# Data Mining

3. Exploring Data

# What is data exploration?

**A preliminary exploration of the data to better understand its characteristics.**

- Key motivations of data exploration include
  - Helping to select the right tool for preprocessing or analysis
  - Making use of humans' abilities to recognize patterns
    - People can recognize patterns not captured by data analysis tools

- Related to the area of Exploratory Data Analysis (EDA)
  - Created by statistician John Tukey
  - Seminal book is Exploratory Data Analysis by Tukey
  - A nice online introduction can be found in Chapter 1 of the NIST Engineering Statistics Handbook

  http://www.itl.nist.gov/div898/handbook/index.htm

# Techniques Used In Data Exploration

- In EDA, as originally defined by Tukey
  - The focus was on visualization
  - Clustering and anomaly detection were viewed as exploratory techniques
  - In data mining, clustering and anomaly detection are major areas of interest, and not thought of as just exploratory

- In our discussion of data exploration, we focus on
  - Summary statistics
  - Visualization
  - Online Analytical Processing (OLAP)

# Iris Sample Data Set

☐ Many of the exploratory data techniques are illustrated with the Iris Plant data set.

– Can be obtained from the UCI Machine Learning Repository http://www.ics.uci.edu/~mlearn/MLRepository.html

– From the statistician Douglas Fisher

– Three flower types (classes):
  ◆ Setosa
  ◆ Virginica
  ◆ Versicolour

– Four (non-class) attributes
  ◆ Sepal width and length
  ◆ Petal width and length



Virginica. Robert H. Mohlenbrock. USDA NRCS. 1995. Northeast wetland flora: Field office guide to plant species. Northeast National Technical Center, Chester, PA. Courtesy of USDA NRCS Wetland Science Institute.

iris setosa — petal, sepal
iris versicolor — petal, sepal
iris virginica — petal, sepal

# Summary Statistics

☐ Summary statistics  are numbers that summarize properties of the data

– Summarized properties include frequency, location and spread

◆ Examples:   location - mean
                 spread - standard deviation

– Most summary statistics can be calculated in a single pass through the data

# Frequency and Mode

- The frequency of an attribute value is the percentage of time the value occurs in the data set

  - For example, given the attribute 'gender' and a representative population of people, the gender 'female' occurs about 50% of the time.

- The mode of an attribute is the most frequent attribute value

- The notions of frequency and mode are typically used with categorical data

# Percentiles

- For continuous data, the notion of a percentile is more useful.

Given an ordinal or continuous attribute $x$ and a number $p$ between 0 and 100, the $p$th percentile $x_p$ is a value of $x$ such that $p\%$ of the observed values of $x$ are less than $x_p$.

- For instance, the 50th percentile is the value $x_{50\%}$ such that 50% of all values of $x$ are less than $x_{50\%}$.

# Measures of Location: Mean and Median

- The mean is the most common measure of the location of a set of points.

- However, the mean is very sensitive to outliers.

- Thus, the median or a trimmed mean is also commonly used.

$$\text{mean}(x) = \overline{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$$

$$\text{median}(x) = \begin{cases} x_{(r+1)} & \text{if } m \text{ is odd, i.e., } m = 2r + 1 \\ \frac{1}{2}(x_{(r)} + x_{(r+1)}) & \text{if } m \text{ is even, i.e., } m = 2r \end{cases}$$

# Measures of Spread: Range and Variance

Range is the difference between the max and min

The variance or standard deviation is the most common measure of the spread of a set of points.

$$\text{variance}(x) = s_x^2 = \frac{1}{m-1} \sum_{i=1}^{m} (x_i - \overline{x})^2$$

However, this is also sensitive to outliers, so that other measures are often used.

Absolute Average Deviation (AAD),

Median Absolute Deviation (MAD),

Interquartile Range (IQ)

$$\text{AAD}(x) = \frac{1}{m} \sum_{i=1}^{m} |x_i - \overline{x}|$$

$$\text{MAD}(x) = median\left(\{|x_1 - \overline{x}|, \dots, |x_m - \overline{x}|\}\right)$$

$$\text{interquartile range}(x) = x_{75\%} - x_{25\%}$$

| Measure of central tendency m(X) | Mean absolute deviation |
|---|---|
| Mean = 5 | \|2-5\|+\|2-5\|+\|3-5\|+\|4-5\|+\|14-5\| / 5 = 3.6 |
| Median = 3 | \|2-3\|+\|2-3\|+\|3-3\|+\|4-3\|+\|14-3\| / 5 = 2.8 |
| Mode = 2 | \|2-2\|+\|2-2\|+\|3-2\|+\|4-2\|+\|14-2\| / 5 = 3.0 |

Seq = { 2  2  3  4  14 }.

# Basic Statistical Descriptions of Data

- Motivation
  - To better understand the data: central tendency, variation and spread
- Data dispersion characteristics
  - median, max, min, quantiles, outliers, variance, etc.
- Numerical dimensions correspond to sorted intervals
  - Data dispersion: analyzed with multiple granularities of precision
  - Boxplot or quantile analysis on sorted intervals
- Dispersion analysis on computed measures
  - Folding measures into numerical dimensions
  - Boxplot or quantile analysis on the transformed cube

# Measuring the Central Tendency

- Mean (algebraic measure) (sample vs. population):

  Note: $n$ is sample size and $N$ is population size.

  $$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \qquad \mu = \frac{\sum x}{N}$$

  Weighted arithmetic mean:

  $$\overline{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

  Trimmed mean: chopping extreme values

# Measuring the Central Tendency

- Median:
    - Middle value if odd number of values, or average of the middle two values otherwise
    - Estimated by interpolation (for *grouped data*):

      **Estimated Median** $= L + \left( \dfrac{(n/2) - B}{G} \right) \times w$

| age | frequency |
|-----|-----------|
| 1–5 | 200 |
| 6–15 | 450 |
| 16–20 | 300 |
| 21–50 | 1500 ← |
| 51–80 | 700 |
| 81–110 | 44 |

3194

where:
- **L** is the lower class boundary of the group containing the median  20.5
- **n** is the total number of values  3194
- **B** is the cumulative frequency of the groups before the median group 200 + 450 + 300 = 950
- **G** is the frequency of the median group  1500
- **w** is the group width  = 5

Median # 1597

- Mode

  $mean - mode = 3 \times (mean - median)$

    - Value that occurs most frequently in the data
    - Unimodal, bimodal, trimodal

  Empirical formula:

Given two school classes:

Morning class (20 students) = 62, 67, 71, 74, 76, 77, 78, 79, 79, 80, 80, 81, 81, 82, 83, 84, 86, 89, 93, 98  (Mean = 80)

Afternoon class (30 students) = 81, 82, 83, 84, 85, 86, 87, 87, 88, 88, 89, 89, 89, 90, 90, 90, 90, 91, 91, 91, 92, 92, 93, 93, 94, 95, 96, 97, 98, 99  (Mean = 90)

The unweighted mean of the (80 + 90) / 2  = 85

Do not account for the difference in number of students in each class (20 versus 30);
Hence the value of 85 does not reflect the average student grade (independent of class).

The average student grade can be obtained by : adding all the grades up and divide by the total number of students:

$$\bar{x} = 4300 / 50 = 86$$

Or, this can be accomplished by weighting the class means by the number of students in each class. The larger class is given more "weight":

$$\bar{x} = \frac{(20 \times 80) + (30 \times 90)}{20 + 30} = 86$$

Thus, the weighted mean makes it possible to find the mean average student grade without knowing each student's score. Only the class means and the number of students in each class are needed.

# Symmetric vs. Skewed Data

- Median, mean and mode of symmetric, positively and negatively skewed data

symmetric

Mean
Median
Mode

positively skewed

Mode
Mean
Median

negatively skewed

Mean    Mode
Median

# Measuring the Dispersion of Data

- ☐ Quartiles, outliers and boxplots

  - **Quartiles**: $Q_1$ ($25^{th}$ percentile), $Q_3$ ($75^{th}$ percentile)

  - **Inter-quartile range**: IQR = $Q_3 - Q_1$

  - **Five number summary**: min, $Q_1$, median, $Q_3$, max

  - **Boxplot**: ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually

  - **Outlier**: usually, a value higher/lower than 1.5 x IQR

- ☐ Variance and standard deviation (*sample: s, population: σ*)

  - **Variance**: (algebraic, scalable computation)

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

  - **Standard deviation** *s (or σ)* is the square root of variance $s^2$ *(or $\sigma^2$)*

# Boxplot Analysis



- **Five-number summary** of a distribution
  - Minimum, Q1, Median, Q3, Maximum

- **Boxplot**
  - Data is represented with a box
  - The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
  - The median is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
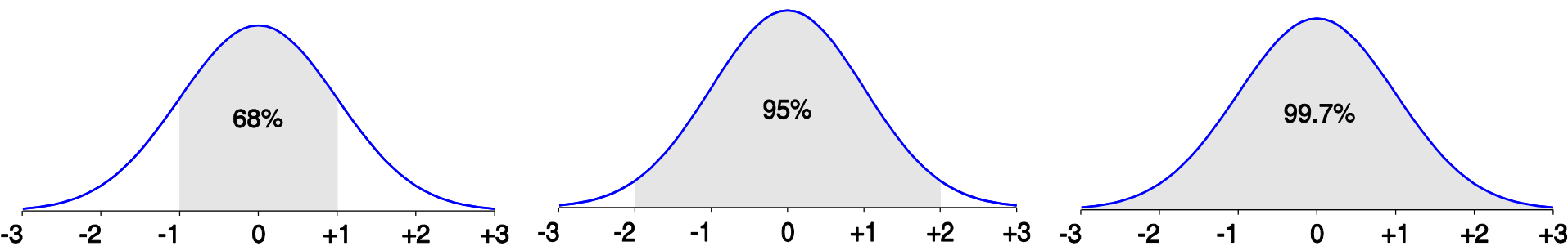  - Outliers: points beyond a specified outlier threshold, plotted individually

# Visualization of Data Dispersion: 3-D Boxplots

# **Properties of Normal Distribution Curve**

☐ The normal (distribution) curve

- From μ–σ to μ+σ: contains about 68% of the measurements (μ: mean, σ: standard deviation)

- From μ–2σ to μ+2σ: contains about 95% of it

- From μ–3σ to μ+3σ: contains about 99.7% of it

| 68% | 95% | 99.7% |

-3 -2 -1 0 +1 +2 +3    -3 -2 -1 0 +1 +2 +3    -3 -2 -1 0 +1 +2 +3

# 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

# Visualization

Visualization is the conversion of data into a visual or tabular format so that the characteristics of the data and the relationships among data items or attributes can be analyzed or reported.

☐ Visualization of data is one of the most powerful and appealing techniques for data exploration.

- Humans have a well developed ability to analyze large amounts of information that is presented visually

- Can detect general patterns and trends

- Can detect outliers and unusual patterns

# 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

# 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

# Example: Sea Surface Temperature

☐ The following shows the Sea Surface Temperature (SST) for July 1982

  – Tens of thousands of data points are summarized in a single figure

# Representation

- Is the mapping of information to a <span style="color:red">visual format</span>
- Data <span style="color:red">objects, their attributes, and the relationships</span> among data objects are translated into graphical elements such as <span style="color:red">points, lines, shapes, and colors.</span>
- Example:
  - Objects are often represented as points
  - Their attribute values can be represented as the position of the points or the characteristics of the points, e.g., color, size, and shape
  - If position is used, then the relationships of points, i.e., whether they form groups or a point is an outlier, is easily perceived.

# Arrangement

- Is the placement of visual elements within a display

- Can make a large difference in how easy it is to understand the data

- Example:  9 Objects and 6 Attributes

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 1 | 0 | 0 | 1 |
| 3 | 0 | 1 | 0 | 1 | 1 | 0 |
| 4 | 1 | 0 | 1 | 0 | 0 | 1 |
| 5 | 0 | 1 | 0 | 1 | 1 | 0 |
| 6 | 1 | 0 | 1 | 0 | 0 | 1 |
| 7 | 0 | 1 | 0 | 1 | 1 | 0 |
| 8 | 1 | 0 | 1 | 0 | 0 | 1 |
| 9 | 0 | 1 | 0 | 1 | 1 | 0 |

Re-arranged
Permuted →

|   | 6 | 1 | 3 | 2 | 5 | 4 |
|---|---|---|---|---|---|---|
| 4 | 1 | 1 | 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 0 | 0 | 0 |
| 6 | 1 | 1 | 1 | 0 | 0 | 0 |
| 8 | 1 | 1 | 1 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 1 | 1 | 1 |
| 3 | 0 | 0 | 0 | 1 | 1 | 1 |
| 9 | 0 | 0 | 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 | 1 | 1 | 1 |
| 7 | 0 | 0 | 0 | 1 | 1 | 1 |

3-1's &
3-0's

3- 0's &
3- 1's

# Selection

- Is the <span style="color:red">elimination or the de-emphasis</span> of certain objects and attributes
- Selection may involve the <span style="color:red">choosing a subset</span> of attributes
  - Dimensionality reduction is often used to reduce the number of dimensions to two or three
  - Alternatively, pairs of attributes can be considered
- Selection may also involve choosing a subset of objects
  - A region of the screen can only show so many points
  - Can sample, but want to preserve points in sparse areas

# Visualization Techniques: Histograms

- Histogram
  - Usually shows the distribution of values of a single variable
  - Divide the values into bins and show a bar plot of the number of objects in each bin.
  - The height of each bar indicates the number of objects
  - Shape of histogram depends on the number of bins

- Example: Petal Width (10 and 20 bins, respectively)

# 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

# Two-Dimensional Histograms

- Show the joint distribution of the values of two attributes
- Example: petal width and petal length
    - What does this tell us?

# Visualization Techniques: Box Plots

- Box Plots
  - Invented by J. Tukey
  - Another way of displaying the distribution of data
  - Following figure shows the basic part of a box plot



outlier

10th percentile

75th percentile

50th percentile

25th percentile

10th percentile

# Example of Box Plots

☐ Box plots can be used to compare attributes

# Visualization Techniques: Scatter Plots

- Scatter plots
  - Attributes values determine the position
  - Two-dimensional scatter plots most common, but can have three-dimensional scatter plots
  - Often additional attributes can be displayed by using the size, shape, and color of the markers that represent the objects
  - It is useful to have arrays of scatter plots can compactly summarize the relationships of several pairs of attributes
    - See example on the next slide

# Scatter Plot Array of Iris Attributes

# 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 $x_i$ sorted in increasing order, $f_i$ indicates that approximately $f_i*100\%$ 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.
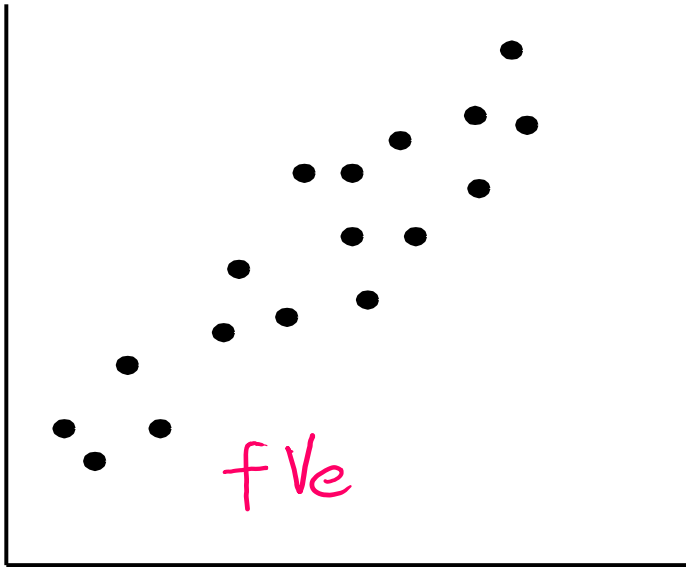
# Scatter plot

- 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
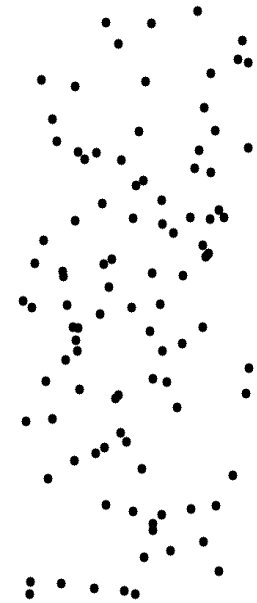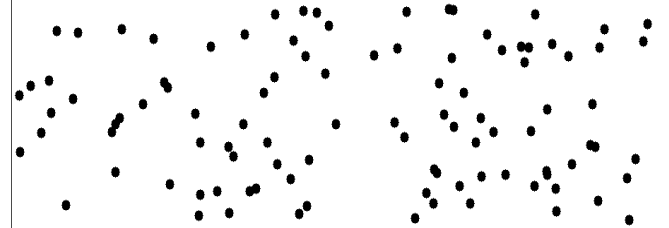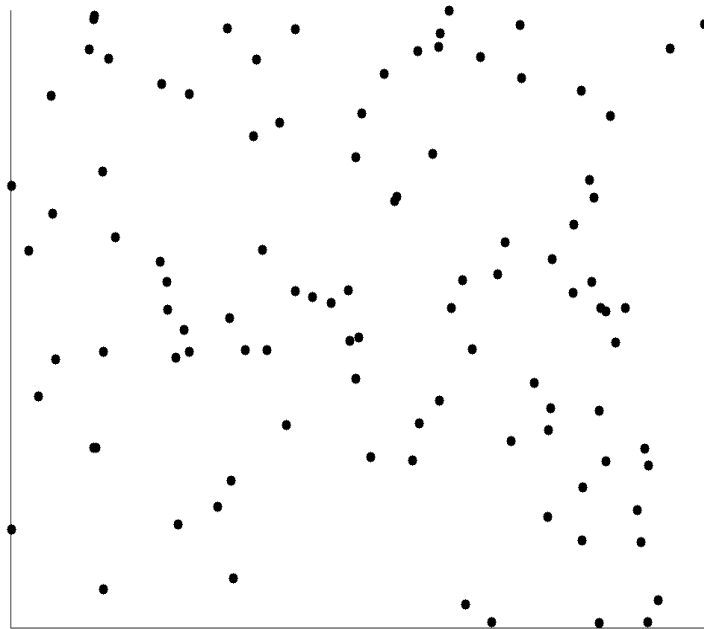
# Positively and Negatively Correlated Data



The handwritten annotations read: "+Ve" (top left plot), "−Ve" (top right plot), and "strong - Non-Linear" (bottom left plot).

- The left half fragment is positively correlated
- The right half is negative correlated
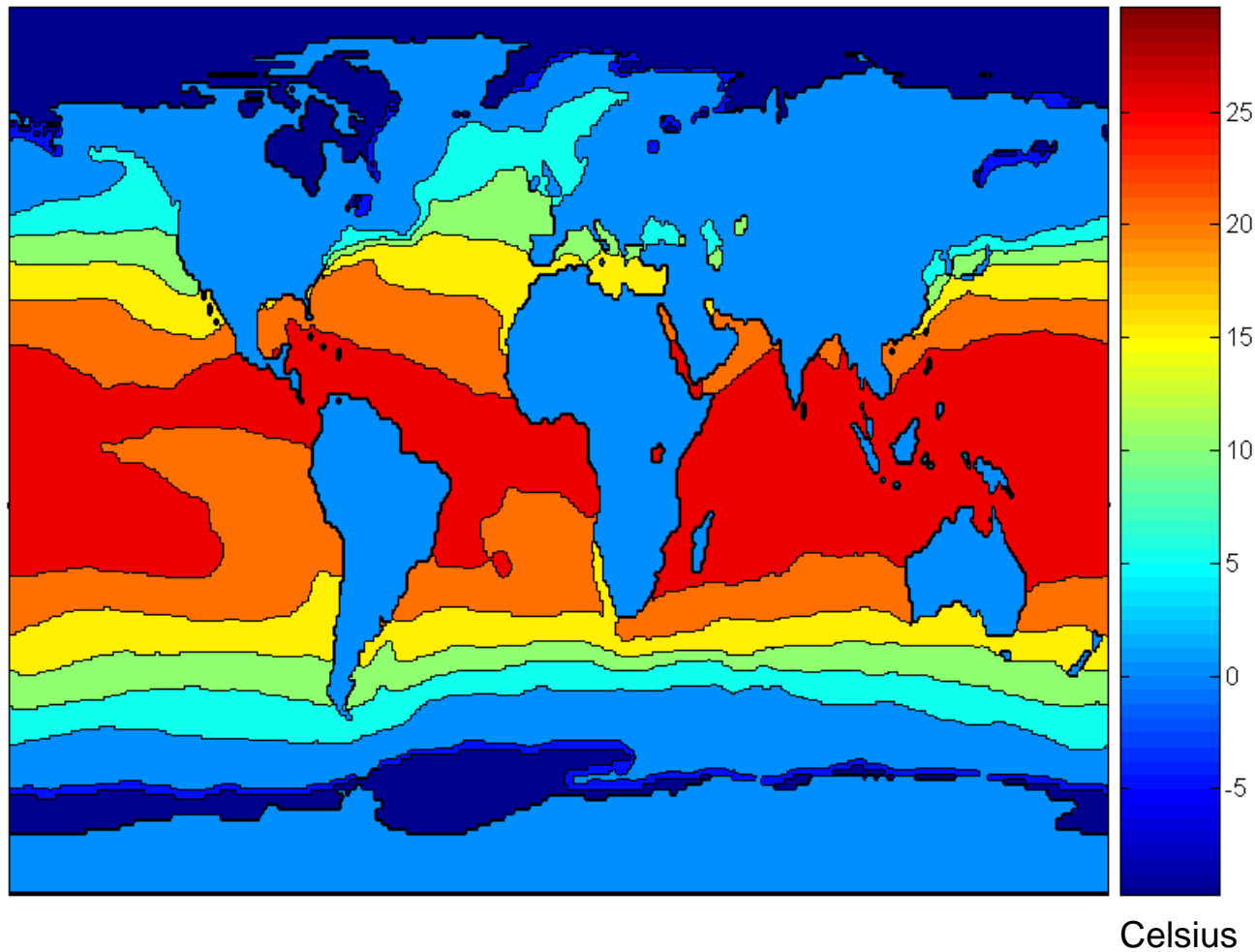
# Uncorrelated Data

# Visualization Techniques: Contour Plots

- Contour plots
  - Useful when a continuous attribute is measured on a spatial grid
  - They partition the plane into regions of similar values
  - The contour lines that form the boundaries of these regions connect points with equal values
  - The most common example is contour maps of elevation
  - Can also display temperature, rainfall, air pressure, etc.
    - An example for Sea Surface Temperature (SST) is provided on  the next slide

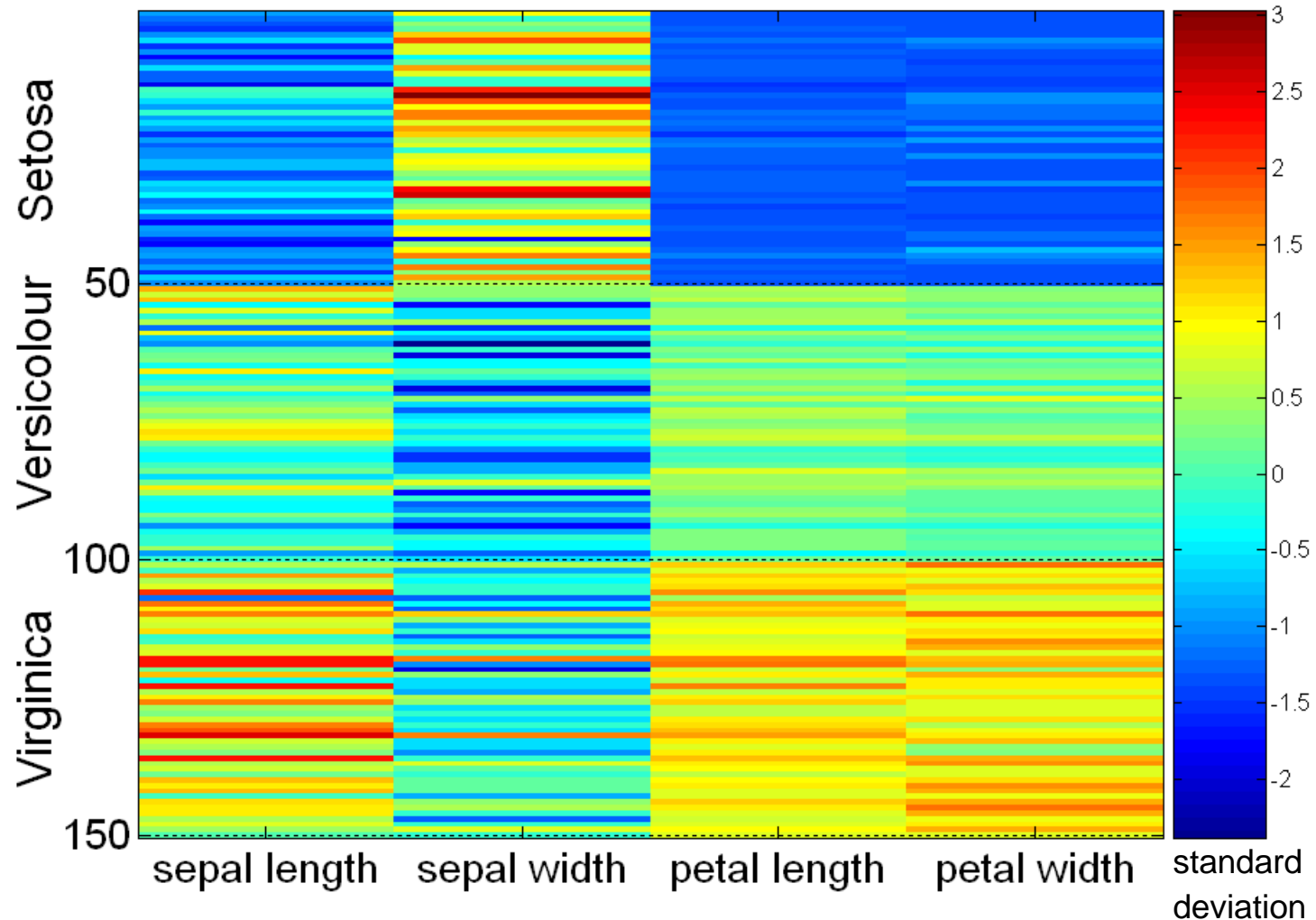# Contour Plot Example: SST Dec, 1998
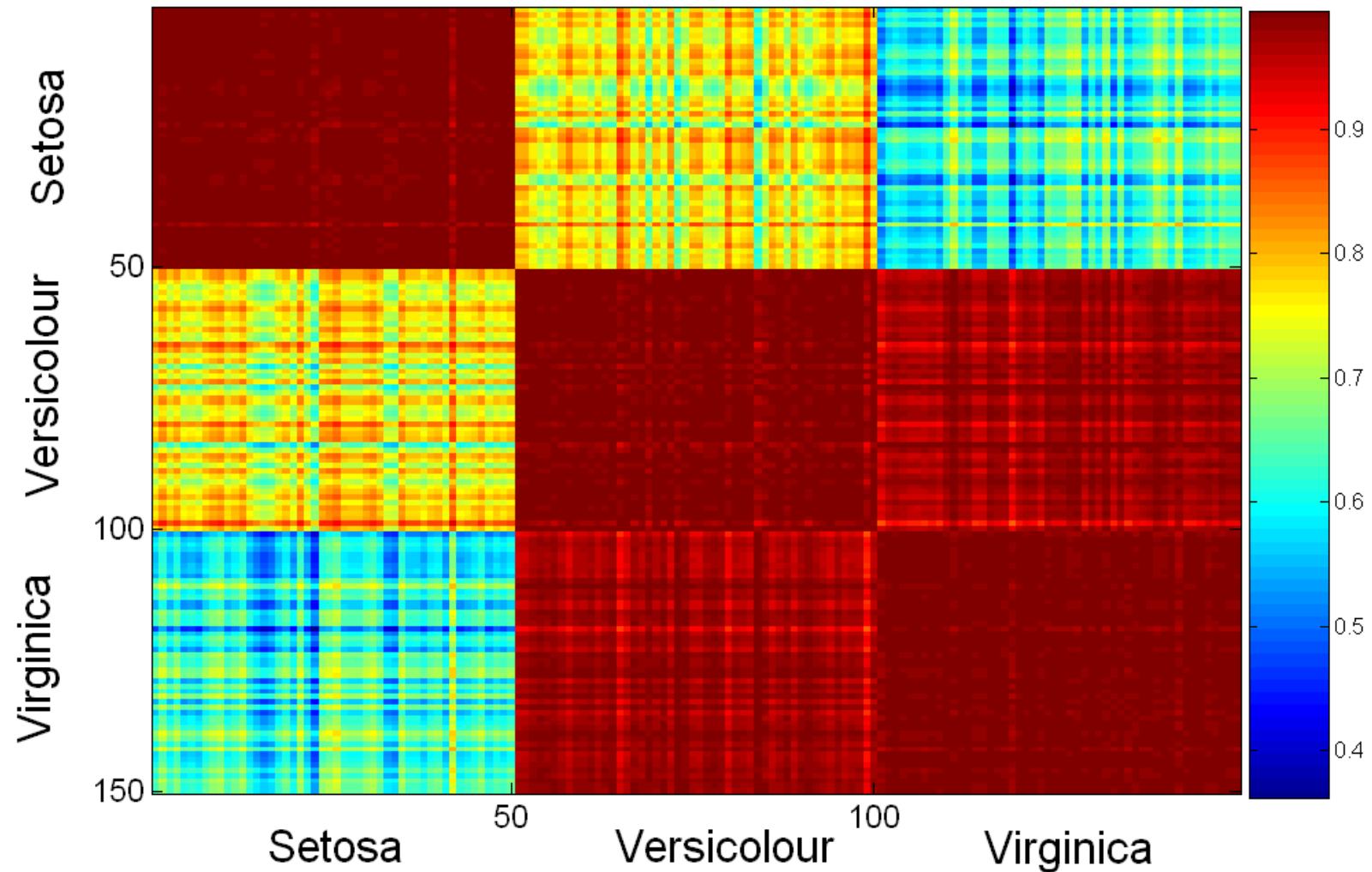


Celsius

# Visualization Techniques: Matrix Plots

☐ Matrix plots

- Can plot the data matrix

- This can be useful when objects are sorted according to class

- Typically, the attributes are normalized to prevent one attribute from dominating the plot

- Plots of similarity or distance matrices can also be useful for visualizing the relationships between objects

- Examples of matrix plots are presented on the next two slides

# Visualization of the Iris Data Matrix

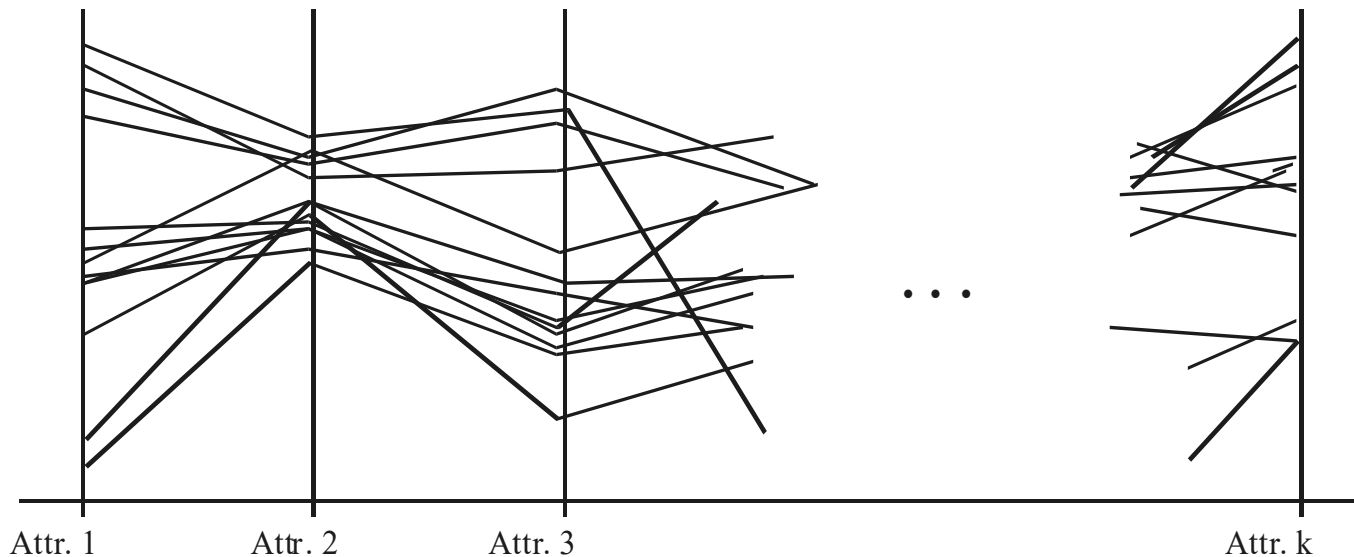# Visualization of the Iris Correlation Matrix

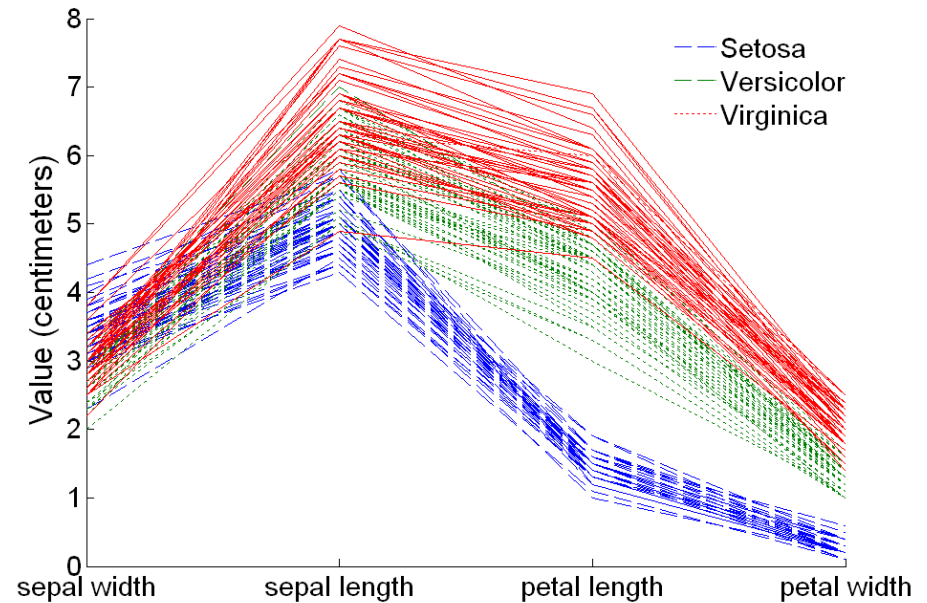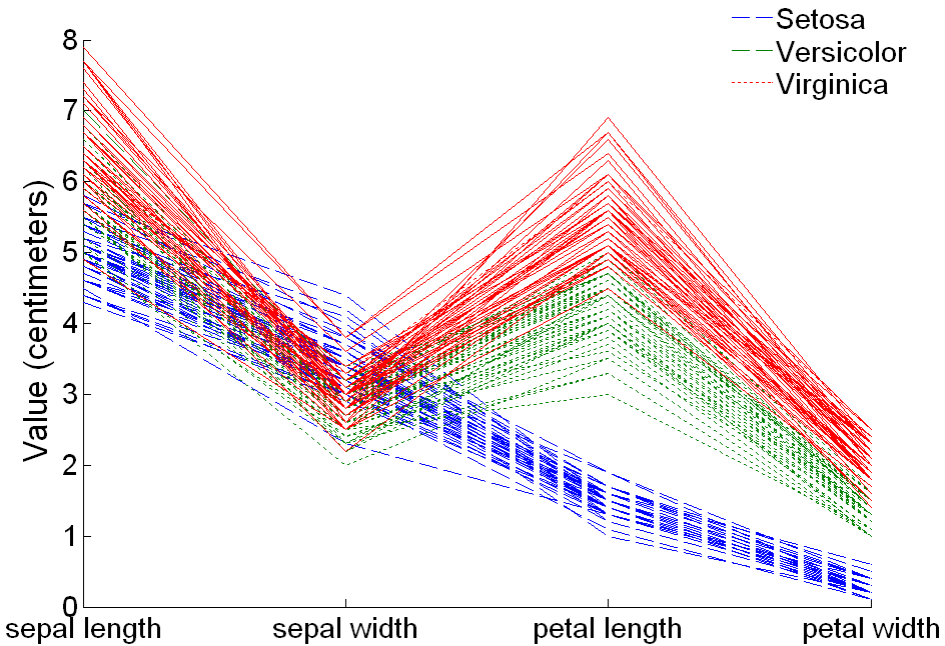# Visualization Techniques: Parallel Coordinates

- Parallel Coordinates
  - Used to plot the attribute values of high-dimensional data
  - Instead of using perpendicular axes, use a set of parallel axes
  - The attribute values of each object are plotted as a point on each corresponding coordinate axis and the points are connected by a line
  - Thus, each object is represented as a line
  - Often, the lines representing a distinct class of objects group together, at least for some attributes
  - Ordering of attributes is important in seeing such groupings
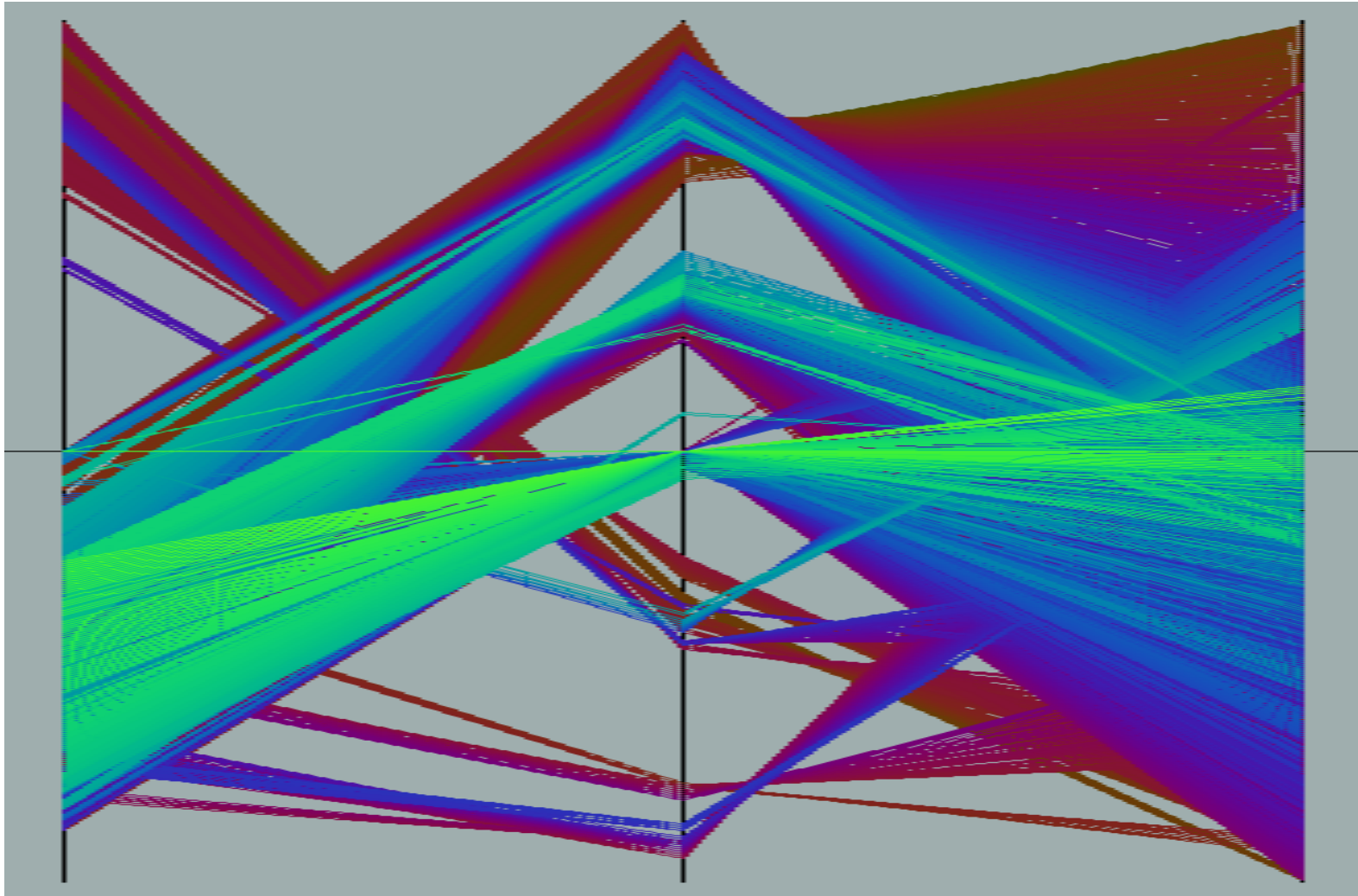
# Parallel Coordinates

- 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



Attr. 1          Attr. 2          Attr. 3                              Attr. k

# Parallel Coordinates Plots for Iris Data

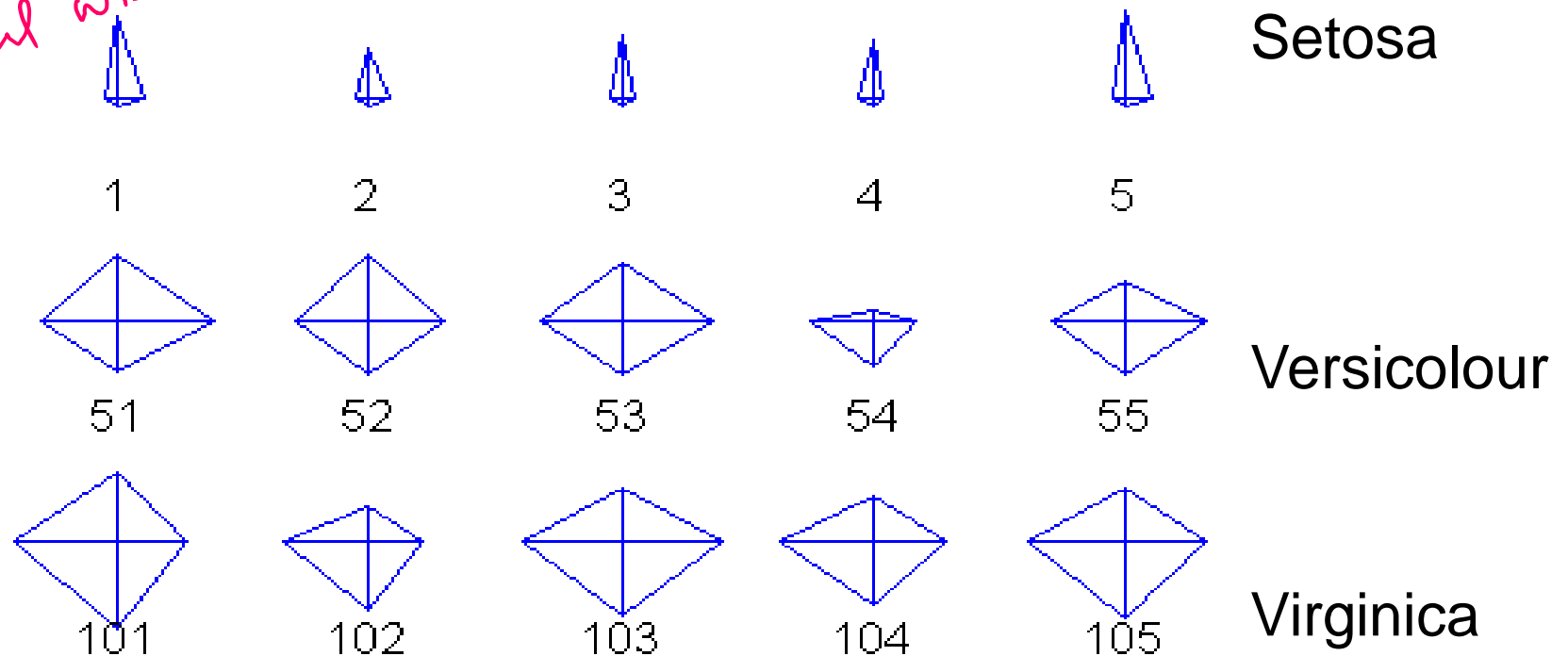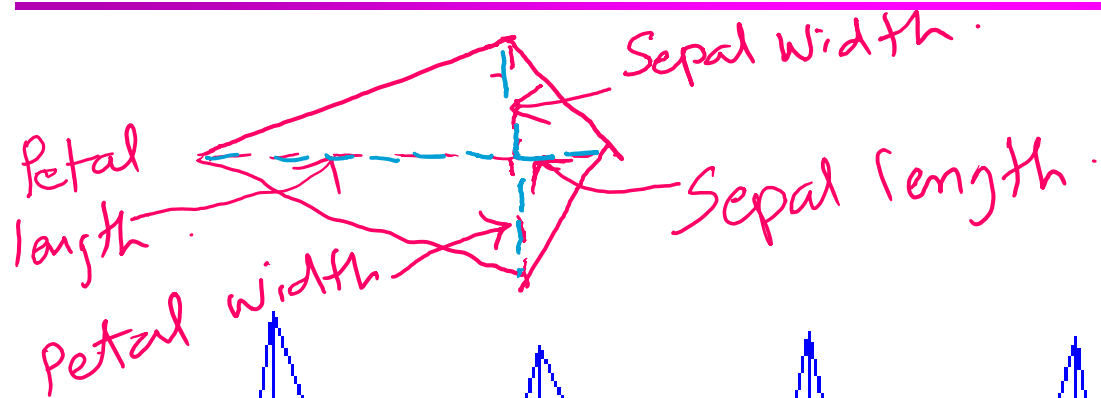# Parallel Coordinates of a Data Set

# Other Visualization Techniques

- Star Plots
  - Similar approach to parallel coordinates, but axes radiate from a central point
  - The line connecting the values of an object is a polygon

- Chernoff Faces
  - Approach created by Herman Chernoff
  - This approach associates each attribute with a characteristic of a face
  - The values of each attribute determine the appearance of the corresponding facial characteristic
  - Each object becomes a separate face
  - Relies on human's ability to distinguish faces

# Star Plots for Iris Data

# Chernoff Faces for Iris Data

Plot of 15 Iris flowers using Chernoff Faces
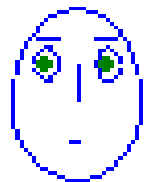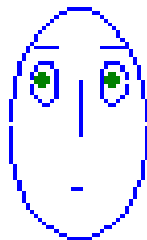


Setosa

1      2      3      4      5

Versicolour

51     52     53     54     55

Virginica

101    102    103    104    105

# 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

# 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 (Lighter = small value)

- Ex: Income vs Others : Credit limit increases with Income, Mid income people more likely to purchase more items, No correlation between Income and Age

(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

# OLAP

- On-Line Analytical Processing (OLAP) was proposed by E. F. Codd, the father of the relational database.

- Relational databases put data into tables, while OLAP uses a multidimensional array representation.

  – Such representations of data previously existed in statistics and other fields

- There are a number of data analysis and data exploration operations that are easier with such a data representation.

# Creating a Multidimensional Array

- Two key steps in converting tabular data into a multidimensional array.
  - First, identify which attributes are to be the dimensions and which attribute is to be the target attribute whose values appear as entries in the multidimensional array.
    - The attributes used as dimensions must have discrete values
    - The target value is typically a count or continuous value, e.g., the cost of an item
    - Can have no target variable at all except the count of objects that have the same set of attribute values
  - Second, find the value of each entry in the multidimensional array by summing the values (of the target attribute) or count of all objects that have the attribute values corresponding to that entry.

# Example: Iris data

- We show how the attributes, petal length, petal width, and species type can be converted to a multidimensional array
  - First, we discretized the petal width and length to have categorical values: *low*, *medium*, and *high*
  - We get the following table - note the count attribute

| Petal Length | Petal Width | Species Type | Count |
|---|---|---|---|
| low | low | Setosa | 46 |
| low | medium | Setosa | 2 |
| medium | low | Setosa | 2 |
| medium | medium | Versicolour | 43 |
| medium | high | Versicolour | 3 |
| medium | high | Virginica | 3 |
| high | medium | Versicolour | 2 |
| high | medium | Virginica | 3 |
| high | high | Versicolour | 2 |
| high | high | Virginica | 44 |

# Example: Iris data (continued)

- Each unique tuple of petal width, petal length, and species type identifies one element of the array.

- This element is assigned the corresponding count value.

- The figure illustrates the result.

- All non-specified tuples are 0.

# Example: Iris data (continued)

☐ Slices of the multidimensional array are shown by the following cross-tabulations

☐ What do these tables tell us?

|        | Width |        |      |
|--------|-------|--------|------|
| **Length** | low | medium | high |
| low    | 46    | 2      | 0    |
| medium | 2     | 0      | 0    |
| high   | 0     | 0      | 0    |

*Setosa*

|        | Width |        |      |
|--------|-------|--------|------|
| **Length** | low | medium | high |
| low    | 0     | 0      | 0    |
| medium | 0     | 43     | 3    |
| high   | 0     | 2      | 2    |

*Versicolour*

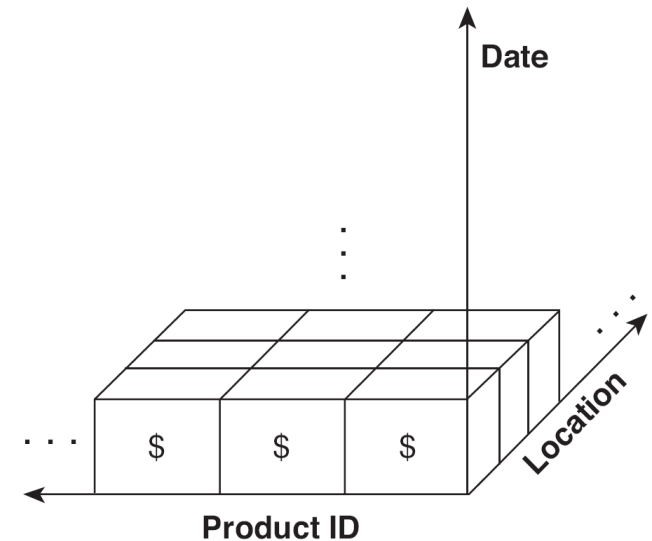|        | Width |        |      |
|--------|-------|--------|------|
| **Length** | low | medium | high |
| low    | 0     | 0      | 0    |
| medium | 0     | 0      | 3    |
| high   | 0     | 3      | 44   |

*Virginica*

# OLAP Operations: Data Cube

- The key operation of a OLAP is the <span style="color:red">formation of a data cube</span>

- A data cube is a multidimensional representation of data, together with all possible aggregates.

- By all possible aggregates, we mean the aggregates that result by selecting a proper subset of the dimensions and summing over all remaining dimensions.

- For example, if we choose the species type dimension of the Iris data and sum over all other dimensions, the result will be a one-dimensional entry with three entries, each of which gives the number of flowers of each type.

# Data Cube Example

- Consider a data set that records the sales of products at a number of company stores at various dates.

- This data can be represented as a 3 dimensional array

- There are 3 two-dimensional aggregates (3 choose 2 ), 3 one-dimensional aggregates, and 1 zero-dimensional aggregate (the overall total)

# Data Cube Example (continued)

☐ The following figure table shows one of the two dimensional aggregates, along with two of the one-dimensional aggregates, and the overall total

|  | date | | | | |
|---|---|---|---|---|---|
| product ID | Jan 1, 2004 | Jan 2, 2004 | ... | Dec 31, 2004 | total |
| 1 | $1,001 | $987 | ... | $891 | $370,000 |
| ⋮ | ⋮ | | | ⋮ | ⋮ |
| 27 | $10,265 | $10,225 | ... | $9,325 | $3,800,020 |
| ⋮ | ⋮ | | | ⋮ | ⋮ |
| total | $527,362 | $532,953 | ... | $631,221 | $227,352,127 |

# OLAP Operations: Slicing and Dicing

- Slicing is selecting a group of cells from the entire multidimensional array by specifying a specific value for one or more dimensions.

- Dicing involves selecting a subset of cells by specifying a range of attribute values.

    – This is equivalent to defining a subarray from the complete array.

- In practice, both operations can also be accompanied by aggregation over some dimensions.

# OLAP Operations: Roll-up and Drill-down

- Attribute values often have a hierarchical structure.

  - Each date is associated with a year, month, and week.

  - A location is associated with a continent, country, state (province, etc.), and city.

  - Products can be divided into various categories, such as clothing, electronics, and furniture.

- Note that these categories often nest and form a tree or lattice

  - A year contains months which contains day

  - A country contains a state which contains a city

# OLAP Operations: Roll-up and Drill-down

☐ This hierarchical structure gives rise to the roll-up and drill-down operations.

– For sales data, we can aggregate (roll up) the sales across all the dates in a month.

– Conversely, given a view of the data where the time dimension is broken into months, we could split the monthly sales totals (drill down) into daily sales totals.

– Likewise, we can drill down or roll up on the location or product ID attributes.