

Measurements and Data

Sargur Srihari

University at Buffalo
The State University of New York

Topics

- Types of Data
- Distance Measurement
- Data Transformation
- Forms of Data
- Data Quality

Importance of Measurement

- Aim of mining structured data is to discover relationships that exist in the real world
 - business, physical, conceptual
- Instead of looking at real world we look at data describing it
- 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., color, 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 (multiply by constant) does not change ratio
- Interval
 - e.g., temperature, calendar time
 - Unit of measurement is arbitrary, as well as origin

Operational Measurement

- Measuring Programming Effort (Halstead 1977)

Programming effort $e = am(n+m)\log(a+b)/2b$

a = no of unique operators

b = no of unique operands

n = no of total operator occurrences

m = no of operand occurrences

- Defines programming effort as well as a way of measuring it.
- Operational measurements are concerned with prediction whereas non-operational measurements are concerned with description


Distance and Similarity

- Many data mining techniques are based on similarity measures between objects
 - nearest-neighbor classification
 - cluster analysis,
 - multi-dimensional scaling
- $s(i,j)$: similarity, $d(i,j)$: dissimilarity
- Possible transformations:
 - $d(i,j) = 1 - s(i,j)$ or
 - $d(i,j) = \sqrt{2 * (1 - s(i,j))}$
- Proximity is a general term to indicate similarity and dissimilarity
- Distance is used to indicate dissimilarity

Metric Properties

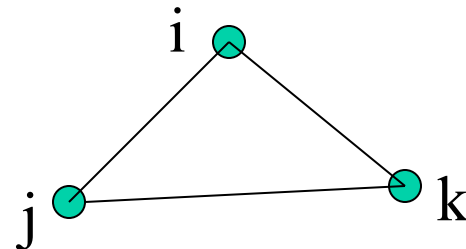
A metric is a dissimilarity (distance) measure that satisfies:

1. $d(i,j) \geq 0$
2. $d(i,j) = d(j,i)$
3. $d(i,j) \leq d(i,k) + d(k,j)$

Positivity 

Commutativity

Triangle Inequality

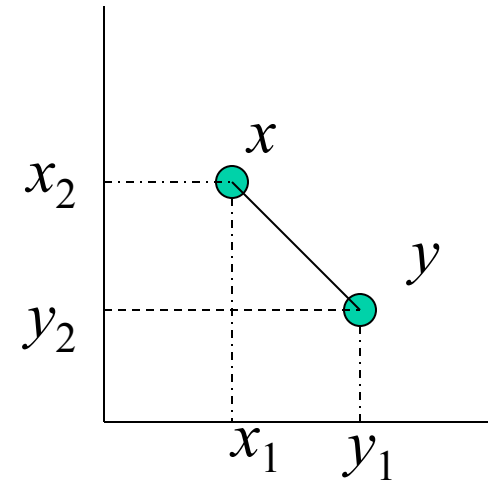


Examples of Metrics

- Euclidean Distance d_E
 - Standardized (divide by variance)
 - Weighted d_{WE}
- Minkowski measure
 - Manhattan Distance
- Mahanalobis Distance d_M
 - Use of Covariance
- Binary data Distances

Euclidean Distance between Vectors

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



- 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

when variables are not commensurate

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

$$\sigma_k = \left(\frac{1}{n} \sum_{i=1}^n (x_k(i) - \mu_k)^2 \right)^{\frac{1}{2}}$$

where

$$\mu_k = \frac{1}{n} \sum_{i=1}^n x_k(i)$$

- Updated value that removes the effect of scale:


$$x_k' = \frac{x_k}{\sigma_k}$$

Weighted Euclidean Distance

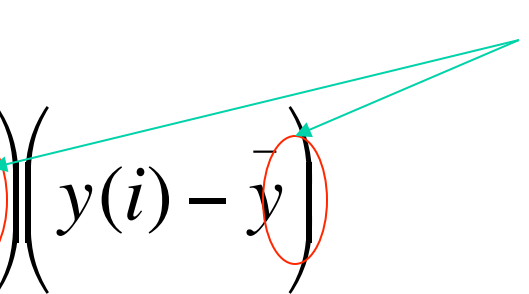
- If we know relative importance of variables

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

Use of Covariance in Distance

- Similarities between cup
- Suppose we measure cup-height 100 times and diameter only once
 - 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*
 - approach is to not only to standardize data in each direction but also to use covariance between variables

Covariance between two Scalar Variables

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


The diagram shows two red circles around the sample mean terms \bar{x} and \bar{y} in the formula. Two green arrows originate from the text "Sample means" on the right and point to these two circles.

- A scalar value to measure how x and y vary together
- Obtained by
 - multiplying for each sample its mean-centered value of x with mean-centered value of y
 - and then adding over all samples
- 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
- With d variables can construct a $d \times d$ matrix of covariances
 - Such a covariance matrix is symmetric.

For Vectors: Covariance Matrix and Data Matrix

- Let $X = n \times d$ data matrix
- Rows of X are the data vectors $x(i)$
- Definition of covariance:

$$Cov(i, j) = \frac{1}{n} \sum_{k=1}^n \left(x_k(i) - \bar{x} \right) \left(y_k(i) - \bar{y} \right)$$

- If values of X are mean-centered
 - i.e., value of each variable is relative to the sample mean of that variable
 - then $V = X^T X$ is the $d \times d$ covariance matrix

Correlation Coefficient

Value of Covariance is dependent upon ranges of x and y

Dependency is removed by

dividing values of x by their standard deviation
and values of y by their standard deviation

$$\rho(X, Y) = \frac{\sum_{i=1}^n (x(i) - \bar{x})(y(i) - \bar{y})}{\sigma_x \sigma_y}$$

With p variables, can form a $d \times d$ correlation matrix

Correlation Matrix

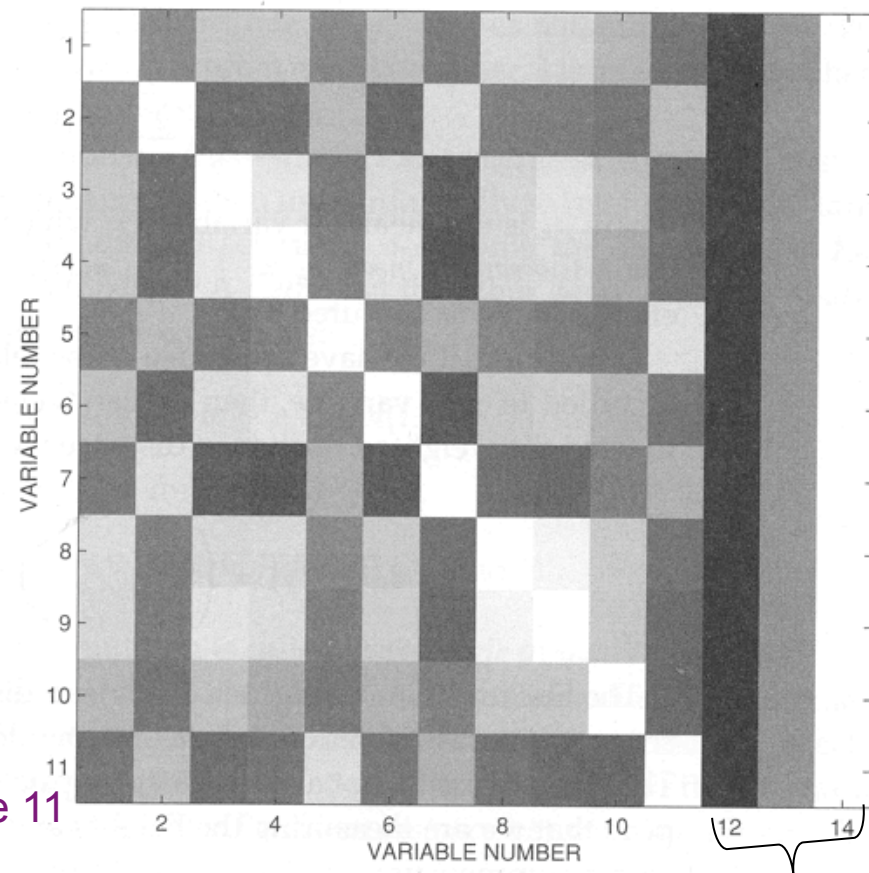
Housing related variables

across city suburbs ($d=11$)

11 x 11 pixel image (White 1, Black -1)

Columns 12-14 have values -1,0,1 for pixel intensity reference

Remaining represent correlation matrix



Reference for -1, 0,+1

Variables 3 and 4 are highly negatively correlated with Variable 2

Variable 5 is positively correlated with Variable 11

Variables 8 and 9 are highly correlated

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 11×11 pixels represent the 11×11 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 per dwelling, (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.

Incorporating Covariance Matrix in Distance

Mahalanobis Distance between samples $x(i)$ and $x(j)$ is:

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

T is transpose

Σ is $d \times d$ covariance matrix

Σ^{-1} standardizes data relative to Σ

Matrix multiplication
yields a scalar value

d_M discounts the effect of several highly correlated variables

Generalizing Euclidean Distance

Minkowski or L_λ metric $\left(\sum_{k=1}^p (x_k(i) - x_k(j))^\lambda \right)^{\frac{1}{\lambda}}$

- $\lambda = 2$ gives the Euclidean metric
- $\lambda = 1$ gives the Manhattan or City-block metric

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

- $\lambda = \infty$ yields

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

Srihari

Distance Measures for Binary Data

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

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

Proportion of variables
on which objects have same value

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

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

Example: two documents
do not have certain terms

- *Dice Coefficient* extends this argument
 - If 00 matches are irrelevant then 10 and 01 matches should have half relevance
- Generalization to discrete values (non-binary)
 - Score 1 for if two objects agree and 0 otherwise
- Adaptation to mixed data types
 - Use additive distance measures

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

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

Srihari

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

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

Srihari

Weighted Dissimilarity Measures for Binary Vectors

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

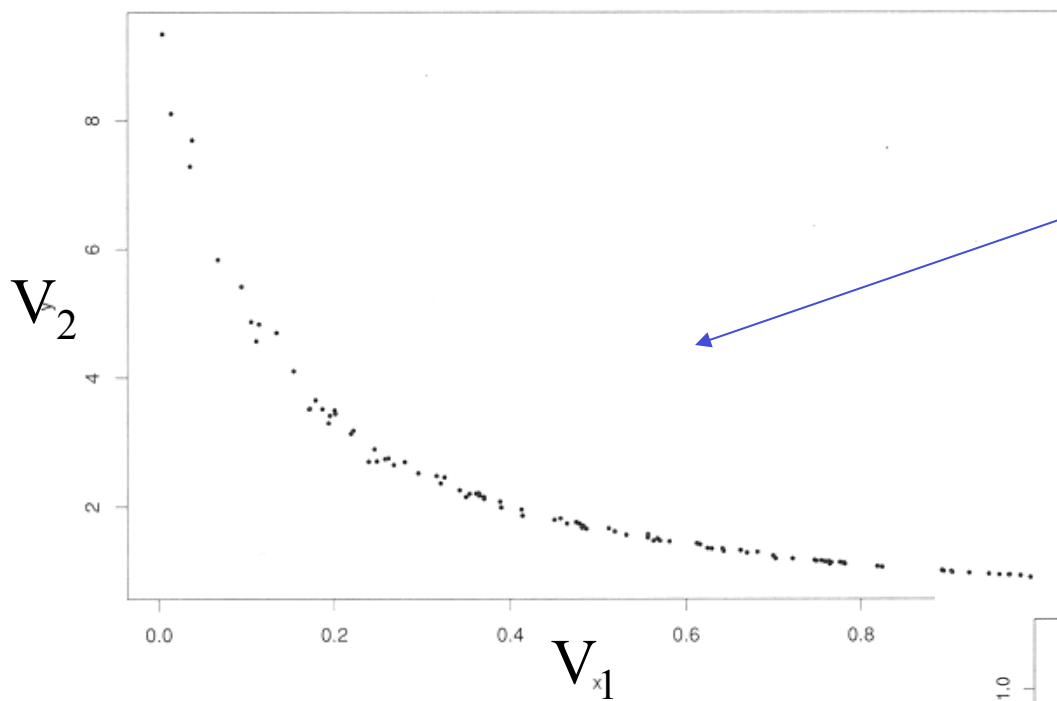
$$D^{sm}(X, Y) = 1 - \frac{S_{11} + \beta \cdot S_{00}}{N}$$

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

Transforming the Data

Model depends on form of data

If Y is a function of X^2 then we could use quadratic function or choose $U = X^2$ and use a linear fit



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

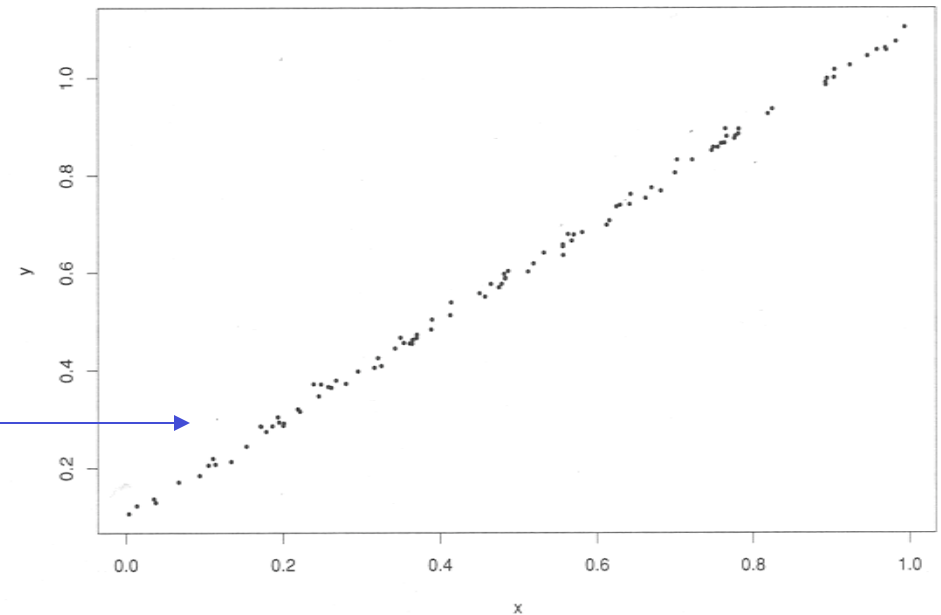
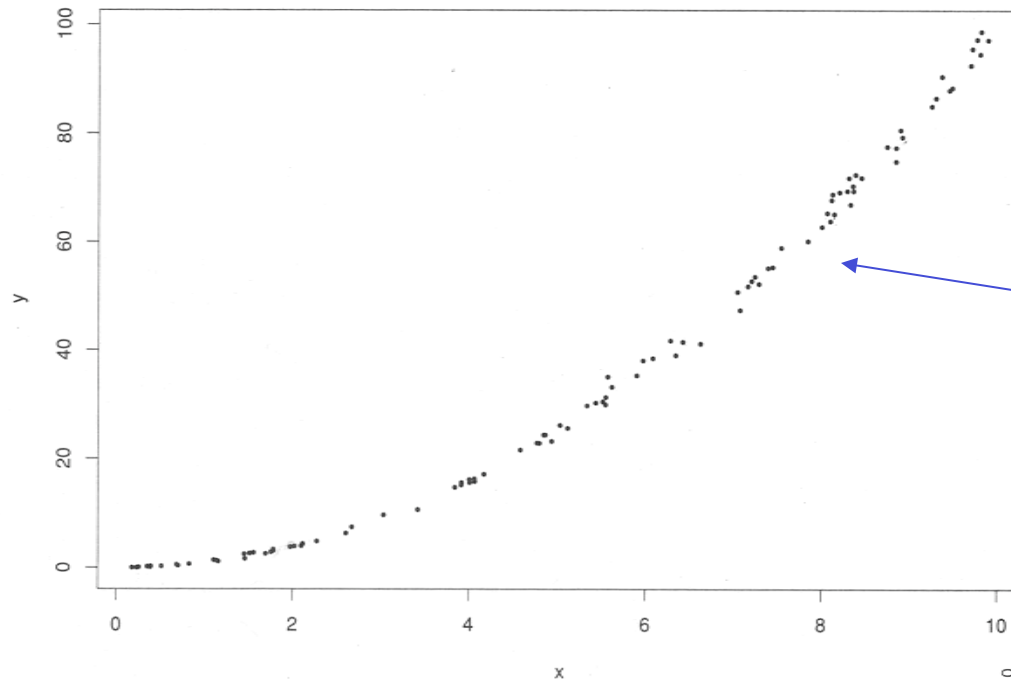


Figure 2.3 The data of figure 2.2 after the simple transformation of V_2 to $1/V_2$.



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

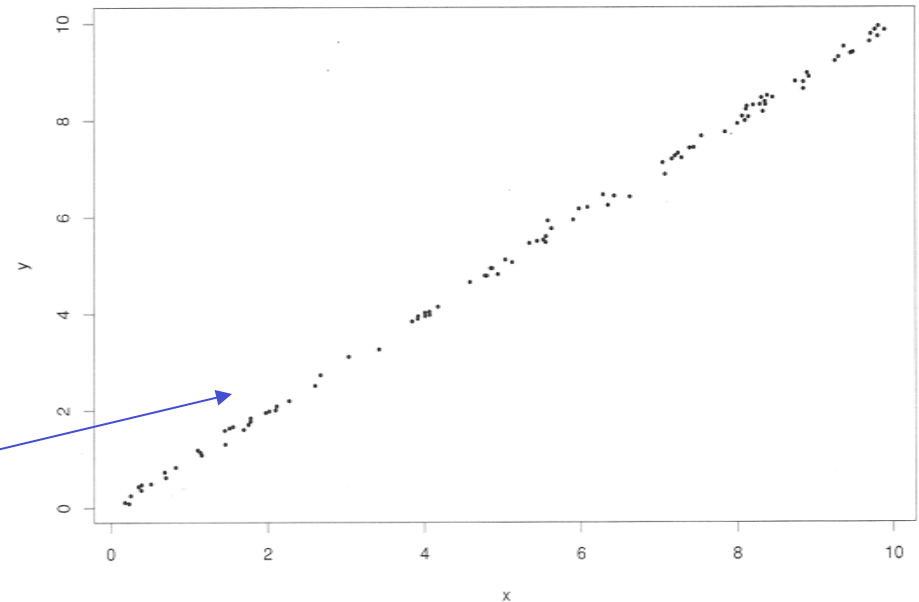


Figure 2.5 The data of figure 2.4 after a simple square root transformation of V_2 . Now the variance of V_2 is relatively constant as V_1 increases.

Forms of Data

Standard Data (Data Matrix)

Multirelational Data

String

Event Sequence

Hierarchical Data

Data Matrix

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

Multirelational Data (multiple data matrices)

Payroll Database

Name	Department Name	Age	Salary

Department Table

Department Name	Budget		Manager

Can be combined together to form a data matrix with fields
name, department-name, age, salary, budget, manager
Or create as many rows as department-names
Flattening requires needless replication (Storage issues)

String Data

- Sequence of symbols from a finite alphabet
 - Standard matrix form is unsuitable
- Sequence of values from a categorical variable
 - Standard English text (alphanumeric characters, spaces, punctuation marks)
 - Protein and DNA/RNA sequences (A,C,G,T)

Event Sequence Data

- Sequence of pairs of the form {event, occurrence time}
- A string where each sequence item is tagged with an occurrence time
 - Telecommunication alarm log
 - Transaction data (records of retail or financial)
 - Can occur asynchronously

Data Quality

Data Quality for Individual Measurements

- Data Mining Depends on Quality of data
- Many interesting patterns discovered may result from measurement inaccuracies.
- Sources of error
 - Errors in measurement
 - Carelessness
 - Instrumentation failure
 - Inadequate definition of what we are measuring

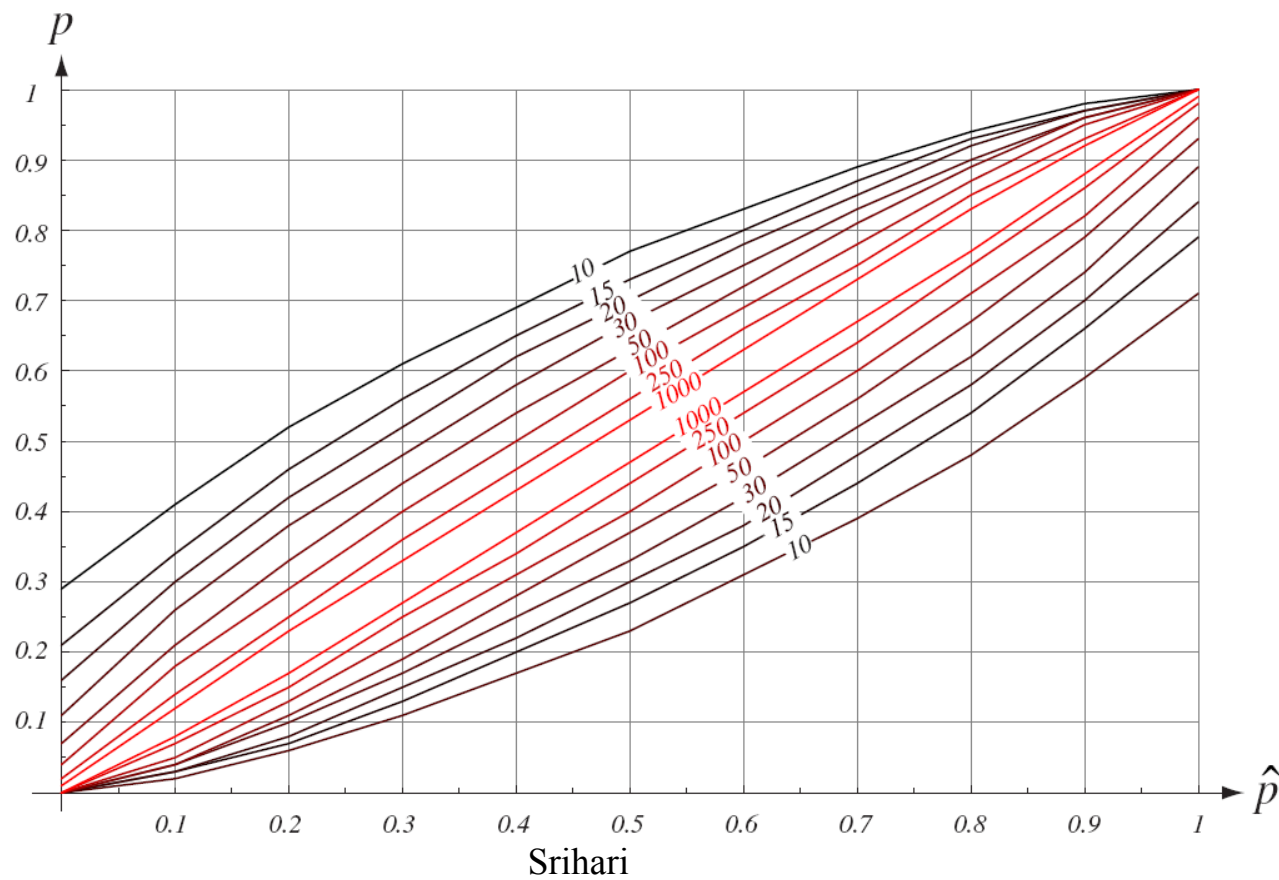
Precision and Accuracy

- Precise Measurement
 - Small variability (measured by variance)
 - Repeated measurements yield same value
 - Many digits of precision is not necessarily accurate (results of calculations give many digits)
- Accurate
 - Not only small variability but close to true value
 - Precise measurement of height with shoes will not give an accurate measurement
 - Mean of repeated measurements and true value is “Bias”

Data Quality for Collections of Data

- Collections of Data
 - Much of statistics is concerned with inference from a sample to a population
 - How to infer things from a fraction about entire population
 - Two sources of error:
 - sample size and bias

- Confidence Intervals



Biased Sample

- Inappropriate samples
 - To calculate average weight of people in New York it would be inappropriate to restrict samples to women, or to office workers
- Random sample is key to make valid inferences
 - Stratification (gender, age, education, occupation)
 - Proportional representation

Anomalous Observations

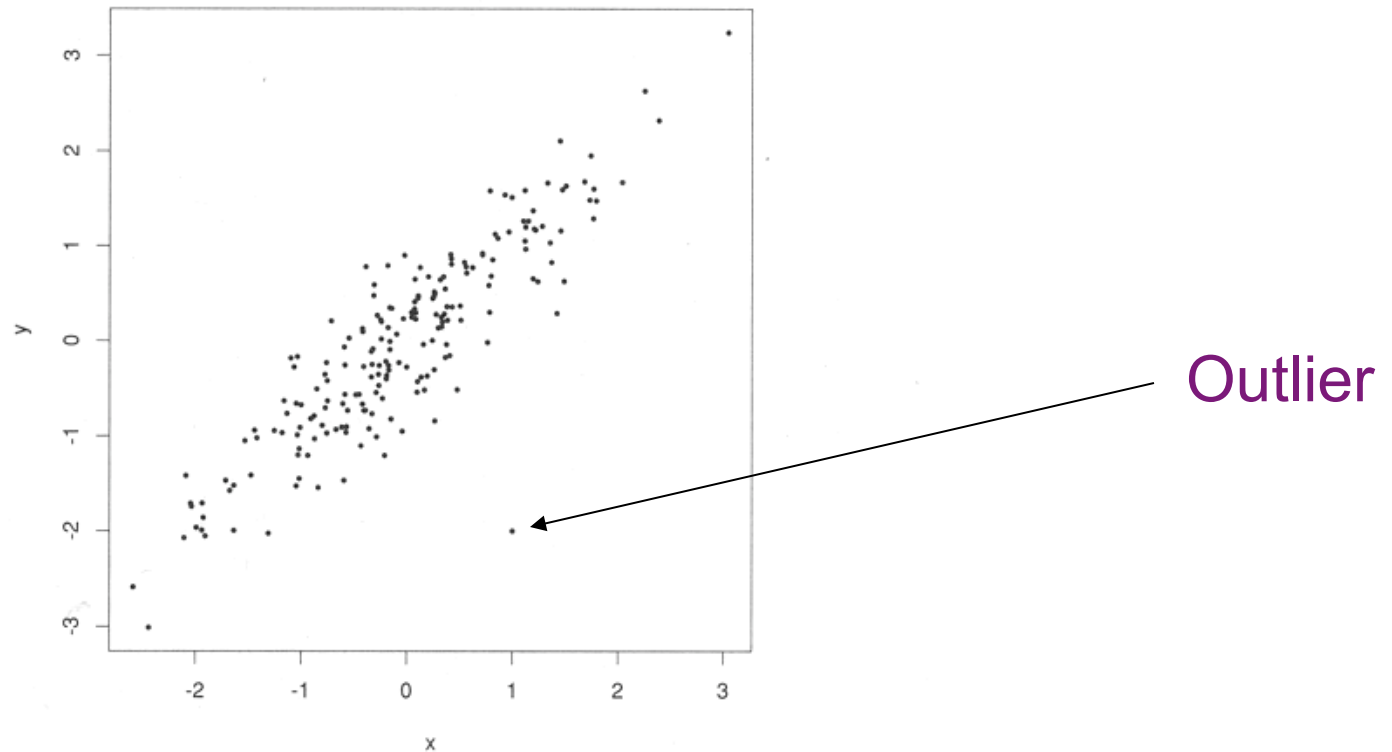


Figure 2.6 A plot of 200 points from highly positively correlated bivariate data (from a bivariate normal distribution), with a single easily identifiable outlier.