

# Measurement and Data

# Data describes the real world

- Data maps entities in the domain of interest to symbolic representation by means of a measurement procedure
- Numerical relationships between variables capture relationships between objects
- Measurement process is crucial

# Types of Measurement

- **Ordinal**, e.g., excellent=5, very good=4, good=3...
- **Nominal**, e.g., religion, profession
  - Need non-metric methods
- **Ratio**, e.g., weight
  - has concatenation property, two weights add to balance a third:  $2+3 = 5$
  - changing scale (multiplying values by a constant) does not change ratio
- **Interval**, e.g., temperature, calendar time
  - Unit of measurement is arbitrary, as well as origin

# Distance Measures

- Many data mining techniques (e.g., nn-classification, cluster analysis) are based on similarity measures between objects
- $s(i,j)$ : similarity,  $d(i,j)$ : dissimilarity
- Possible transformations:  $d(i,j) = 1 - s(i,j)$  or  $d(i,j) = \sqrt{2 * (1 - s(i,j))}$

# Metric Properties

1.  $d(i,j) \geq 0$ : Positivity
2.  $d(i,j) = d(j,i)$ : Commutativity
3.  $d(i,j) \leq d(i,k) + d(k,j)$ : Triangle Inequality

# Euclidean Distance between vectors

$$d_E(x, y) = \left( \sum_{k=1}^p (x_k - y_k)^2 \right)^{1/2}$$

# Commensurability

- Euclidean distance assumes variables are commensurate
- E.g., each variable a measure of length
- If one were weight and other was length there is no obvious choice of units
- Altering units would change which variables are important

# Standardizing the Data

- Divide each variable by its standard deviation
- Standard deviation for the  $k^{\text{th}}$  variable is

$$\mathbf{s}_k = \left( \frac{1}{n} \sum_{i=1}^n (x_k(i) - \mathbf{m}_k)^2 \right)^{\frac{1}{2}}$$

where

$$\mathbf{m}_k = \frac{1}{n} \sum_{i=1}^n x_k(i)$$

# Weighted Euclidean Distance

- If we know relative importance of variables

$$d_{WE}(i, j) = \left( \sum_{k=1}^p w_k ((x_k(i) - x_k(j))^2) \right)^{\frac{1}{2}}$$

# Need for Covariance in distance measure

- Suppose we measured a cup's height 100 times and diameter only once
- Clearly height will dominate although 99 of the height measurements are not contributing anything
- They are very highly correlated
- To eliminate redundancy we need a *data-driven method*

# Sample Covariance between X and Y

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n \left( x(i) - \bar{x} \right) \left( y(i) - \bar{y} \right)$$

- Measure of how X and Y vary together
- Large positive value if large values of X tend to be associated with large values of Y and small values of X with small values of Y
- Large negative value if large values of X tend to be associated with small values of Y

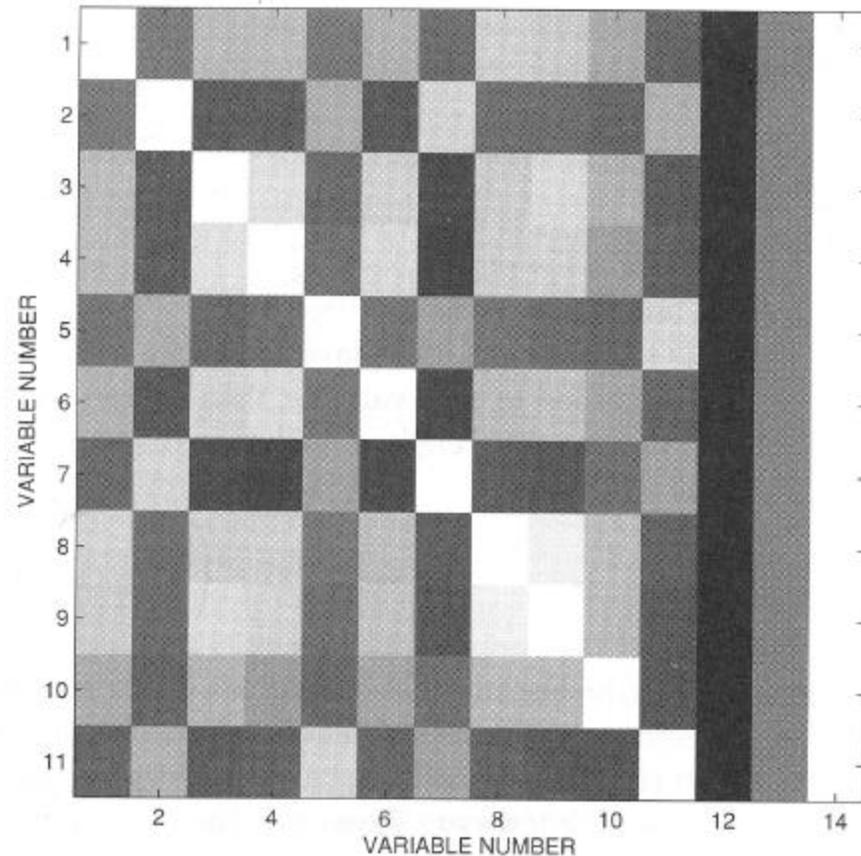
# Correlation Coefficient

Value of Covariance is dependent upon ranges of X and Y

Removed by dividing values of X by their standard deviation and values of Y by their standard deviation

$$r(X, Y) = \frac{\sum_{i=1}^n (x(i) - \bar{x})(y(i) - \bar{y})}{s_x s_y}$$

# Correlation Matrix



**Figure 2.1** A sample correlation matrix plotted as a pixel image. White corresponds to +1 and black to -1. The three rightmost columns contain values of -1, 0, and +1 (respectively) to provide a reference for pixel intensities. The remaining  $11 \times 11$  pixels represent the  $11 \times 11$  correlation matrix. The data come from a well-known data set in the regression research literature, in which each data vector is a suburb of Boston and each variable represents a certain general characteristic of a suburb. The variable names are (1) per-capita crime rate, (2) proportion of area zoned for large residential lots, (3) proportion of non-retail business acres, (4) nitric oxide concentration, (5) average number of rooms perdwelling, (6) proportion of pre-1940 homes, (7) distance to retail centers index, (8) accessibility to highways index, (9) property tax rate, (10) pupil-to-teacher ratio, and (11) median value of owner-occupied homes.

# Mahalanobis Distance

$$d_M(i, j) = [(x(i) - x(j))^T \Sigma^{-1} (x(i) - x(j))]^{\frac{1}{2}}$$

# Generalizing Euclidean Distance

- Minkowski or  $L_p$  metric

$$\left( \sum_{k=1}^p (x_k(i) - x_k(j))^p \right)^{\frac{1}{p}}$$

- $p = 2$  gives the Euclidean metric

# Minkowski metric

- $p = 1$  is the Manhattan or city block metric

$$\sum_{k=1}^p |x_k(i) - x_k(j)|$$

- $p = \text{infinity}$  yields

$$\max_k |x_k(i) - x_k(j)|$$

# Multivariate Binary Data

- Most obvious measure is *Hamming Distance* normalized by number of bits

$$\frac{S_{11} + S_{00}}{S_{11} + S_{10} + S_{01} + S_{00}}$$

- If we don't care about irrelevant properties had by neither object we have *Jaccard Coefficient*

$$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$$

- *Dice Coefficient* extends this argument. If 00 matches are irrelevant then 10 and 01 matches should have half relevance

# Some Similarity/Dissimilarity Measures for N-dim Binary Vectors

Measure	$S(X, Y)$	$D(X, Y)$
Jaccard-Needham	$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{S_{11} + S_{10} + S_{01}}$
Dice	$\frac{S_{11}}{2S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{2S_{11} + S_{10} + S_{01}}$
Correlation	$\frac{S_{11}S_{00} - S_{10}S_{01}}{\sigma}$	$\frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2\sigma}$
Yule	$\frac{S_{11}S_{00} - S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$	$\frac{S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$
* Russell-Rao	$\frac{S_{11}}{N}$	$\frac{N - S_{11}}{N}$
* Sokal-Michener	$\frac{S_{11} + S_{00}}{N}$	$1 - \frac{S_{11} + S_{00}}{N}$
Rogers-Tanimoto	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2S_{10} + 2S_{01}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$
* Rogers-Tanimoto-a	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2(N - S_{11} - S_{00})}{2N - S_{11} - S_{00}}$
Kulzinsky	$\frac{S_{11}}{(S_{10} + S_{01})}$	$\frac{S_{10} + S_{01} - S_{11} + N}{(S_{10} + S_{01} + N)}$

where  $\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}$

$$S_{00}(X, Y) = \bar{X} \cdot \bar{Y}, S_{01}(X, Y) = \bar{X} \cdot Y \quad S_{10}(X, Y) = X \cdot \bar{Y}, S_{11}(X, Y) = X \cdot Y$$

# Some Similarity/Dissimilarity Measures for N-dim Binary Vectors

Measure	$S(X, Y)$	$D(X, Y)$
Jaccard-Needham	$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{S_{11} + S_{10} + S_{01}}$
Dice	$\frac{S_{11}}{2S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{2S_{11} + S_{10} + S_{01}}$
Correlation	$\frac{S_{11}S_{00} - S_{10}S_{01}}{\sigma}$	$\frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2\sigma}$
Yule	$\frac{S_{11}S_{00} - S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$	$\frac{S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$
* Russell-Rao	$\frac{S_{11}}{N}$	$\frac{N - S_{11}}{N}$
* Sokal-Michener	$\frac{S_{11} + S_{00}}{N}$	$1 - \frac{S_{11} + S_{00}}{N}$
Rogers-Tanmoto	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2S_{10} + 2S_{01}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$
* Rogers-Tanmoto-a	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2(N - S_{11} - S_{00})}{2N - S_{11} - S_{00}}$
Kulzinsky	$\frac{S_{11}}{(S_{10} + S_{01})}$	$\frac{S_{10} + S_{01} - S_{11} + N}{(S_{10} + S_{01} + N)}$

where  $\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}$

$$S_{00}(X, Y) = \bar{X} \cdot \bar{Y}, S_{01}(X, Y) = \bar{X} \cdot Y \quad S_{10}(X, Y) = X \cdot \bar{Y}, S_{11}(X, Y) = X \cdot Y$$

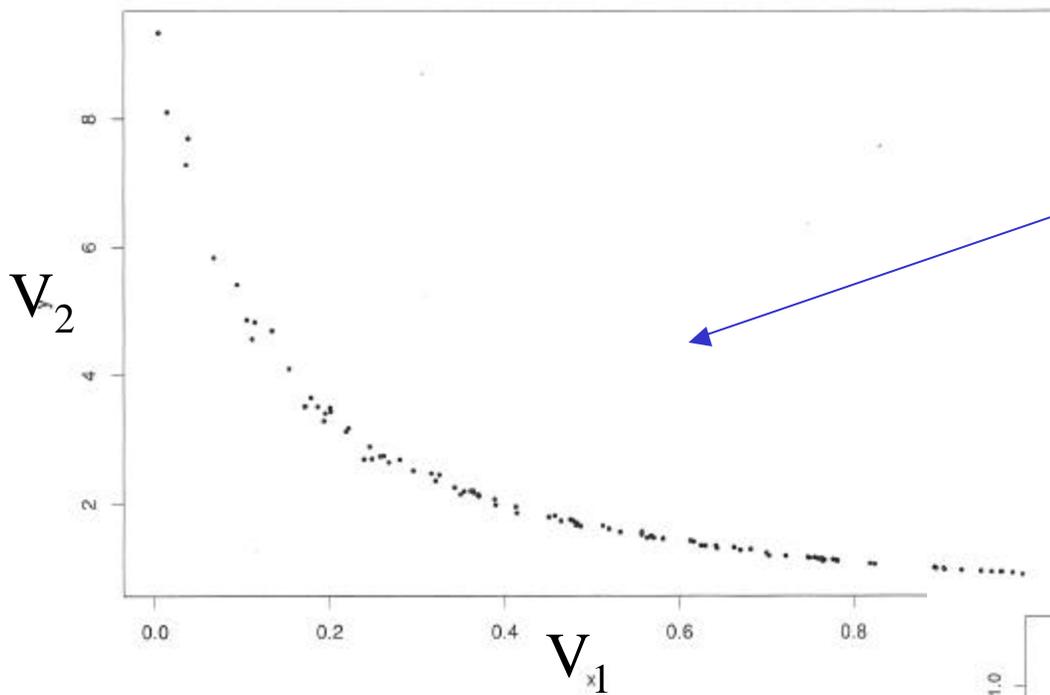
# Weighted Dissimilarity Measures for Binary Vectors

- Unequal importance to '0' matches and '1' matches
- Multiply  $S_{00}$  with  $\beta$  ( $[0,1]$ )
- Examples:

$$D^{sm}(X, Y) = 1 - \frac{S_{11} + \mathbf{b} \cdot S_{00}}{N}$$

$$D^{rta}(X, Y) = \frac{2(N - S_{11} - \mathbf{b} \cdot S_{00})}{2N - S_{11} - \mathbf{b} \cdot S_{00}}$$

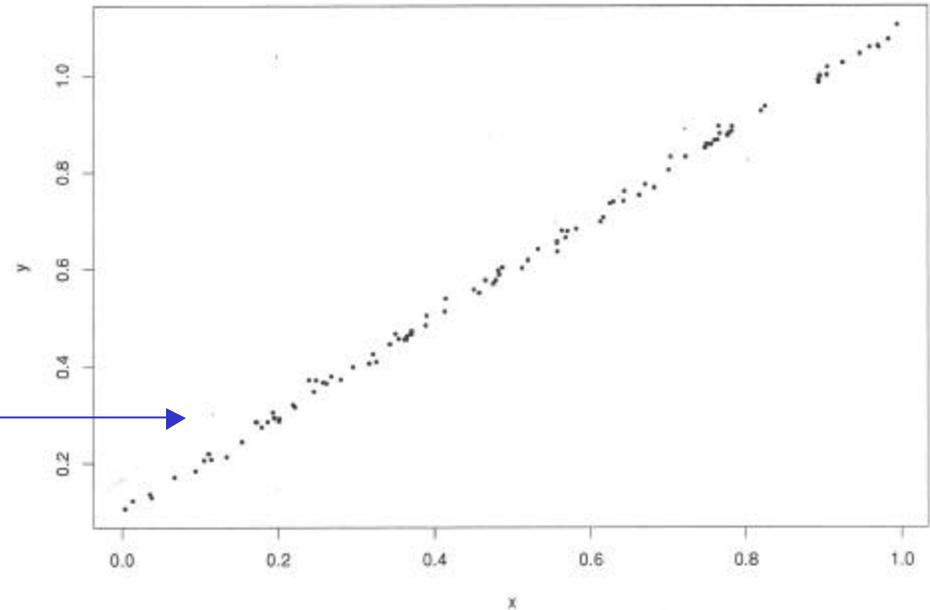
# Transforming the Data



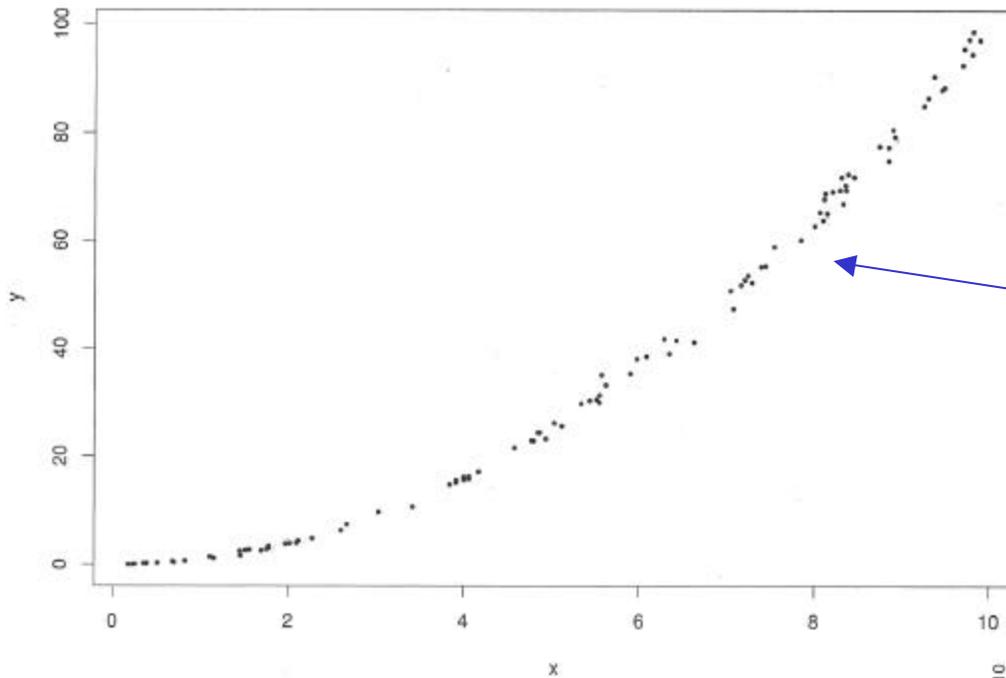
$V_1$  is non-linearly Related to  $V_2$

**Figure 2.2** A simple nonlinear relationship between variable  $V_1$  and subsequent figures  $V_1$  and  $V_2$  are on the  $X$  and  $Y$  axes respectively.

$V_3=1/V_2$  is linearly related to  $V_1$



**Figure 2.3** The data of figure 2.2 after the simple transformation of  $V_2$  to  $1/V_2$ .



Variance increases  
(regression assumes  
variance is constant)

Figure 2.4 Another simple nonlinear relationship. Here the variance of  $V_2$  as  $V_1$  increases.

Square root transformation  
keeps the variance constant

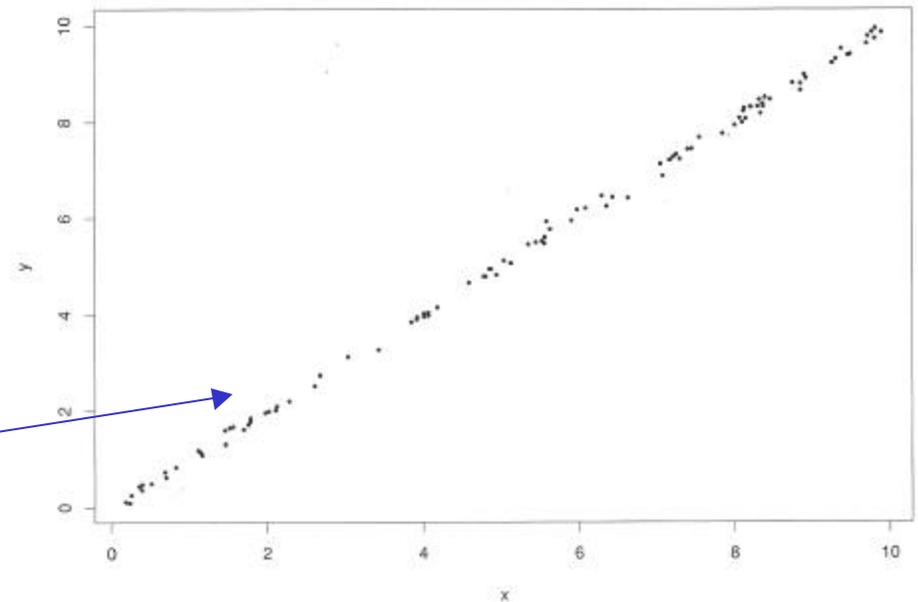


Figure 2.5 The data of figure 2.4 after a simple square root transformation of  $V_2$ . Now the variance of  $V_2$  is relatively constant as  $V_1$  increases.

# Form of Data

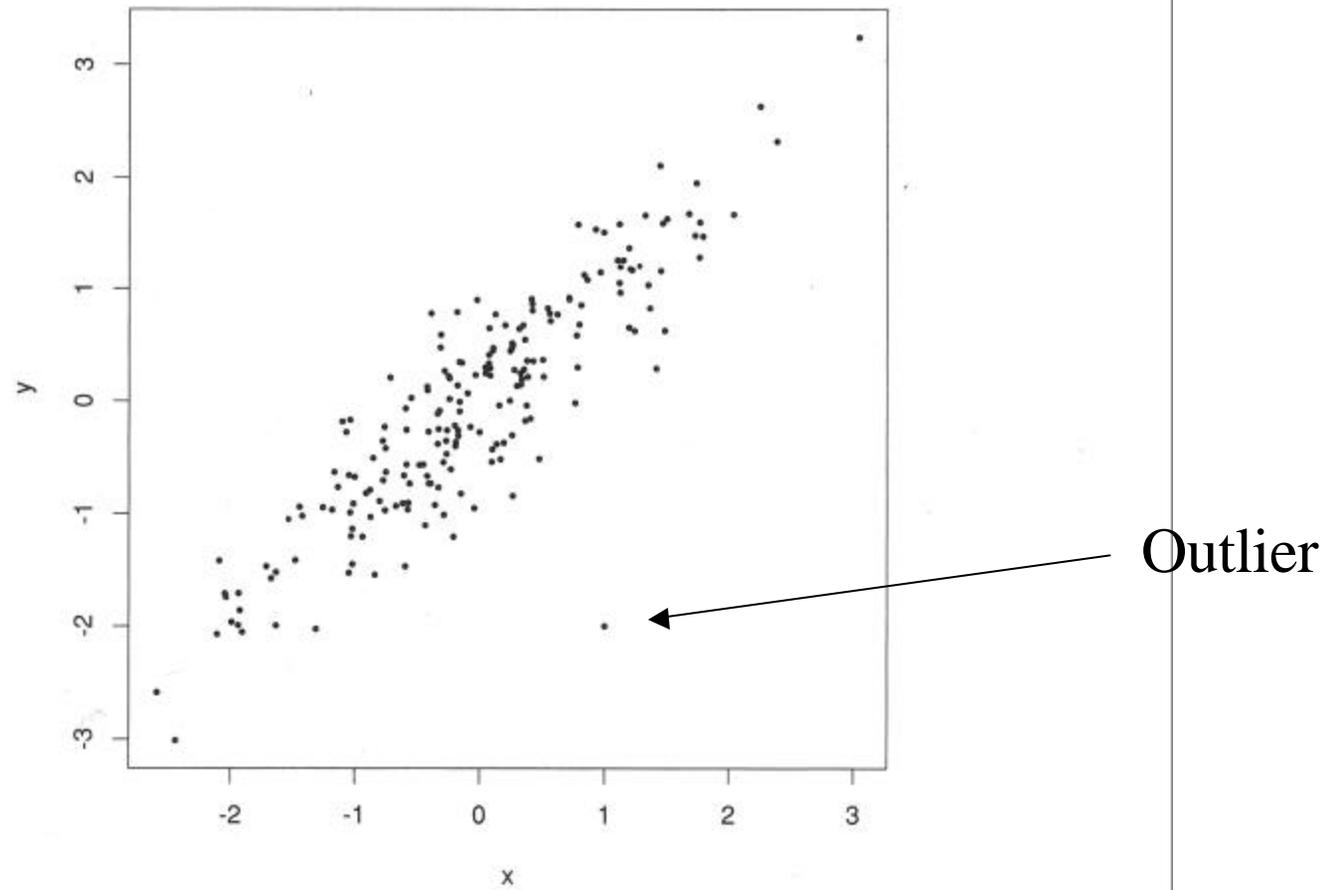
# Data Matrix

- A set of  $p$  measurements on objects  $o(1)\dots o(n)$
- $n$  rows and  $p$  columns
- Also called *standard data*, *data matrix* or *table*

# Multirelational Data

- Payroll database has
  - **Employees table**: **name, department-name, age, salary**
  - **Department table**: **department-name, budget, manager**
- The tables are connected to each other by the **department-name** field and the fields **name** and **manager**
- Can be combined together, e.g., with fields **name, department-name, age, salary, budget, manager**
- Or create as many rows as department-names
- Flattening may require needless replication of values

# Data Quality



**Figure 2.6** A plot of 200 points from highly positively correlated bivariate data (from a bivariate normal distribution), with a single easily identifiable outlier.