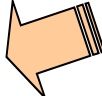


Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data 
- Data Visualization
- Measuring Data Similarity and Dissimilarity
- Summary

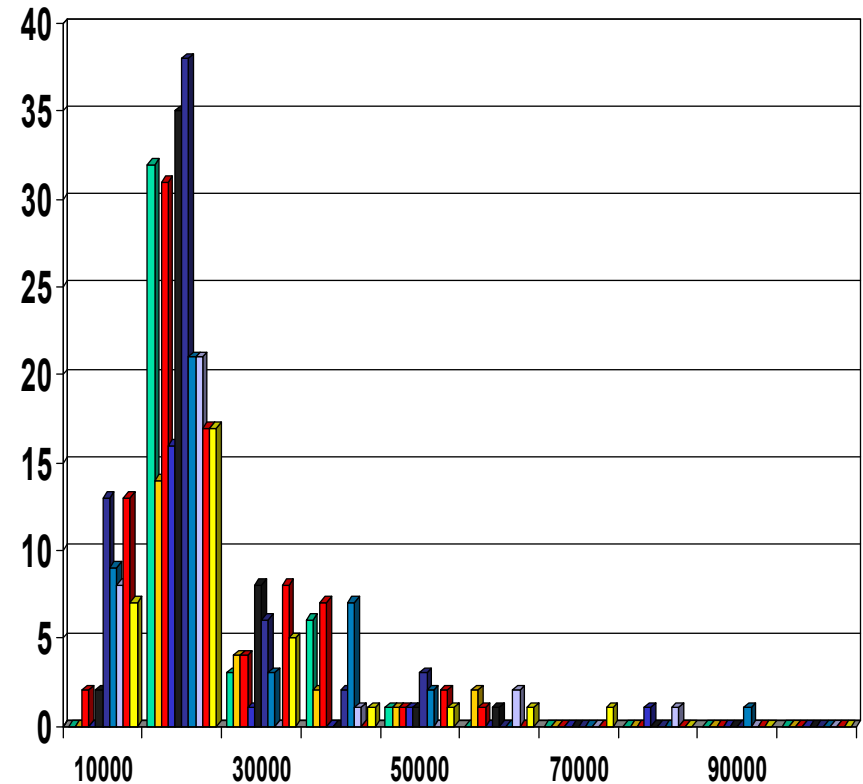


Graphic Displays of Basic Statistical Descriptions

- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis indicates values, y-axis does 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 a point in the plane

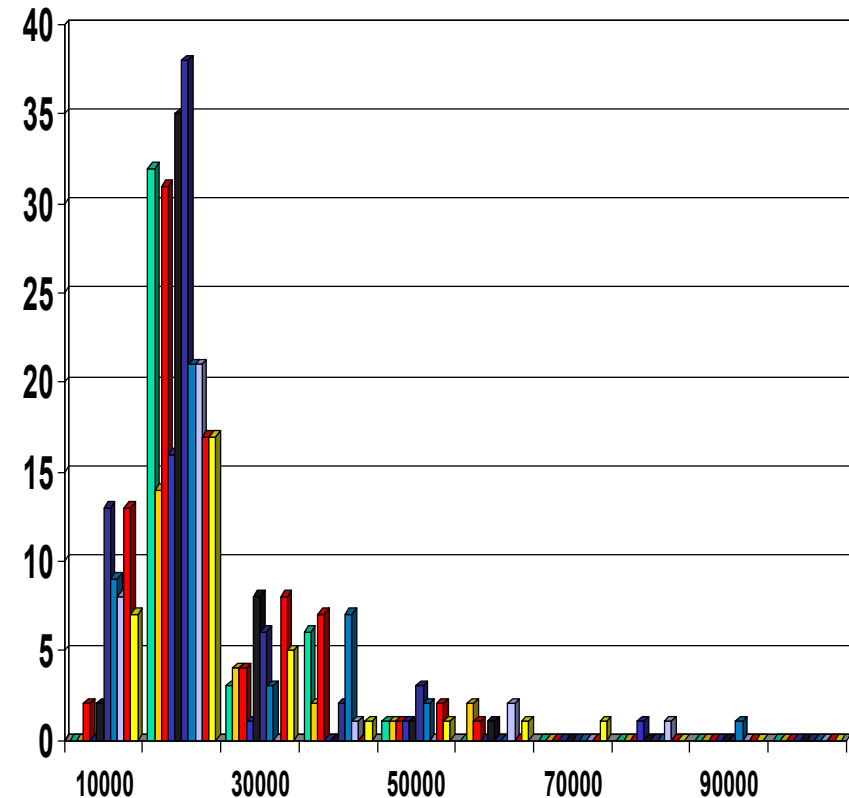
Histogram Analysis

- Histogram: Graph display of frequencies shown as bars
- It shows what proportion of cases fall into each of several categories
 - The categories are usually specified as non-overlapping intervals of some variable
 - The categories (bars) must be adjacent

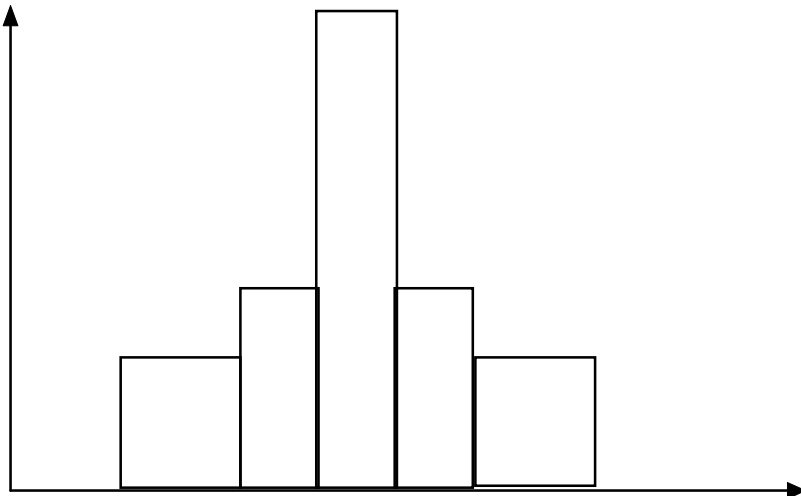
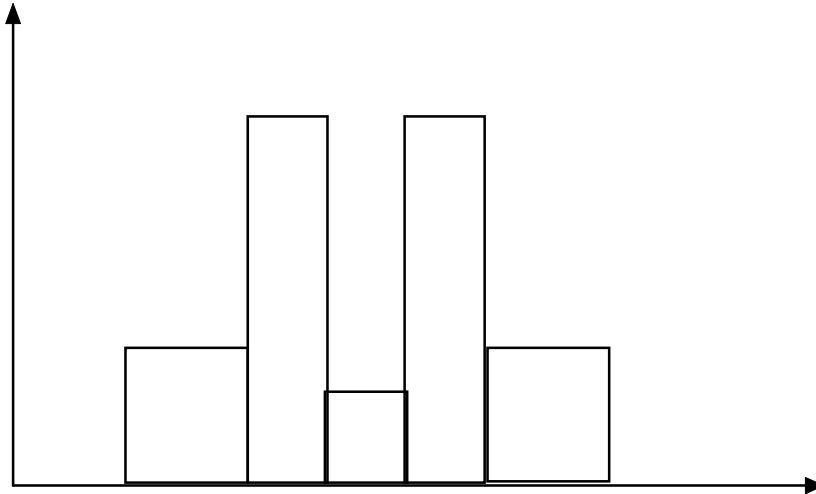


Histogram Analysis

- Differs from a bar chart
 - The *area* of the bar denotes the value (histogram)
 - The *height* denotes the value (bar chart)
 - A crucial distinction when the categories are not of uniform width



Histograms Often Tell More than Boxplots

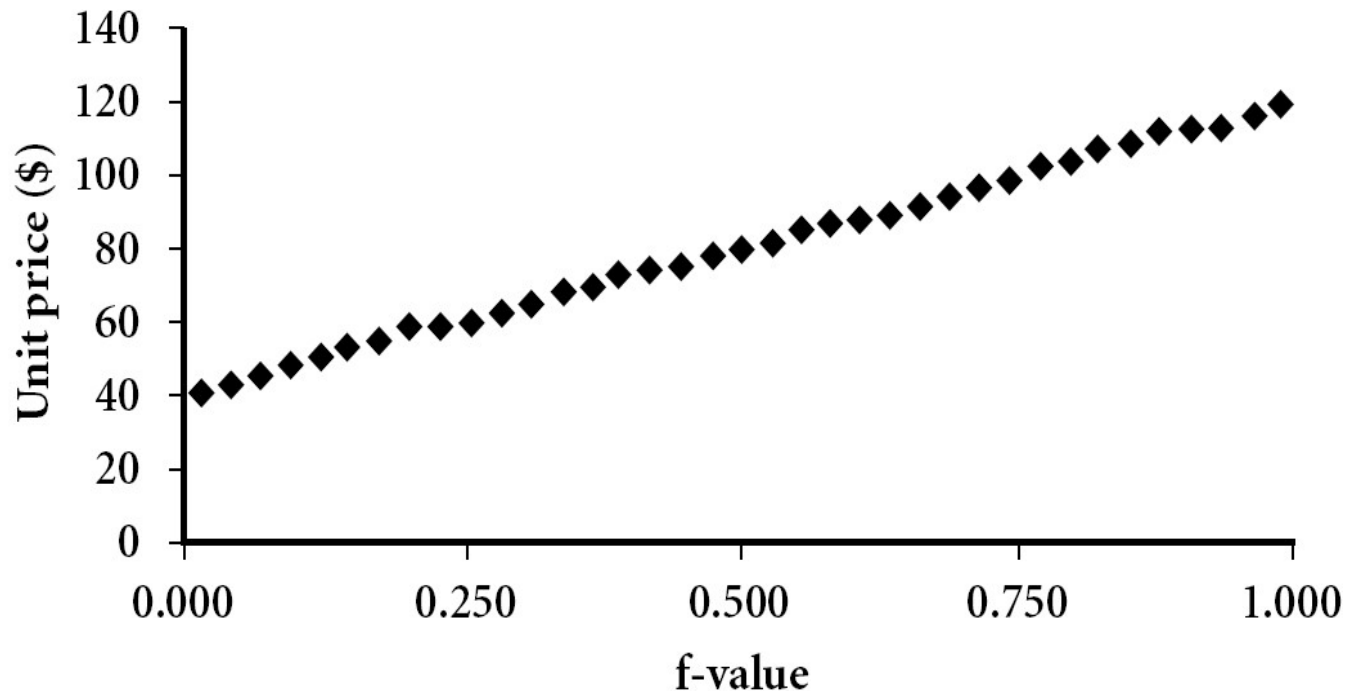


- The two histograms shown in the left may have the same boxplot representation
 - The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions



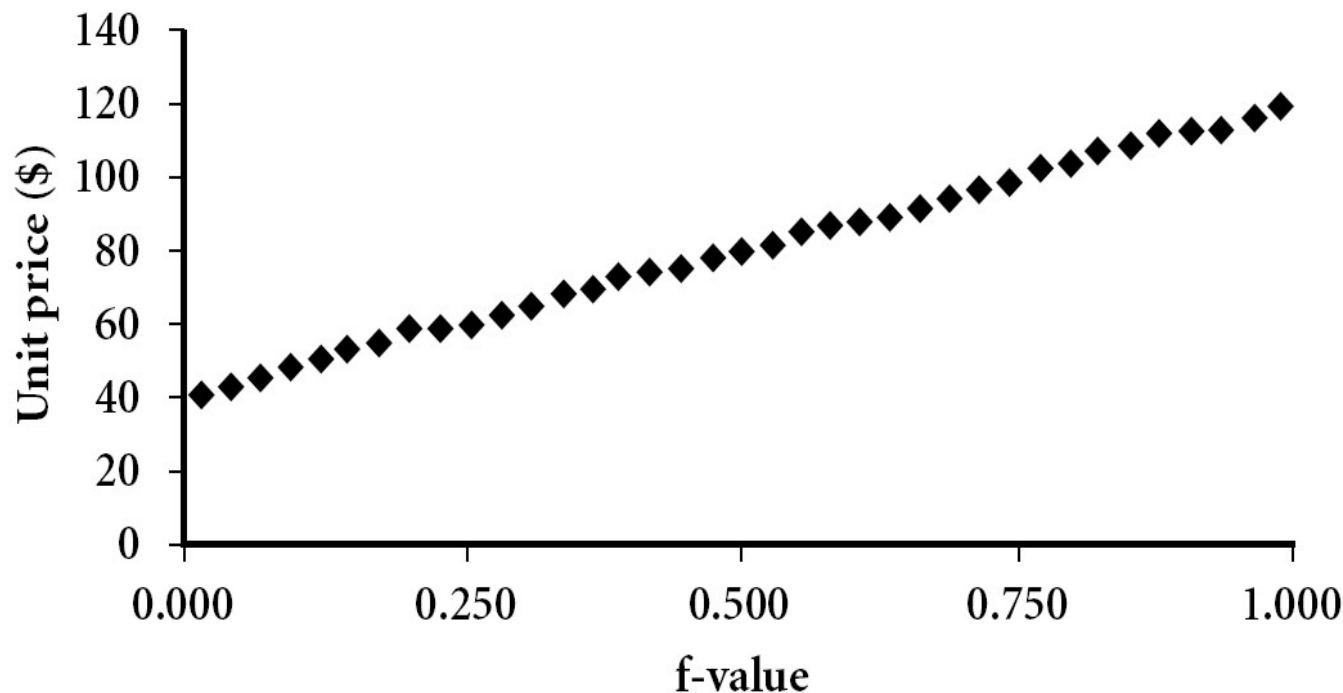
Quantile Plot

- Displays all of the data
 - Allowing the user to assess both the overall behavior and unusual occurrences



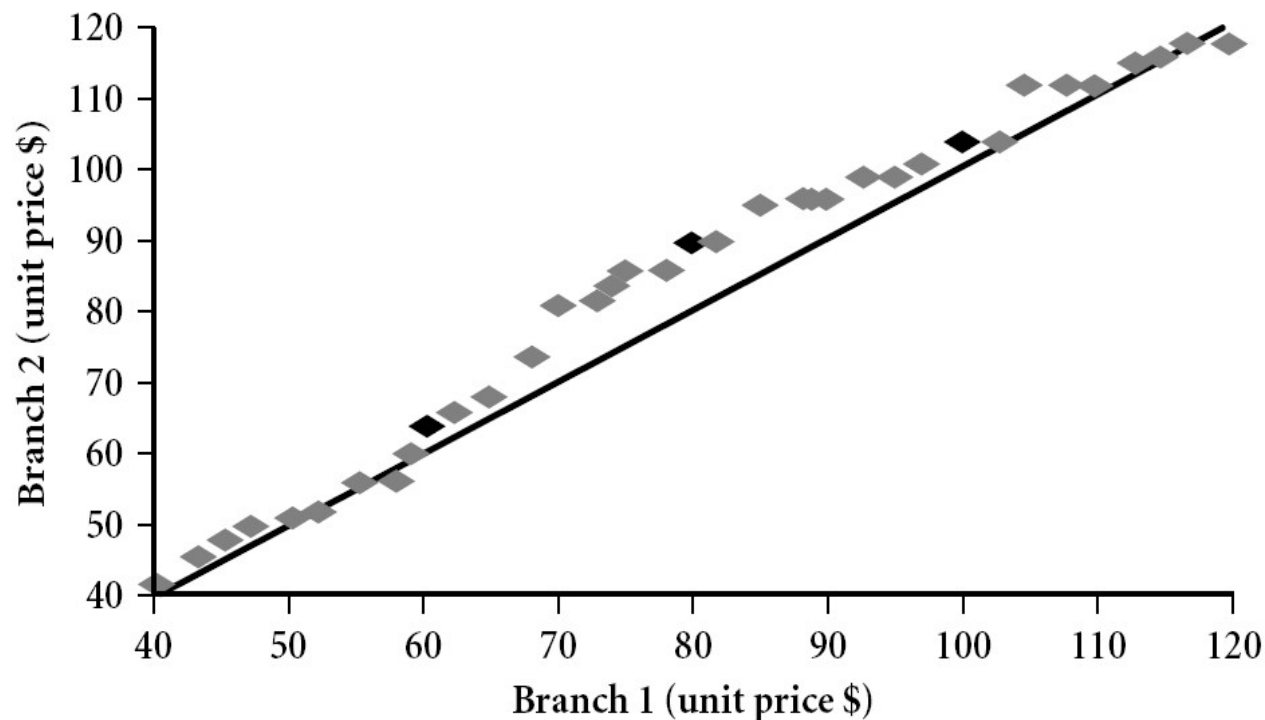
Quantile Plot

- 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



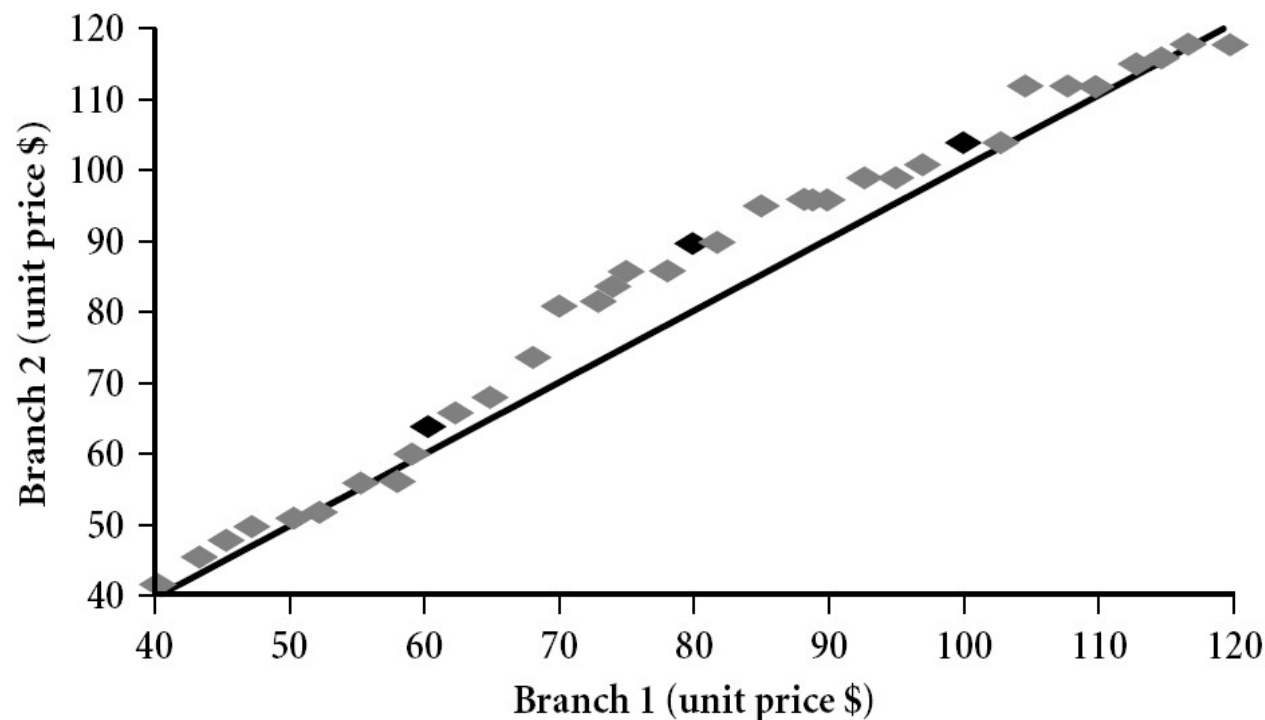
Quantile-Quantile (Q-Q) Plot

- Displays 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?



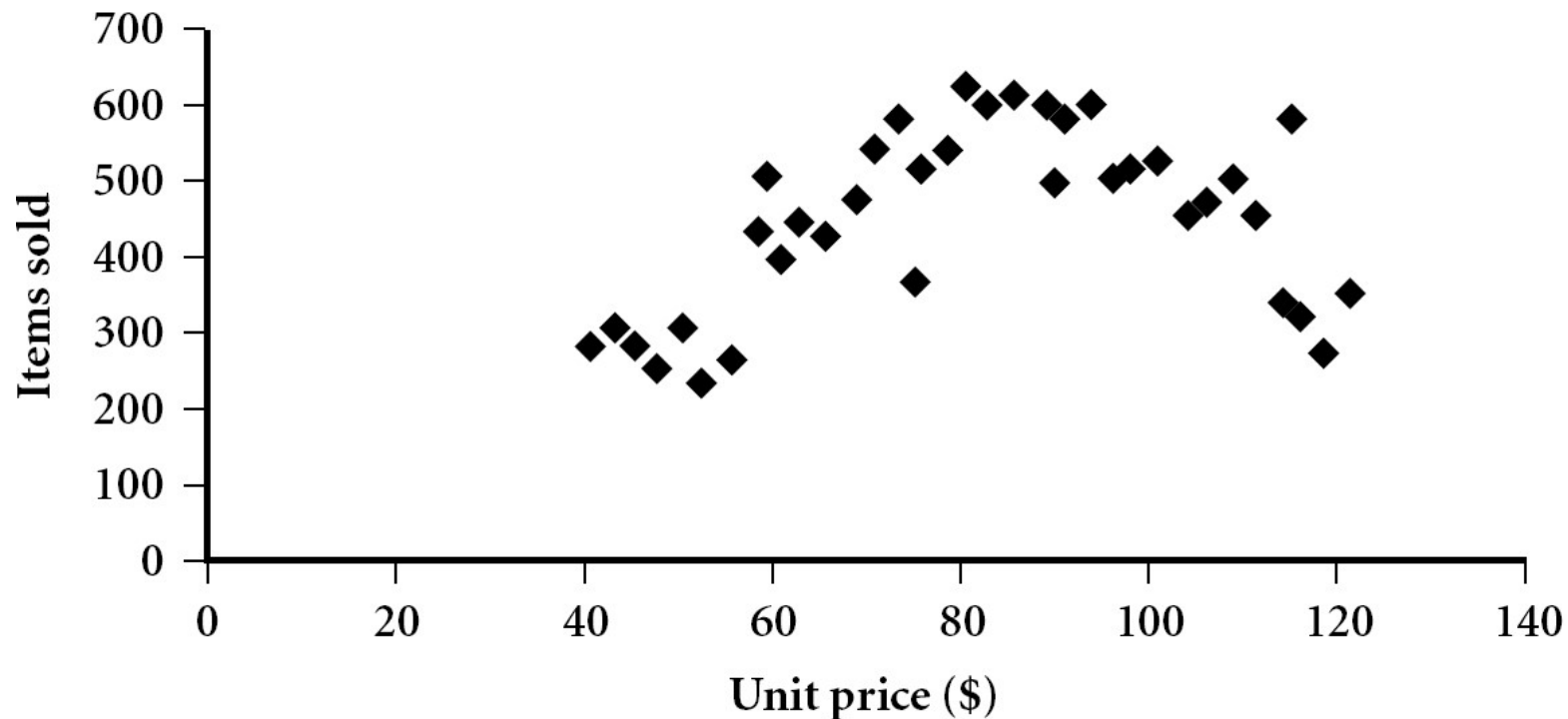
Quantile-Quantile (Q-Q) Plot

- 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.

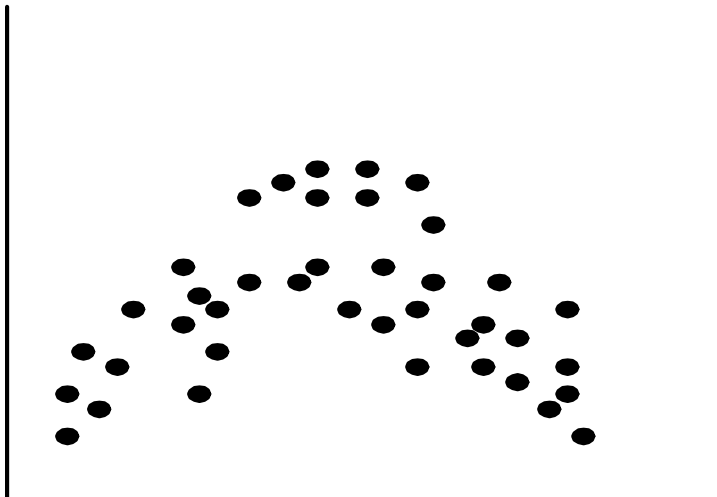
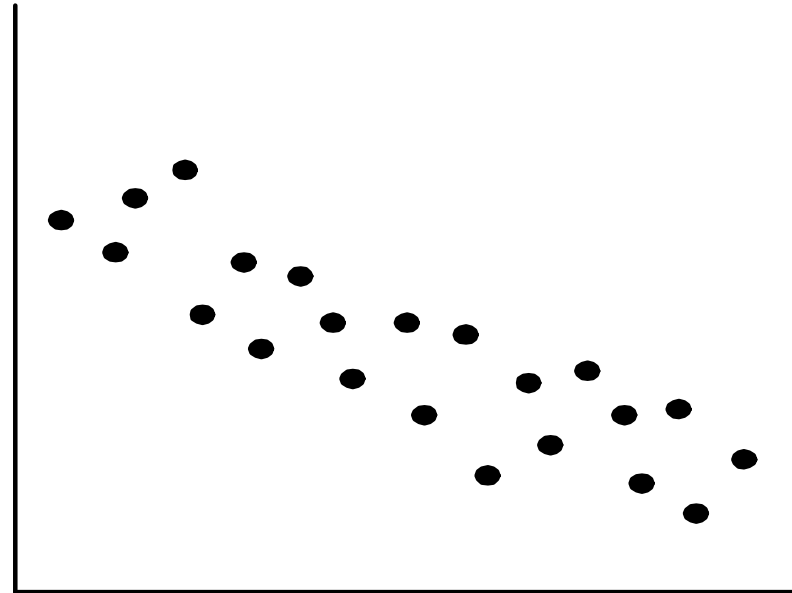
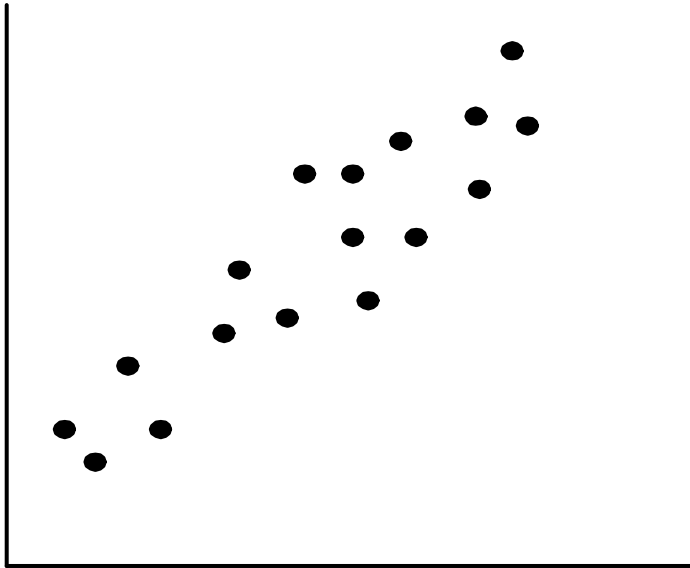


Scatter plot

- Each pair of values is treated as a pair of coordinates and plotted as a point in the plane
- Provides a first look at bivariate data to see clusters of points, outliers, etc



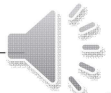
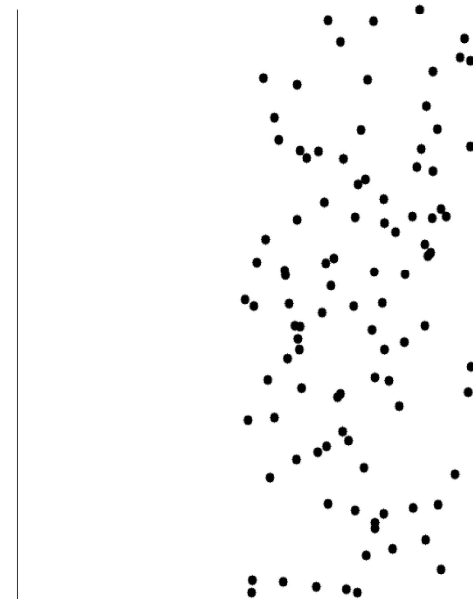
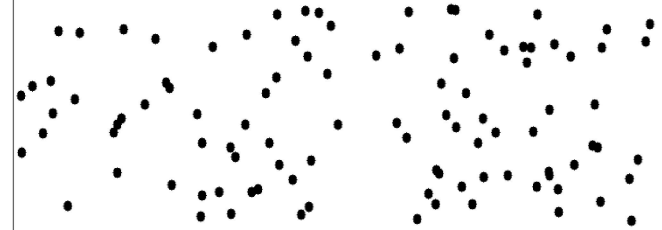
Positively and Negatively Correlated Data



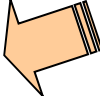
- The left half fragment is positively correlated
- The right half is negative correlated



Uncorrelated Data



Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Measuring Data Similarity and Dissimilarity 
- Summary

Similarity and Dissimilarity

- **Similarity**

- Numerical measure of **how much alike** two data objects are
- This value is higher when objects are more alike
- Often falls in the range $[0,1]$

- **Dissimilarity** (e.g., distance)

- Numerical measure of **how much different** two data objects are
- Lower when objects are more alike
- Minimum dissimilarity is often 0
- Upper limit varies

- **Proximity** refers to a similarity or dissimilarity

Data Matrix and Dissimilarity Matrix

■ Data matrix

- n data points with p dimensions (attributes)
- Two modes

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

■ Dissimilarity matrix

- n data points, but registers only the distance
- A triangular matrix
- Single mode

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

Proximity Measure for Nominal Attributes

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)
- Method 1: Simple matching
 - m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- Method 2: Use a large number of binary attributes
 - creating a new binary attribute for each of the M nominal states

Proximity Measure for Binary Attributes

- A contingency table for binary data

		O j		
		1	0	sum
O i	1	q	r	q + r
	0	s	t	s + t
sum		q + s	r + t	p

- Distance measure for **symmetric** binary variables:
- Distance measure for **asymmetric** binary variables:
- Jaccard coefficient (**similarity** measure for **asymmetric** binary variables):

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$d(i, j) = \frac{r + s}{q + r + s}$$

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

Dissimilarity between Binary Variables

■ Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- Gender is a symmetric attribute (**ignored** in this case)
- The remaining attributes are **asymmetric** binary
- Let the values **Y** and **P** be 1, and the value N 0

$$d (jack , mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d (jack , jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d (jim , mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$



Standardizing Numeric Data

- Z-score:
$$Z = \frac{x - \mu}{\sigma}$$
 - X: raw score to be standardized, μ : mean of the population, σ : standard deviation
 - Meaning: the distance between the raw score and the population mean in units of the standard deviation
 - “-” when the raw score is **below** the mean
 - “+” when the raw score is **above** the mean

Standardizing Numeric Data

- An alternative way: Calculate the **mean absolute deviation**

$$s_f = \frac{1}{n} (|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

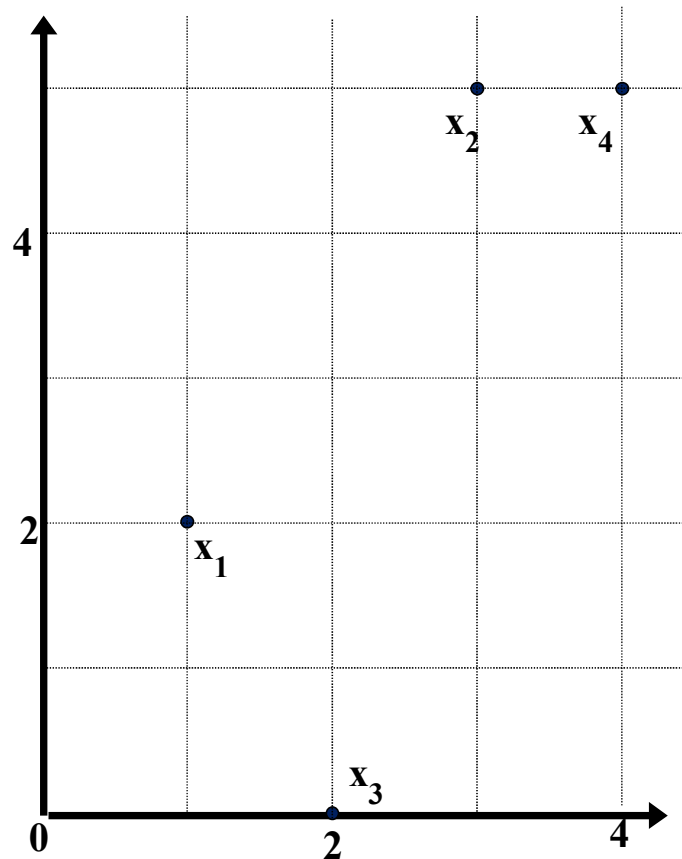
$$m_f = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf}).$$

- standardized measure (*z-score*): $z_{if} = \frac{x_{if} - m_f}{s_f}$

- Using mean absolute deviation is more robust than using standard deviation (when outliers exist)

Example:

Data Matrix and Dissimilarity Matrix



Data Matrix

point	attribute1	attribute2
$x1$	1	2
$x2$	3	5
$x3$	2	0
$x4$	4	5

Dissimilarity Matrix
(with Euclidean Distance)

	$x1$	$x2$	$x3$	$x4$
$x1$	0			
$x2$	3.61	0		
$x3$	5.1	5.1	0	
$x4$	4.24	1	5.39	0

