

# Exploring Data

Instructor: Dr.Eng. Rahmat Widyanto

Dr.Eng. Rahmat Widyanto

1

## 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/SEMATECH e-Handbook of Statistical Methods*  
<http://www.itl.nist.gov/div898/handbook/index.htm>

Dr.Eng. Rahmat Widyanto

2

## 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)

Dr.Eng. Rahmat Widyanto

3

## 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.

Dr.Eng. Rahmat Widyanto

4

### 3.1 The Iris Data Set

In the following discussion, we will often refer to the Iris data set that is available from the University of California at Irvine (UCI) Machine Learning Repository. It consists of information on 150 Iris flowers, 50 each from one of three Iris species: Setosa, Versicolour, and Virginica. Each flower is characterized by five attributes:

1. sepal length in centimeters
2. sepal width in centimeters
3. petal length in centimeters
4. petal width in centimeters
5. class (Setosa, Versicolour, Virginica)

The sepals of a flower are the outer structures that protect the more fragile parts of the flower, such as the petals. In many flowers, the sepals are green, and only the petals are colorful. For Irises, however, the sepals are also colorful. As illustrated by the picture of a Virginica Iris in Figure 3.1, the sepals of an Iris are larger than the petals and are drooping, while the petals are upright.

### Iris Flower



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.

```
43 44 44 44 45 46 46 46 46 47 47 47 48 48 48 48 49 49 49 49 49 49 49 49 50  
50 50 50 50 50 50 50 50 51 51 51 51 51 51 51 51 51 51 52 52 52 52 53  
54 54 54 54 54 54 55 55 55 55 55 55 55 56 56 56 56 56 56 57 57 57 57  
57 57 57 57 58 58 58 58 58 58 58 59 59 59 60 60 60 60 60 60 61 61 61  
61 61 61 62 62 62 62 63 63 63 63 63 63 63 64 64 64 64 64 64 64  
65 65 65 65 65 66 66 67 67 67 67 67 67 68 68 68 69 69 69 69 70  
71 72 72 72 73 74 76 77 77 77 77 77 79
```

**Figure 3.4.** Sepal length data from the Iris data set.

```
4 : 34444566667788888999999  
5 : 00000000001111111122223444445555556666667777777888888999  
6 : 0000011111222233333333444444555556677777778889999  
7 : 0122234677779
```

**Figure 3.5.** Stem and leaf plot for the sepal length from the Iris data set.

```

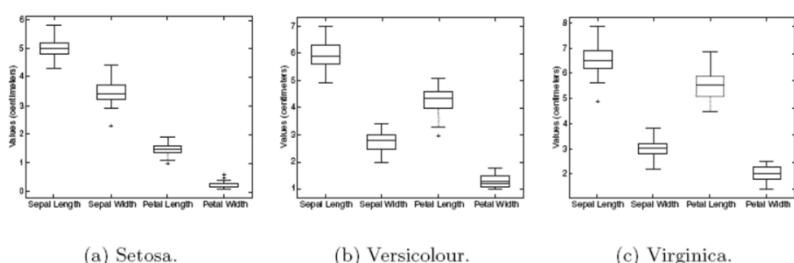
4 : 3444
4 : 566667788888999999
5 : 000000000111111112222344444
5 : 55555566666677777888888999
6 : 0000011111222233333333444444
6 : 5555667777778889999
7 : 0122234
7 : 677779

```

**Figure 3.6.** Stem and leaf plot for the sepal length from the Iris data set when buckets corresponding to digits are split.

Dr.Eng. Rahmat Widyanto

9



(a) Setosa.

(b) Versicolour.

(c) Virginica.

**Figure 3.12.** Box plots of attributes by Iris species.

Dr.Eng. Rahmat Widyanto

10

## 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

Dr.Eng. Rahmat Widyanto

11

### 3.2.1 Frequencies and the Mode

Given a set of unordered categorical values, there is not much that can be done to further characterize the values except to compute the frequency with which each value occurs for a particular set of data. Given a categorical attribute  $x$ , which can take values  $\{v_1, \dots, v_i, \dots, v_k\}$  and a set of  $m$  objects, the frequency of a value  $v_i$  is defined as

$$\text{frequency}(v_i) = \frac{\text{number of objects with attribute value } v_i}{m}. \quad (3.1)$$

The **mode** of a categorical attribute is the value that has the highest frequency.

Dr.Eng. Rahmat Widyanto

12

## 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

Dr.Eng. Rahmat Widyanto

13

**Example 3.1.** Consider a set of students who have an attribute, *class*, which can take values from the set  $\{\text{freshman}, \text{sophomore}, \text{junior}, \text{senior}\}$ . Table 3.1 shows the number of students for each value of the *class* attribute. The mode of the *class* attribute is *freshman*, with a frequency of 0.33. This may indicate dropouts due to attrition or a larger than usual freshman class.

**Table 3.1.** Class size for students in a hypothetical college.

Class	Size	Frequency
freshman	200	0.33
sophomore	160	0.27
junior	130	0.22
senior	110	0.18

Dr.Eng. Rahmat Widyanto

14

Categorical attributes often, but not always, have a small number of values, and consequently, the mode and frequencies of these values can be interesting and useful. Notice, though, that for the Iris data set and the *class* attribute, the three types of flower all have the same frequency, and therefore, the notion of a mode is not interesting.

For continuous data, the mode, as currently defined, is often not useful because a single value may not occur more than once. Nonetheless, in some cases, the mode may indicate important information about the nature of the values or the presence of missing values. For example, the heights of 20 people measured to the nearest millimeter will typically not repeat, but if the heights are measured to the nearest tenth of a meter, then some people may have the same height. Also, if a unique value is used to indicate a missing value, then this value will often show up as the mode.

## 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 is a value  $x_p$  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  $x_{50\%}$ .

**Table 3.2.** Percentiles for sepal length, sepal width, petal length, and petal width. (All values are in centimeters.)

Percentile	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	2.0	1.0	0.1
10	4.8	2.5	1.4	0.2
20	5.0	2.7	1.5	0.2
30	5.2	2.8	1.7	0.4
40	5.6	3.0	3.9	1.2
50	5.8	3.0	4.4	1.3
60	6.1	3.1	4.6	1.5
70	6.3	3.2	5.0	1.8
80	6.6	3.4	5.4	1.9
90	6.9	3.6	5.8	2.2
100	7.9	4.4	6.9	2.5

**Example 3.2.** The percentiles,  $x_{0\%}, x_{10\%}, \dots, x_{90\%}, x_{100\%}$  of the integers from 1 to 10 are, in order, the following: 1.0, 1.5, 2.5, 3.5, 4.5, 5.5, 6.5, 7.5, 8.5, 9.5, 10.0. By tradition,  $x_{0\%} = \min(x)$  and  $x_{100\%} = \max(x)$ . ■

## 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.

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

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

Although the mean is sometimes interpreted as the middle of a set of values, this is only correct if the values are distributed in a symmetric manner. If the distribution of values is skewed, then the median is a better indicator of the middle. Also, the mean is sensitive to the presence of outliers. For data with outliers, the median again provides a more robust estimate of the middle of a set of values.

To overcome problems with the traditional definition of a mean, the notion of a **trimmed mean** is sometimes used. A percentage  $p$  between 0 and 100 is specified, the top and bottom  $(p/2)\%$  of the data is thrown out, and the mean is then calculated in the normal way. The median is a trimmed mean with  $p = 100\%$ , while the standard mean corresponds to  $p = 0\%$ .

**Example 3.3.** Consider the set of values  $\{1, 2, 3, 4, 5, 90\}$ . The mean of these values is 17.5, while the median is 3.5. The trimmed mean with  $p = 40\%$  is also 3.5. ■

**Example 3.4.** The means, medians, and trimmed means ( $p = 20\%$ ) of the four quantitative attributes of the Iris data are given in Table 3.3. The three measures of location have similar values except for the attribute *petal length*.

**Table 3.3.** Means and medians for sepal length, sepal width, petal length, and petal width. (All values are in centimeters.)

Measure	Sepal Length	Sepal Width	Petal Length	Petal Width
mean	5.84	3.05	3.76	1.20
median	5.80	3.00	4.35	1.30
trimmed mean (20%)	5.79	3.02	3.72	1.12

## Measures of Spread: Range and Variance

- Range is the difference between the max and min

$$\text{range}(x) = \max(x) - \min(x) = x_{(m)} - x_{(1)}.$$

- The variance or standard deviation  $s_x^2$  is the most common measure of the spread of a set of points.

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

- Because of outliers, other measures are often used.

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

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

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

Dr.Eng. Rahmat Widyanto

21

**Table 3.4.** Range, standard deviation (std), absolute average difference (AAD), median absolute difference (MAD), and interquartile range (IQR) for sepal length, sepal width, petal length, and petal width. (All values are in centimeters.)

Measure	Sepal Length	Sepal Width	Petal Length	Petal Width
range	3.6	2.4	5.9	2.4
std	0.8	0.4	1.8	0.8
AAD	0.7	0.3	1.6	0.6
MAD	0.7	0.3	1.2	0.7
IQR	1.3	0.5	3.5	1.5

Although the range identifies the maximum spread, it can be misleading if most of the values are concentrated in a narrow band of values, but there are also a relatively small number of more extreme values. Hence, the **variance** is preferred as a measure of spread. The variance of the (observed) values of an attribute  $x$  is typically written as  $s_x^2$  and is defined below. The **standard deviation**, which is the square root of the variance, is written as  $s_x$  and has the same units as  $x$ .

Dr.Eng. Rahmat Widyanto

22

### 3.2.5 Multivariate Summary Statistics

Measures of location for data that consists of several attributes (multivariate data) can be obtained by computing the mean or median separately for each attribute. Thus, given a data set the mean of the data objects,  $\bar{\mathbf{x}}$ , is given by

$$\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_n), \quad (3.9)$$

where  $\bar{x}_i$  is the mean of the  $i^{th}$  attribute  $x_i$ .

For multivariate data, the spread of each attribute can be computed independently of the other attributes using any of the approaches described in Section 3.2.4. However, for data with continuous variables, the spread of the data is most commonly captured by the **covariance matrix  $\mathbf{S}$** , whose  $ij^{th}$  entry  $s_{ij}$  is the covariance of the  $i^{th}$  and  $j^{th}$  attributes of the data. Thus, if  $x_i$  and  $x_j$  are the  $i^{th}$  and  $j^{th}$  attributes, then

$$s_{ij} = \text{covariance}(x_i, x_j). \quad (3.10)$$

In turn,  $\text{covariance}(x_i, x_j)$  is given by

$$\text{covariance}(x_i, x_j) = \frac{1}{m-1} \sum_{k=1}^m (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j), \quad (3.11)$$

where  $x_{ki}$  and  $x_{kj}$  are the values of the  $i^{th}$  and  $j^{th}$  attributes for the  $k^{th}$  object. Notice that  $\text{covariance}(x_i, x_i) = \text{variance}(x_i)$ . Thus, the covariance matrix has the variances of the attributes along the diagonal.

The covariance of two attributes is a measure of the degree to which two attributes vary together and depends on the magnitudes of the variables. A value near 0 indicates that two attributes do not have a (linear) relationship, but it is not possible to judge the degree of relationship between two variables by looking only at the value of the covariance. Because the correlation of two attributes immediately gives an indication of how strongly two attributes are (linearly) related, correlation is preferred to covariance for data exploration. (Also see the discussion of correlation in Section 2.4.5.) The  $ij^{th}$  entry of the **correlation matrix  $\mathbf{R}$** , is the correlation between the  $i^{th}$  and  $j^{th}$  attributes of the data. If  $x_i$  and  $x_j$  are the  $i^{th}$  and  $j^{th}$  attributes, then

$$r_{ij} = \text{correlation}(x_i, x_j) = \frac{\text{covariance}(x_i, x_j)}{s_i s_j}, \quad (3.12)$$

where  $s_i$  and  $s_j$  are the variances of  $x_i$  and  $x_j$ , respectively. The diagonal entries of  $\mathbf{R}$  are  $\text{correlation}(x_i, x_i) = 1$ , while the other entries are between  $-1$  and  $1$ . It is also useful to consider correlation matrices that contain the pairwise correlations of objects instead of attributes.

### 3.2.6 Other Ways to Summarize the Data

There are, of course, other types of summary statistics. For instance, the **skewness** of a set of values measures the degree to which the values are symmetrically distributed around the mean. There are also other characteristics of the data that are not easy to measure quantitatively, such as whether the distribution of values is multimodal; i.e., the data has multiple “bumps” where most of the values are concentrated. In many cases, however, the most effective approach to understanding the more complicated or subtle aspects of how the values of an attribute are distributed, is to view the values graphically in the form of a histogram. (Histograms are discussed in the next section.)

## 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.

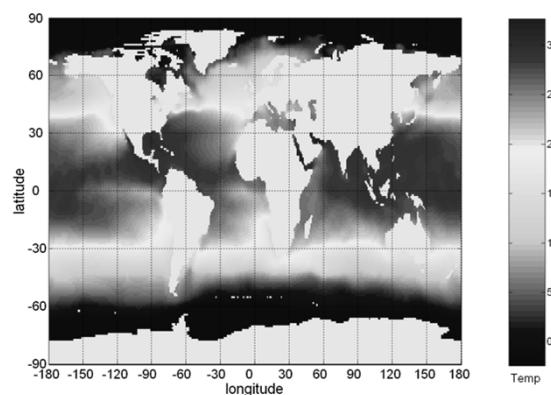
- Visualization of data is one of the most powerful and appealing techniques for data exploration.
  - Humans have a well developed ability to analyze large amounts of information that is presented visually
  - Can detect general patterns and trends
  - Can detect outliers and unusual patterns

Dr.Eng. Rahmat Widyanto

27

## Example: Sea Surface Temperature

- The following shows the Sea Surface Temperature (SST) for July 1982
  - Thousands of data points are summarized in a single figure



Dr.Eng. Rahmat Widyanto

28

## Representation

- Is the mapping of information to a visual format
- Data objects, their attributes, and the relationships among data objects are translated into graphical elements such as points, lines, shapes, and colors.
- 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.

Dr.Eng. Rahmat Widyanto

29

## Arrangement

As discussed earlier, the proper choice of visual representation of objects and attributes is essential for good visualization. The arrangement of items within the visual display is also crucial. We illustrate this with two examples.

**Example 3.5.** This example illustrates the importance of rearranging a table of data. In Table 3.5, which shows nine objects with six binary attributes, there is no clear relationship between objects and attributes, at least at first glance. If the rows and columns of this table are permuted, however, as shown in Table 3.6, then it is clear that there are really only two types of objects in the table—one that has all ones for the first three attributes and one that has only ones for the last three attributes. ■

Dr.Eng. Rahmat Widyanto

30

## 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:

	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

	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

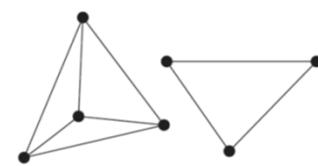
Dr.Eng. Rahmat Widyanto

31

**Example 3.6.** Consider Figure 3.3(a), which shows a visualization of a graph. If the connected components of the graph are separated, as in Figure 3.3(b), then the relationships between nodes and graphs become much simpler to understand. ■



(a) Original view of a graph.



(b) Uncoupled view of connected components of the graph.

**Figure 3.3.** Two visualizations of a graph.

Dr.Eng. Rahmat Widyanto

32

## Selection

- Is the elimination or the de-emphasis of certain objects and attributes
- Selection may involve choosing a subset 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

Dr.Eng. Rahmat Widyanto

33

**Stem and Leaf Plots** Stem and leaf plots can be used to provide insight into the distribution of one-dimensional integer or continuous data. (We will assume integer data initially, and then explain how stem and leaf plots can be applied to continuous data.) For the simplest type of stem and leaf plot, we split the values into groups, where each group contains those values that are the same except for the last digit. Each group becomes a stem, while the last digits of a group are the leaves. Hence, if the values are two-digit integers, e.g., 35, 36, 42, and 51, then the stems will be the high-order digits, e.g., 3, 4, and 5, while the leaves are the low-order digits, e.g., 1, 2, 5, and 6. By plotting the stems vertically and leaves horizontally, we can provide a visual representation of the distribution of the data.

Dr.Eng. Rahmat Widyanto

34

**Example 3.7.** The set of integers shown in Figure 3.4 is the sepal length in centimeters (multiplied by 10 to make the values integers) taken from the Iris data set. For convenience, the values have also been sorted.

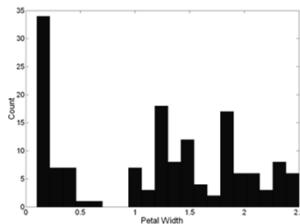
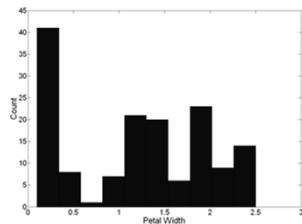
The stem and leaf plot for this data is shown in Figure 3.5. Each number in Figure 3.4 is first put into one of the vertical groups—4, 5, 6, or 7—according to its ten's digit. Its last digit is then placed to the right of the colon. Often, especially if the amount of data is larger, it is desirable to split the stems. For example, instead of placing all values whose ten's digit is 4 in the same “bucket,” the stem 4 is repeated twice; all values 40–44 are put in the bucket corresponding to the first stem and all values 45–49 are put in the bucket corresponding to the second stem. This approach is shown in the stem and leaf plot of Figure 3.6. Other variations are also possible. ■

## 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)**



```
43 44 44 44 45 46 46 46 46 47 47 48 48 48 48 48 49 49 49 49 49 49 49 50  
50 50 50 50 50 50 50 50 51 51 51 51 51 51 51 51 51 51 52 52 52 53  
54 54 54 54 54 55 55 55 55 55 55 55 56 56 56 56 56 56 57 57 57 57  
57 57 57 57 58 58 58 58 58 58 59 59 59 60 60 60 60 60 61 61 61  
61 61 61 62 62 62 62 63 63 63 63 63 63 63 63 64 64 64 64 64 64 64  
65 65 65 65 66 66 67 67 67 67 67 67 67 68 68 68 69 69 69 69 70  
71 72 72 72 73 74 76 77 77 77 77 79
```

**Figure 3.4.** Sepal length data from the Iris data set.

```
4 : 34444566667788888999999  
5 : 0000000001111111122223444445555556666667777777888888999  
6 : 00000111112222333333333444444555556677777778889999  
7 : 0122234677779
```

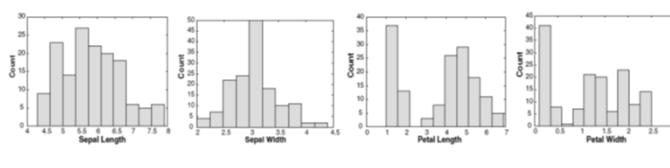
**Figure 3.5.** Stem and leaf plot for the sepal length from the Iris data set.

```
4 : 3444  
4 : 56666778888899999  
5 : 0000000001111111122223444444  
5 : 55555566666777777888888999  
6 : 000001111122223333333334444444  
6 : 555556677777778889999  
7 : 0122234  
7 : 677779
```

**Figure 3.6.** Stem and leaf plot for the sepal length from the Iris data set when buckets corresponding to digits are split.

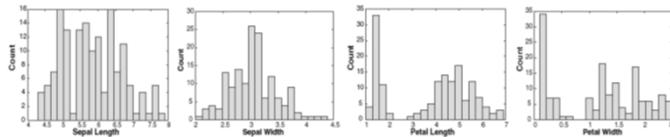
**Example 3.8.** Figure 3.7 shows histograms (with 10 bins) for sepal length, sepal width, petal length, and petal width. Since the shape of a histogram can depend on the number of bins, histograms for the same data, but with 20 bins, are shown in Figure 3.8.

There are variations of the histogram plot. A **relative (frequency) histogram** replaces the count by the relative frequency. However, this is just a change in scale of the  $y$  axis, and the shape of the histogram does not change. Another common variation, especially for unordered categorical data, is the **Pareto histogram**, which is the same as a normal histogram except that the categories are sorted by count so that the count is decreasing from left to right.



(a) Sepal length. (b) Sepal width. (c) Petal length. (d) Petal width.

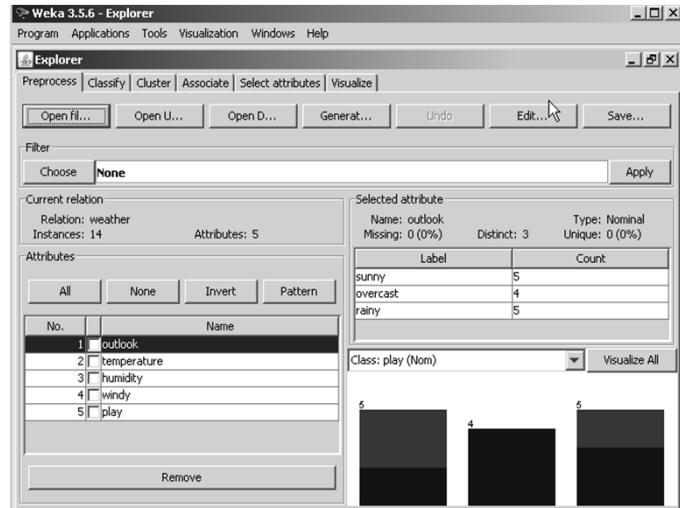
Figure 3.7. Histograms of four Iris attributes (10 bins).



(a) Sepal length. (b) Sepal width. (c) Petal length. (d) Petal width.

Figure 3.8. Histograms of four Iris attributes (20 bins).

## Histogram from Weka

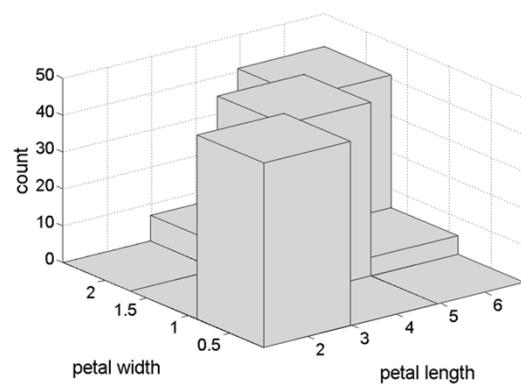


Dr.Eng. Rahmat Widhyanto

41

## Two-Dimensional Histograms

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



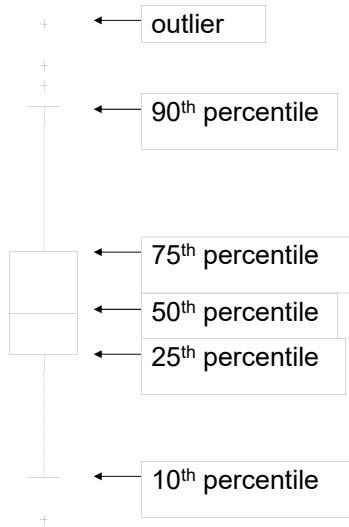
Dr.Eng. Rahmat Widhyanto

42

## 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

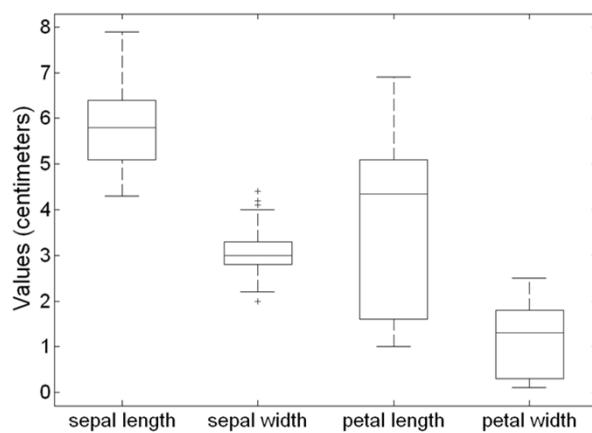


Dr.Eng. Rahmat Widyanto

43

## Example of Box Plots

- Box plots can be used to compare attributes

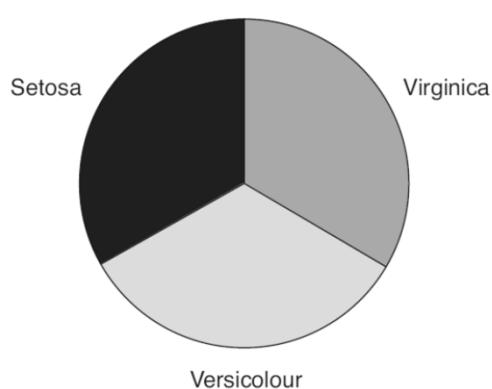


Dr.Eng. Rahmat Widyanto

44

**Pie Chart** A pie chart is similar to a histogram, but is typically used with categorical attributes that have a relatively small number of values. Instead of showing the relative frequency of different values with the area or height of a bar, as in a histogram, a pie chart uses the relative area of a circle to indicate relative frequency. Although pie charts are common in popular articles, they are used less frequently in technical publications because the size of relative areas can be hard to judge. Histograms are preferred for technical work.

**Example 3.11.** Figure 3.13 displays a pie chart that shows the distribution of Iris species in the Iris data set. In this case, all three flower types have the same frequency. ■



**Figure 3.13.** Distribution of the types of Iris flowers.

### Percentile Plots and Empirical Cumulative Distribution Functions

A type of diagram that shows the distribution of the data more quantitatively is the plot of an empirical cumulative distribution function. While this type of plot may sound complicated, the concept is straightforward. For each value of a statistical distribution, a **cumulative distribution function** (CDF) shows the probability that a point is less than that value. For each observed value, an **empirical cumulative distribution function** (ECDF) shows the fraction of points that are less than this value. Since the number of points is finite, the empirical cumulative distribution function is a step function.

**Example 3.12.** Figure 3.14 shows the ECDFs of the Iris attributes. The percentiles of an attribute provide similar information. Figure 3.15 shows the **percentile plots** of the four continuous attributes of the Iris data set from Table 3.2. The reader should compare these figures with the histograms given in Figures 3.7 and 3.8. ■

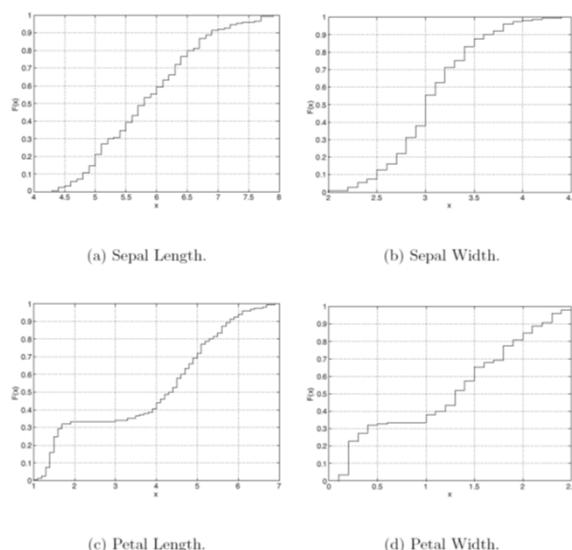


Figure 3.14. Empirical CDFs of four Iris attributes.

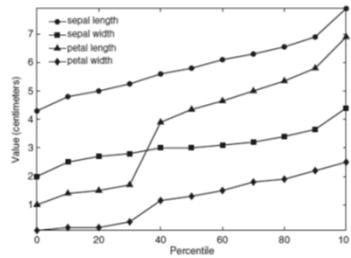


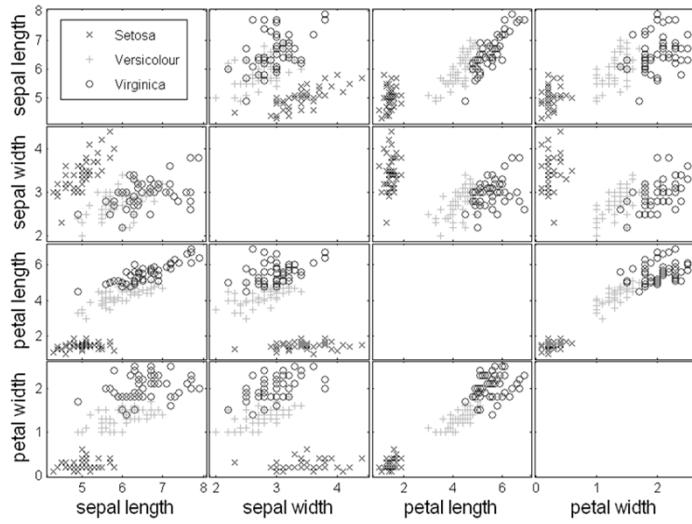
Figure 3.15. Percentile plots for sepal length, sepal width, petal length, and petal width.

## 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



Dr.Eng. Rahmat Widyanto

51

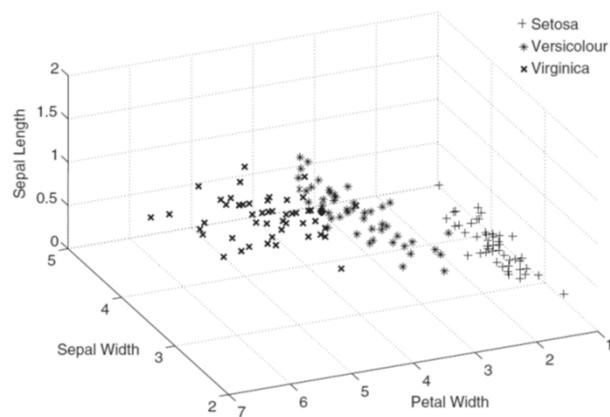
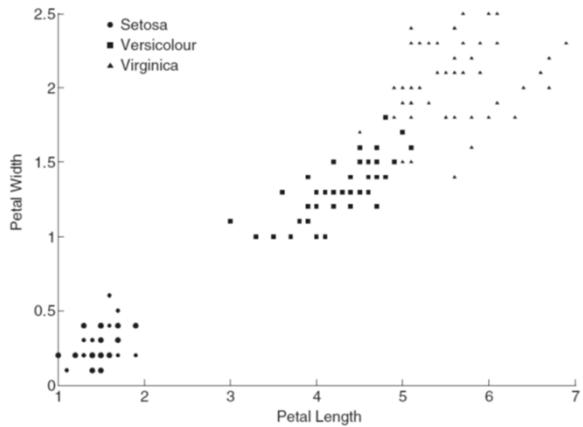


Figure 3.17. Three-dimensional scatter plot of sepal width, sepal length, and petal width.

Dr.Eng. Rahmat Widyanto

52



**Figure 3.18.** Scatter plot of petal length versus petal width, with the size of the marker indicating sepal width.

Dr.Eng. Rahmat Widyanto

53

## Visualization Techniques: Contour Plots

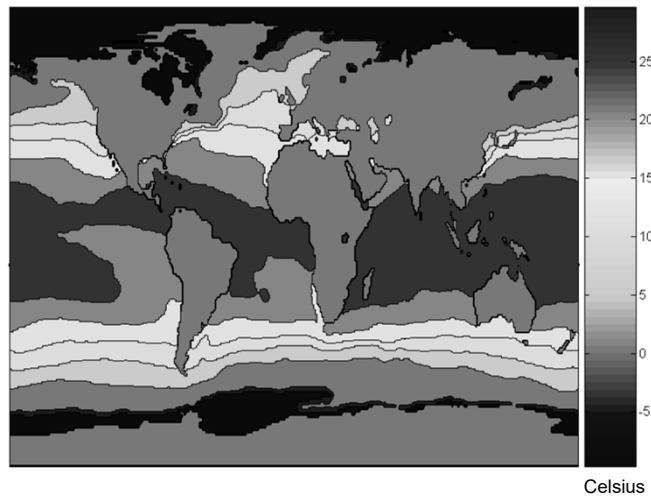
- Contour plots

- 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

Dr.Eng. Rahmat Widyanto

54

## Contour Plot Example: SST Dec, 1998



Dr.Eng. Rahmat Widyanto

55

**Surface Plots** Like contour plots, **surface plots** use two attributes for the  $x$  and  $y$  coordinates. The third attribute is used to indicate the height above the plane defined by the first two attributes. While such graphs can be useful, they require that a value of the third attribute be defined for all combinations of values for the first two attributes, at least over some range. Also, if the surface is too irregular, then it can be difficult to see all the information, unless the plot is viewed interactively. Thus, surface plots are often used to describe mathematical functions or physical surfaces that vary in a relatively smooth manner.

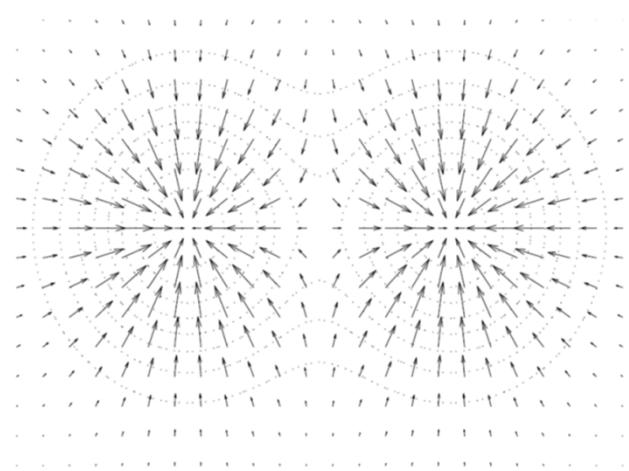
**Example 3.16.** Figure 3.20 shows a surface plot of the density around a set of 12 points. This example is further discussed in Section 9.3.3. ■

Dr.Eng. Rahmat Widyanto

56

**Vector Field Plots** In some data, a characteristic may have both a magnitude and a direction associated with it. For example, consider the flow of a substance or the change of density with location. In these situations, it can be useful to have a plot that displays both direction and magnitude. This type of plot is known as a **vector plot**.

**Example 3.17.** Figure 3.21 shows a contour plot of the density of the two smaller density peaks from Figure 3.20(b), annotated with the density gradient vectors. ■



**Figure 3.21.** Vector plot of the gradient (change) in density for the bottom two density peaks of Figure 3.20.

**Lower-Dimensional Slices** Consider a spatio-temporal data set that records some quantity, such as temperature or pressure, at various locations over time. Such a data set has four dimensions and cannot be easily displayed by the types of plots that we have described so far. However, separate “slices” of the data can be displayed by showing a set of plots, one for each month. By examining the change in a particular area from one month to another, it is possible to notice changes that occur, including those that may be due to seasonal factors.

**Example 3.18.** The underlying data set for this example consists of the average monthly sea level pressure (SLP) from 1982 to 1999 on a  $2.5^\circ$  by  $2.5^\circ$  latitude-longitude grid. The twelve monthly plots of pressure for one year are shown in Figure 3.22. In this example, we are interested in slices for a particular month in the year 1982. More generally, we can consider slices of the data along any arbitrary dimension. ■

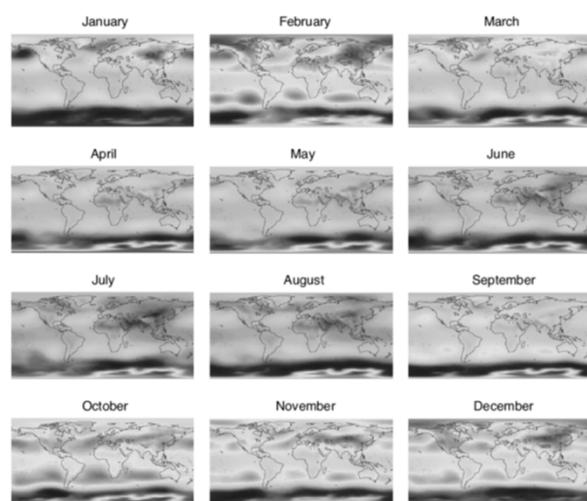


Figure 3.22. Monthly plots of sea level pressure over the 12 months of 1982.

## Visualization Techniques: Matrix Plots

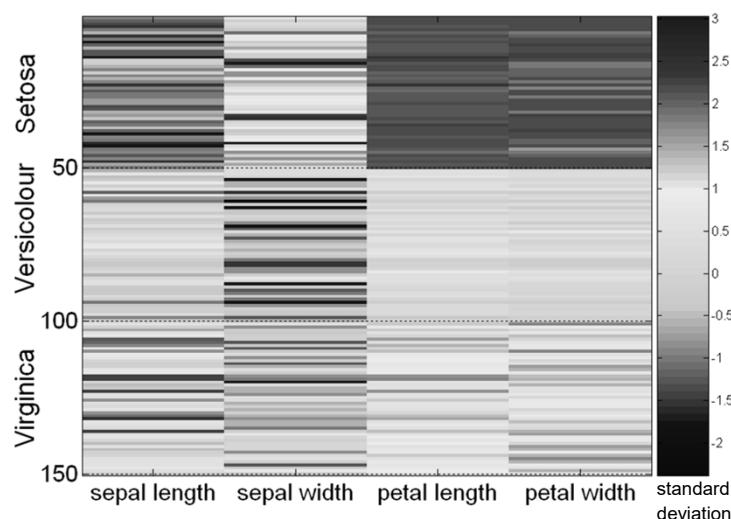
- Matrix plots

- 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

Dr.Eng. Rahmat Widyanto

61

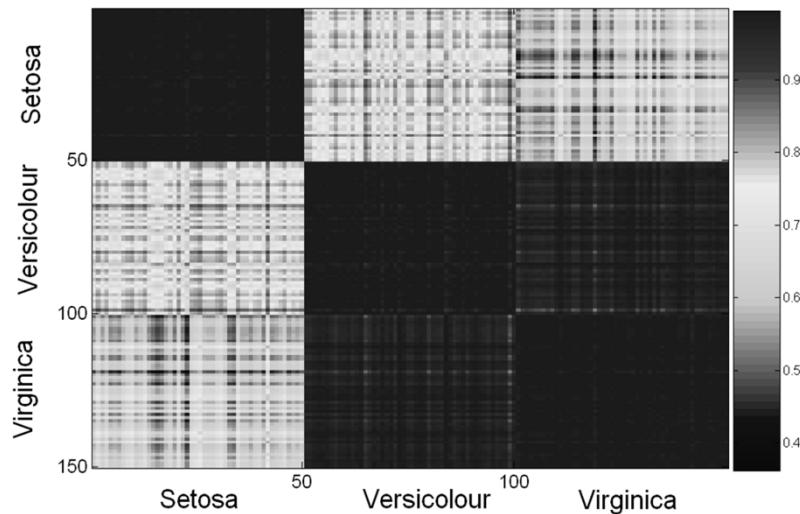
## Visualization of the Iris Data Matrix



Dr.Eng. Rahmat Widyanto

62

## Visualization of the Iris Correlation Matrix



Dr.Eng. Rahmat Widyanto

63

## Visualization Techniques: Parallel Coordinates

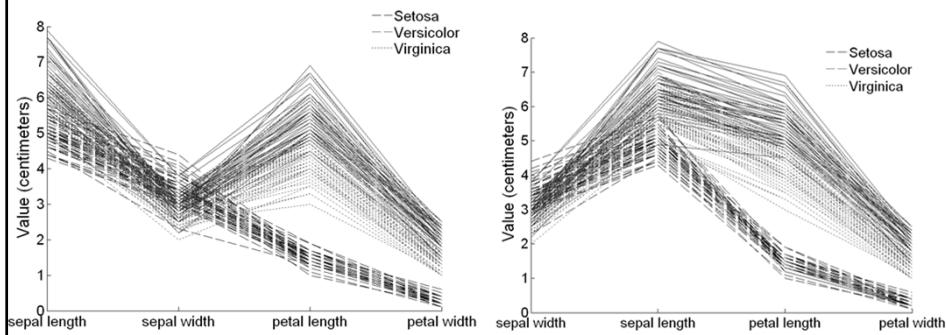
### • Parallel Coordinates

- 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

Dr.Eng. Rahmat Widyanto

64

## Parallel Coordinates Plots for Iris Data



Dr.Eng. Rahmat Widyanto

65

## 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

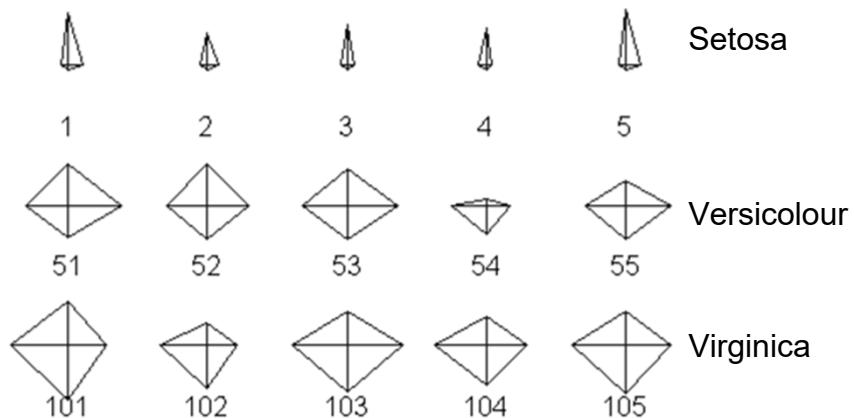
### • Chernoff Faces

- 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

Dr.Eng. Rahmat Widyanto

66

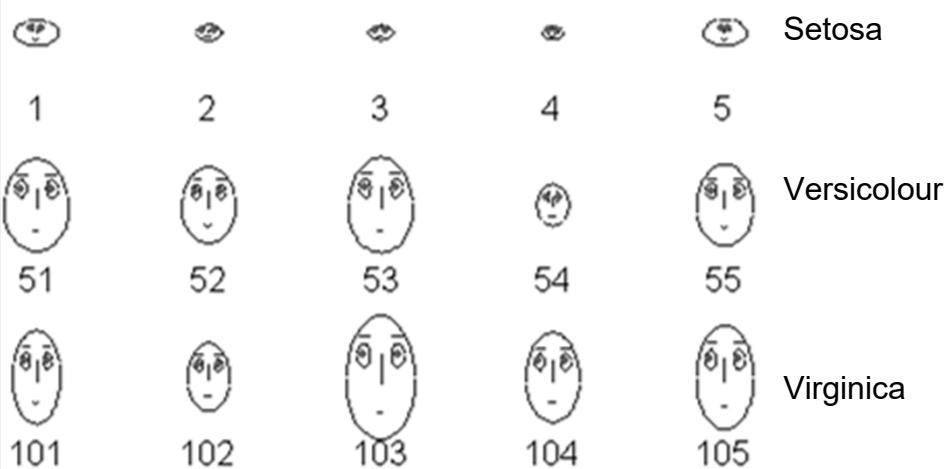
## Star Plots for Iris Data



Dr.Eng. Rahmat Widyanto

67

## Chernoff Faces for Iris Data



Dr.Eng. Rahmat Widyanto

68

### 3.3.5 Do's and Don'ts

To conclude this section on visualization, we provide a short list of visualization do's and don'ts. While these guidelines incorporate a lot of visualization wisdom, they should not be followed blindly. As always, guidelines are no substitute for thoughtful consideration of the problem at hand.

**ACCENT Principles** The following are the *ACCENT* principles for effective graphical display put forth by D. A. Burn (as adapted by Michael Friendly):

**Apprehension** Ability to correctly perceive relations among variables. Does the graph maximize apprehension of the relations among variables?

**Clarity** Ability to visually distinguish all the elements of a graph. Are the most important elements or relations visually most prominent?

**Consistency** Ability to interpret a graph based on similarity to previous graphs. Are the elements, symbol shapes, and colors consistent with their use in previous graphs?

**Efficiency** Ability to portray a possibly complex relation in as simple a way as possible. Are the elements of the graph economically used? Is the graph easy to interpret?

**Necessity** The need for the graph, and the graphical elements. Is the graph a more useful way to represent the data than alternatives (table, text)? Are all the graph elements necessary to convey the relations?

**Truthfulness** Ability to determine the true value represented by any graphical element by its magnitude relative to the implicit or explicit scale. Are the graph elements accurately positioned and scaled?

**Tufte's Guidelines** Edward R. Tufte has also enumerated the following principles for graphical excellence:

- Graphical excellence is the well-designed presentation of interesting data—a matter of *substance*, of *statistics*, and of *design*.

- Graphical excellence consists of complex ideas communicated with clarity, precision, and efficiency.
- Graphical excellence is that which gives to the viewer the greatest number of ideas in the shortest time with the least ink in the smallest space.
- Graphical excellence is nearly always multivariate.
- And graphical excellence requires telling the truth about the data.

## 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

- Converting tabular data into a multidimensional array:
  - Identify which attributes are to be the dimensions and which attribute is to be the target attribute
    - ◆ Attributes used as dimensions must have discrete values
    - ◆ Values of target variable appear as entries in the array
    - ◆ The target value is typically a count or continuous value
    - ◆ Can have no target variable at all except the count of objects that have the same set of attribute values
  - Find the value of each entry in the multidimensional array by summing the values (of the target attribute) or the count of all objects that have the attribute values corresponding to that entry.

Dr.Eng. Rahmat Widyanto

73

## 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*

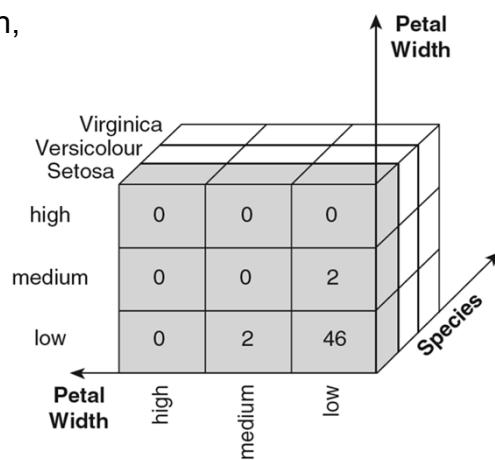
Petal Length	Petal Width	Species Type	Coun
low	low	Setosa	46
low	medium	Setosa	2
medium	low	Setosa	2
medium	medium	Versicolour	43
medium	high	Versicolour	3
medium	high	Virginica	3
high	medium	Versicolour	2
high	medium	Virginica	3
high	high	Versicolour	2
high	high	Virginica	44

Dr.Eng. Rahmat Widyanto

74

## 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.



Dr.Eng. Rahmat Widyanto

75

## Example: Iris data (continued)

- Slices of the multidimensional array are shown by the following cross-tabulations
- What do these tables tell us?

		Width		
		low	medium	high
Length	low	46	2	0
	medium	2	0	0
	high	0	0	0

		Width		
		low	medium	high
Length	low	0	0	0
	medium	0	43	3
	high	0	2	2

		Width		
		low	medium	high
Length	low	0	0	0
	medium	0	0	3
	high	0	3	44

Dr.Eng. Rahmat Widyanto

76

## OLAP Operations: Data Cube

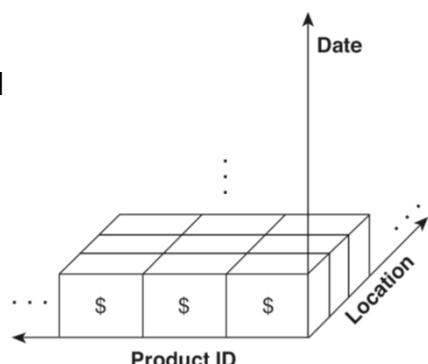
- The key operation of a OLAP is the formation of a data cube
- A data cube is a multidimensional representation of data, together with all possible aggregates.
- By all possible aggregates, we mean the aggregates that result by selecting a proper subset of the dimensions and summing over all remaining 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.

Dr.Eng. Rahmat Widyanto

77

## 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)



Dr.Eng. Rahmat Widyanto

78

## 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

product ID	date				total
	Jan 1, 2004	Jan 2, 2004	...	Dec 31, 2004	
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

Dr.Eng. Rahmat Widyanto

81

## **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.

Dr.Eng. Rahmat Widyanto

82