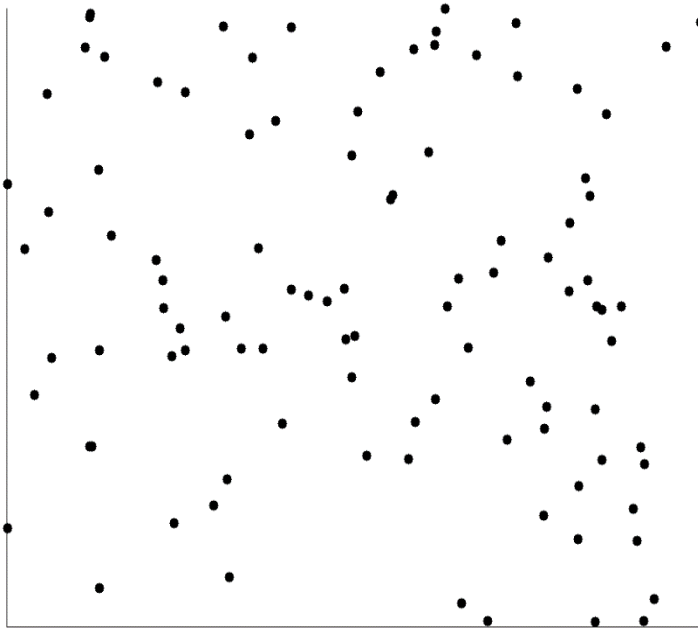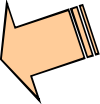
# Positively and Negatively Correlated Data



- The left half fragment is positively correlated
- The right half is negative correlated

# Uncorrelated Data

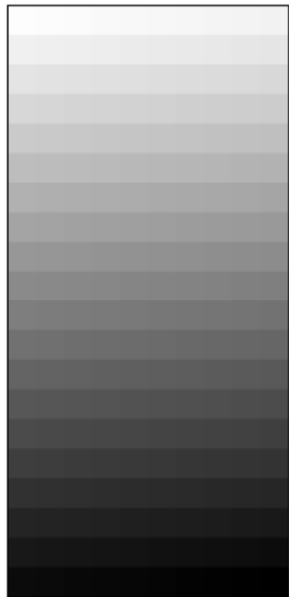# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- Basic Statistical Descriptions of Data

- Data Visualization

- Measuring Data Similarity and Dissimilarity

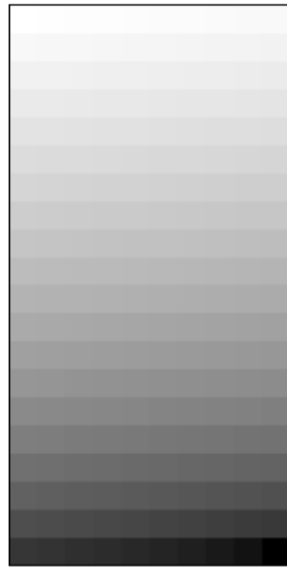- Summary

# Data Visualization

- Why data visualization?
    - Gain insight into an information space by mapping data onto graphical primitives
    - Provide qualitative overview of large data sets
    - Search for patterns, trends, structure, irregularities, relationships among data
    - Help find interesting regions and suitable parameters for further quantitative analysis
    - Provide a visual proof of computer representations derived
- Categorization of visualization methods:
    - Pixel-oriented visualization techniques
    - Geometric projection visualization techniques
    - Icon-based visualization techniques
    - Hierarchical visualization techniques
    - Visualizing complex data and relations
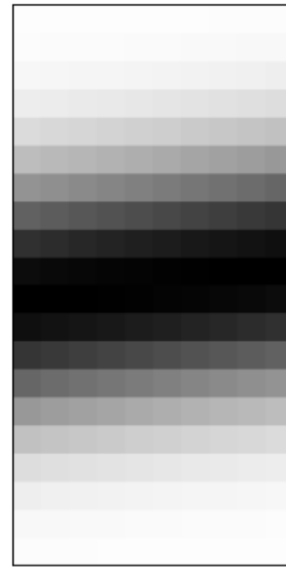
# Pixel-Oriented Visualization Techniques

- For a data set of m dimensions, create m windows on the screen, one for each dimension

- The m dimension values of a record are mapped to m pixels at the corresponding positions in the windows

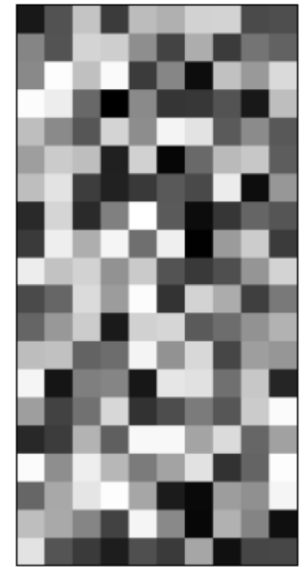- The colors of the pixels reflect the corresponding values
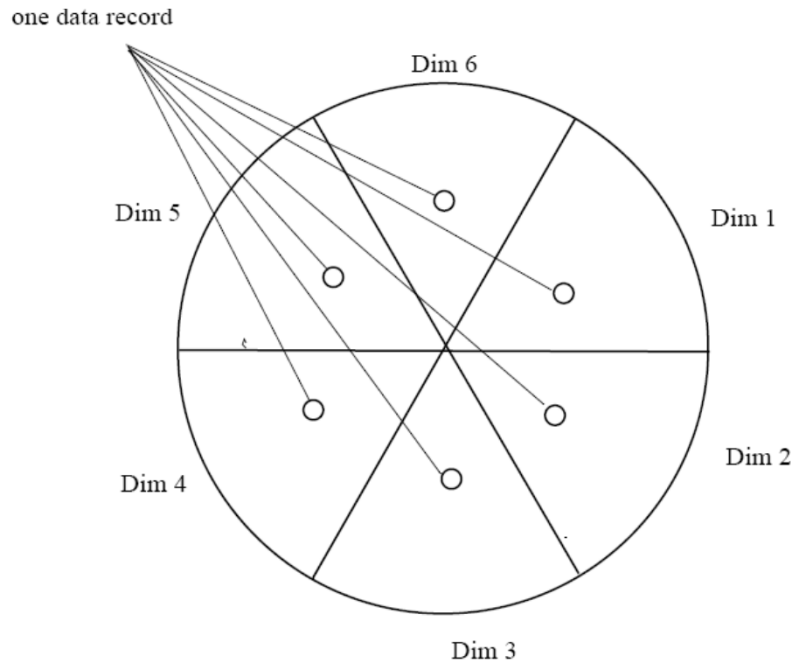


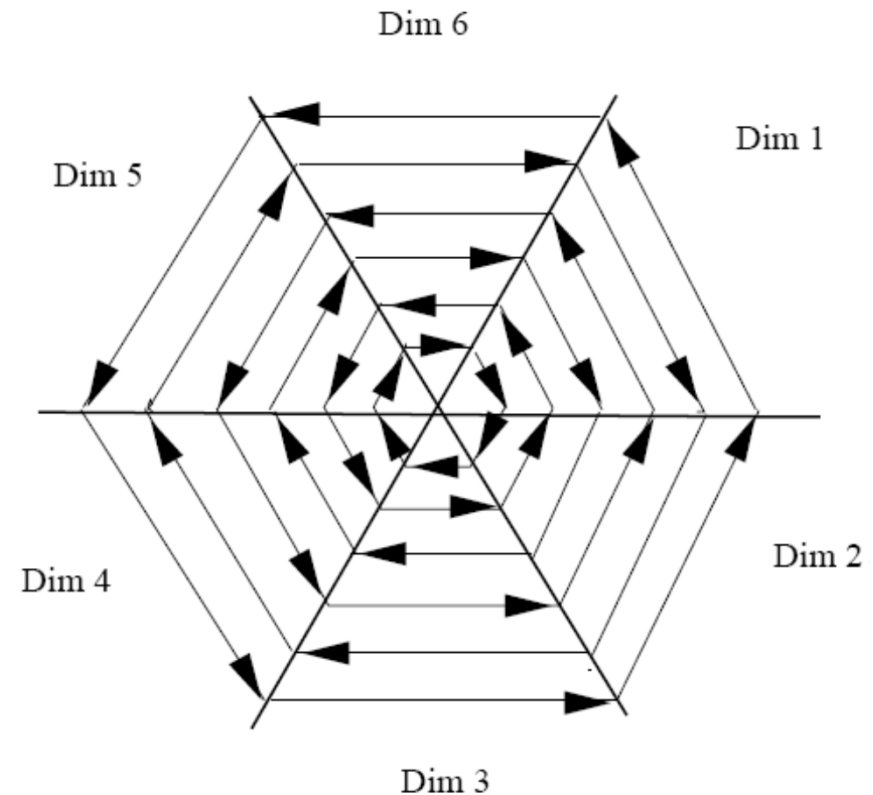(a) Income     (b) Credit Limit     (c) transaction volume     (d) age

# Laying Out Pixels in Circle Segments

- To save space and show the connections among multiple dimensions, space filling is often done in a circle segment
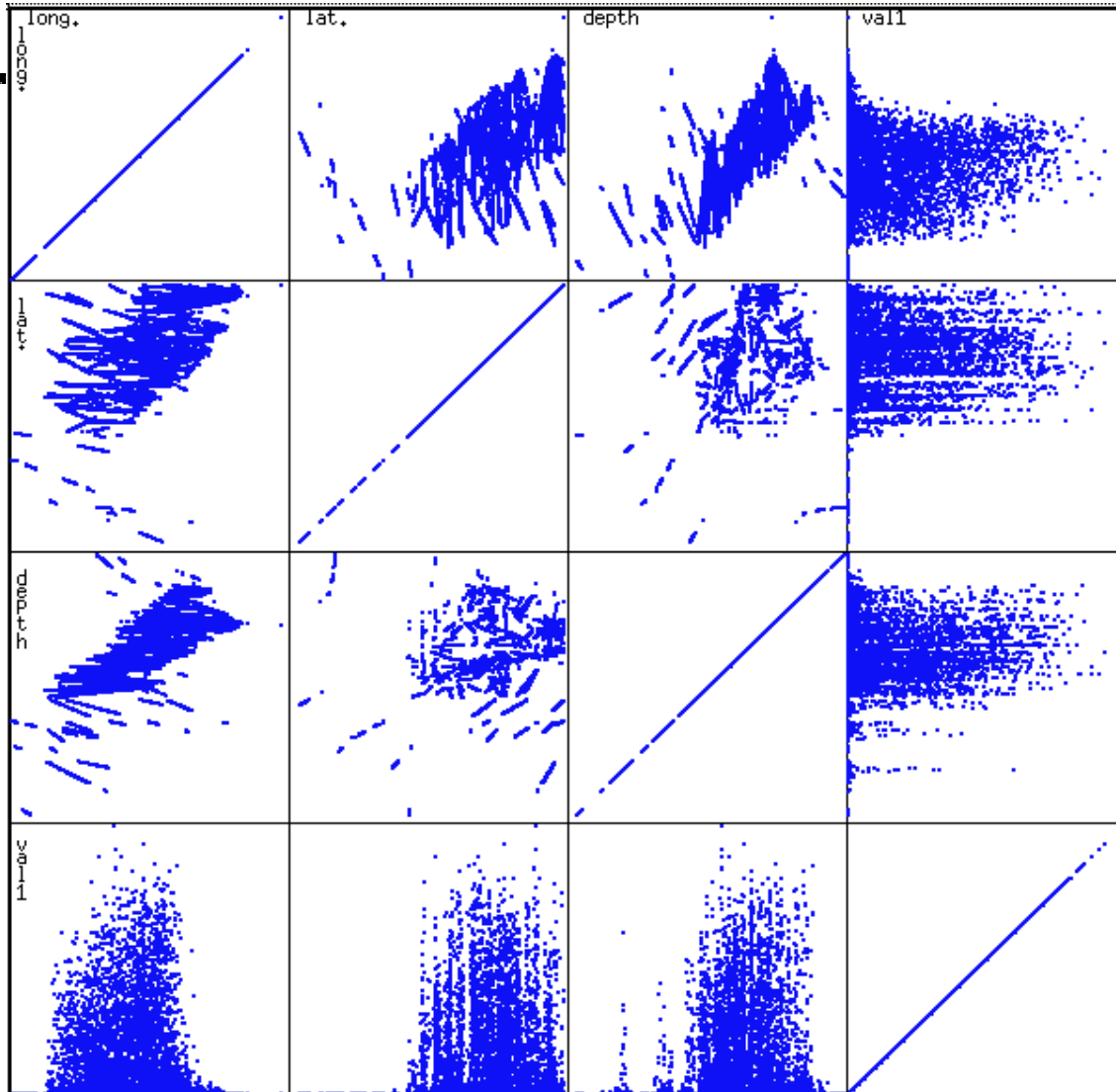


(a) Representing a data record in circle segment

(b) Laying out pixels in circle segment

# Geometric Projection Visualization Techniques

- Visualization of geometric transformations and projections of the data
- Methods
  - Direct visualization
  - Scatterplot and scatterplot matrices
  - Landscapes
  - Projection pursuit technique: Help users find meaningful projections of multidimensional data
  - Prosection views
  - Hyperslice
  - Parallel coordinates
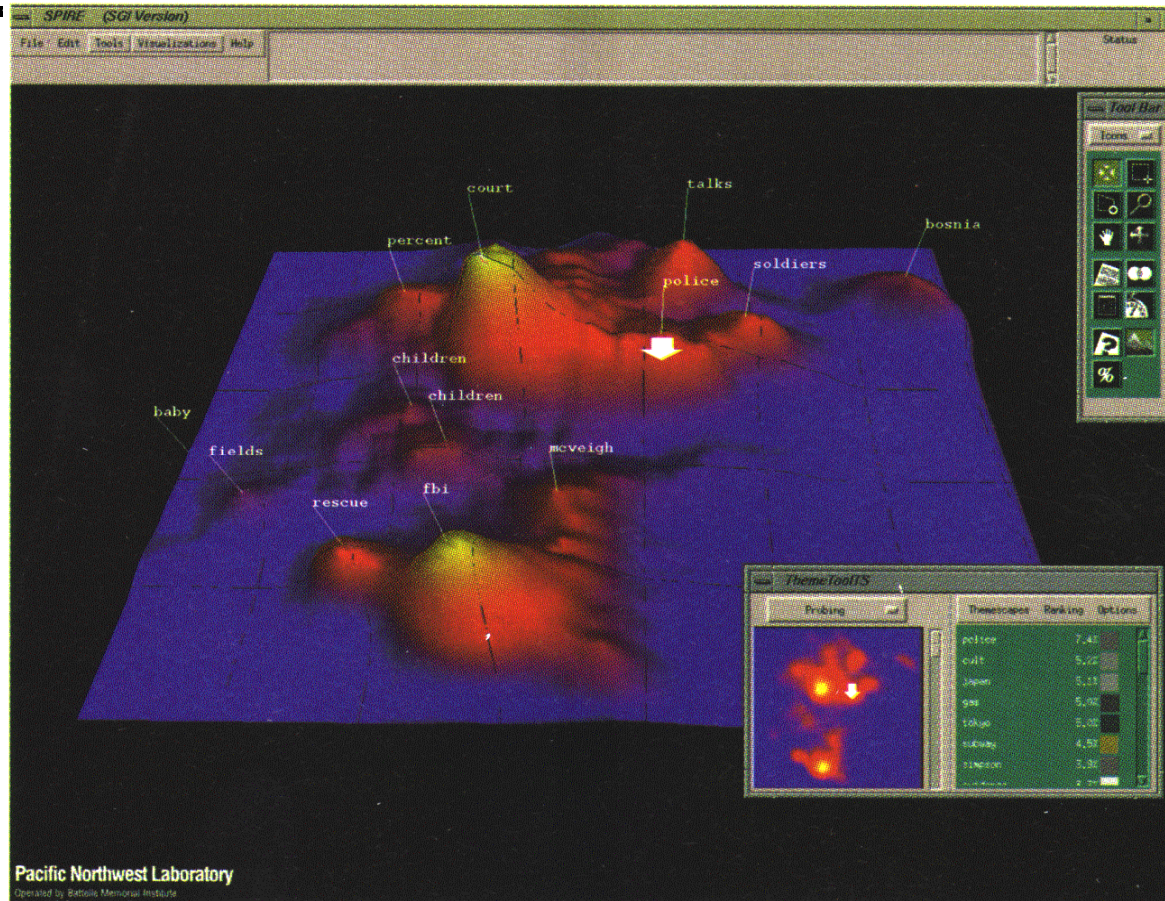
# Scatterplot Matrices



Matrix of scatterplots (x-y-diagrams) of the k-dim. data

30

# Landscapes



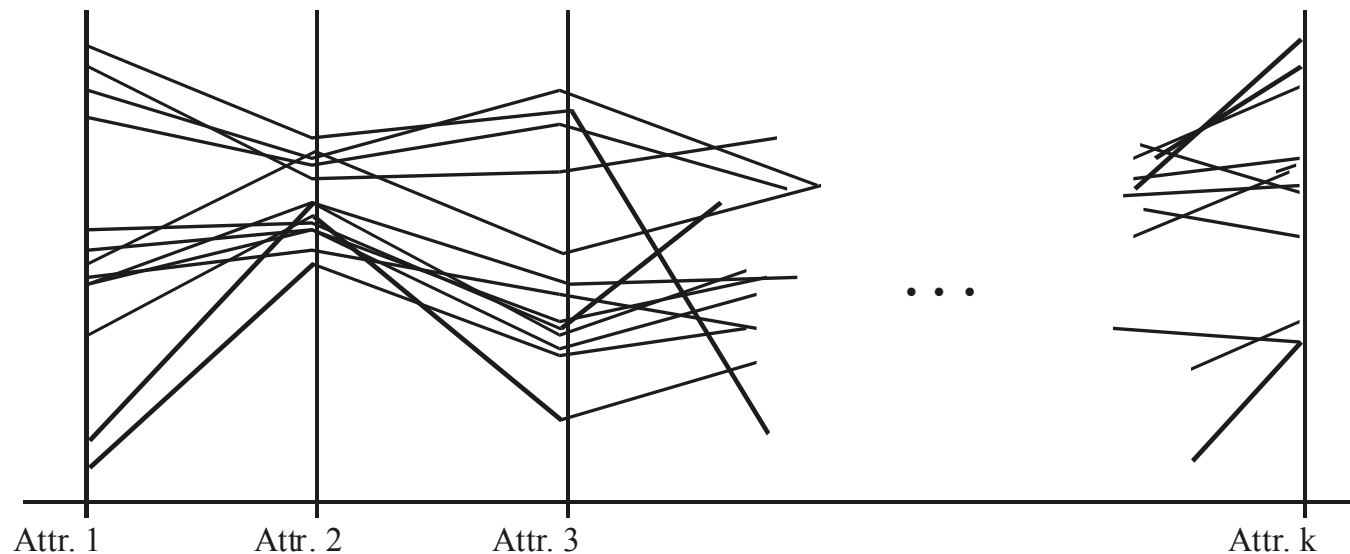Used by permission of B. Wright, Visible Decisions Inc.
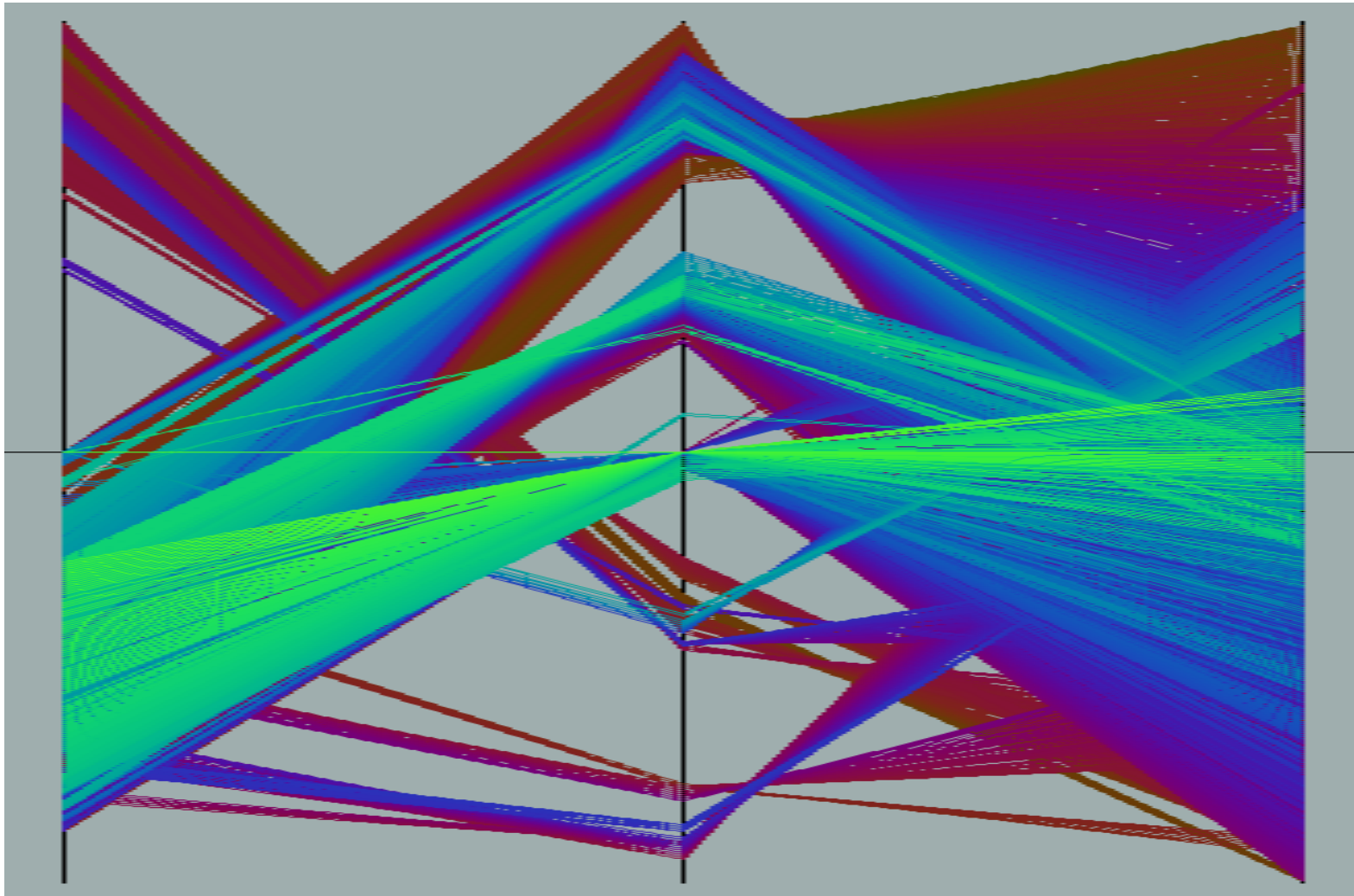
news articles visualized as a landscape

- Visualization of the data as perspective landscape
- The data needs to be transformed into a (possibly artificial) 2D spatial representation which preserves the characteristics of the data

31

# Parallel Coordinates

- n equidistant axes which are parallel to one of the screen axes and correspond to the attributes

- The axes are scaled to the [minimum, maximum]: range of the corresponding attribute

- Every data item corresponds to a polygonal line which intersects each of the axes at the point which corresponds to the value for the attribute
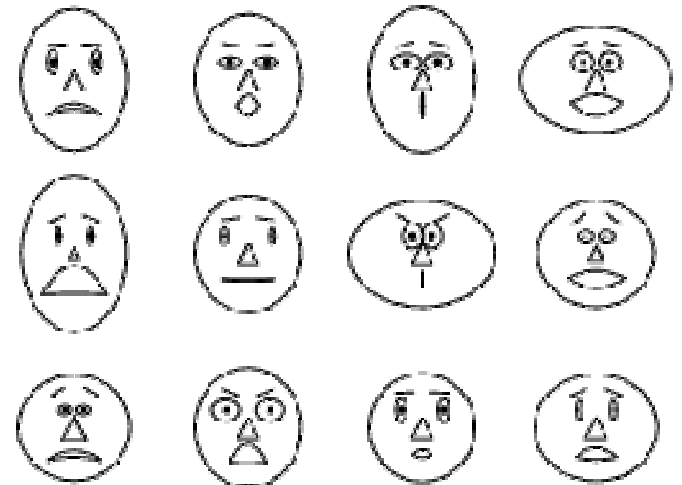
# Parallel Coordinates of a Data Set

# Icon-Based Visualization Techniques

- Visualization of the data values as features of icons

- Typical visualization methods

  - Chernoff Faces

  - Stick Figures

- General techniques

  - Shape coding: Use shape to represent certain information encoding

  - Color icons: Use color icons to encode more information

  - Tile bars: Use small icons to represent the relevant feature vectors in document retrieval
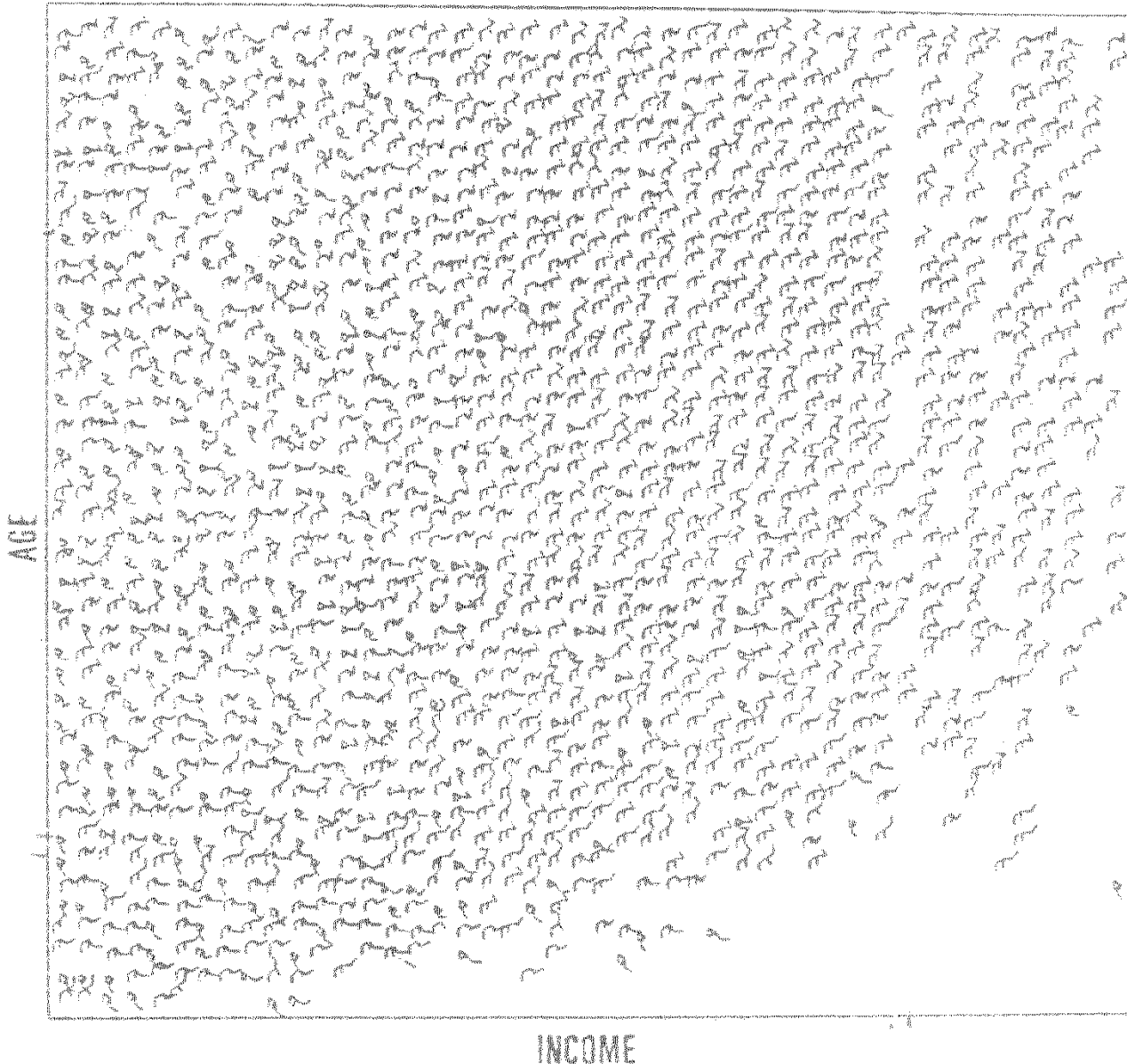
# Chernoff Faces

- A way to display variables on a two-dimensional surface, e.g., let x be eyebrow slant, y be eye size, z be nose length, etc.

- The figure shows faces produced using 10 characteristics--head eccentricity, eye size, eye spacing, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, mouth size, and mouth opening): Each assigned one of 10 possible values, generated using *Mathematica* (S. Dickson)

- REFERENCE: Gonick, L. and Smith, W. *The Cartoon Guide to Statistics.* New York: Harper Perennial, p. 212, 1993

- Weisstein, Eric W. "Chernoff Face." From *MathWorld*--A Wolfram Web Resource. mathworld.wolfram.com/ChernoffFace.html

# Stick Figure



A census data figure showing age, income, gender, education, etc.

A 5-piece stick figure (1 body and 4 limbs w. different angle/length)

INCOME

Two attributes mapped to axes, remaining attributes mapped to angle or length of limbs". Look at texture pattern

36