

Types of Data

Data Mining (2018, Fall Semester)

What is Data?

- ▶ Collection of data objects and their attributes
- ▶ An attribute is a property or characteristic of an object
 - ▶ Examples: eye color of a person, temperature, etc.
 - ▶ Attribute is also known as variable, field, characteristic, or feature
- ▶ A collection of attributes describe **Objects** an object
 - ▶ Object is also known as record, point, case, sample, entity, or instance

Attributes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Attribute Values

- ▶ Attribute values are numbers or symbols assigned to an attribute
- ▶ Distinction between attributes and attribute values
 - ▶ Same attribute can be mapped to different attribute values
 - ▶ Example: height can be measured in feet or meters
 - ▶ Different attributes can be mapped to the same set of values
 - ▶ Example: Attribute values for ID and age are integers
 - ▶ But properties of attribute values can be different
 - ID has no limit but age has a maximum and minimum value



Types of Attributes

- ▶ There are different types of attributes
 - ▶ **Nominal**
 - ▶ Examples: ID numbers, eye color, zip codes
 - ▶ **Ordinal**
 - ▶ Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
 - ▶ **Interval**
 - ▶ Examples: calendar dates, temperatures in Celsius or Fahrenheit.
 - ▶ **Ratio**
 - ▶ Examples: temperature in Kelvin, length, time, counts



Properties of Attribute Values

- ▶ The type of an attribute depends on which of the following properties it possesses:
 - ▶ Distinctness: = ≠
 - ▶ Order: < >
 - ▶ Addition: + -
 - ▶ Multiplication: * /
- ▶ Nominal attribute: distinctness
- ▶ Ordinal attribute: distinctness & order
- ▶ Interval attribute: distinctness, order & addition
- ▶ Ratio attribute: all 4 properties



Attribute Type	Description	Examples	Operations
Nominal	The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. (=, ≠)	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, χ^2 test
Ordinal	The values of an ordinal attribute provide enough information to order objects. (<, >)	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Interval	For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. (+, -)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests
Ratio	For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current	geometric mean, harmonic mean, percent variation



Discrete and Continuous Attributes

▶ Discrete Attribute

- ▶ Has only a finite or countably infinite set of values
- ▶ Examples: zip codes, counts, or the set of words in a collection of documents
- ▶ Often represented as integer variables.
- ▶ Note: binary attributes are a special case of discrete attributes

▶ Continuous Attribute

- ▶ Has real numbers as attribute values
- ▶ Examples: temperature, height, or weight.
- ▶ Practically, real values can only be measured and represented using a finite number of digits.
- ▶ Continuous attributes are typically represented as floating-point variables.



Types of data sets

▶ Record

- ▶ Data Matrix
- ▶ Document Data
- ▶ Transaction Data

▶ Graph

- ▶ World Wide Web
- ▶ Molecular Structures

▶ Ordered

- ▶ Spatial Data
- ▶ Temporal Data
- ▶ Sequential Data
- ▶ Genetic Sequence Data



Important Characteristics of Structured Data

- ▶ **Dimensionality**
 - ▶ **Curse of Dimensionality**
- ▶ **Sparsity**
 - ▶ **Only presence counts**
- ▶ **Resolution**
 - ▶ **Patterns depend on the scale**



Record Data

- ▶ Data that consists of a collection of records, each of which consists of a fixed set of attributes

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Data Matrix

- ▶ If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute
- ▶ Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1



Document Data

- ▶ Each document becomes a 'term' vector,
 - ▶ each term is a component (attribute) of the vector,
 - ▶ the value of each component is the number of times the corresponding term occurs in the document.

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0



Transaction Data

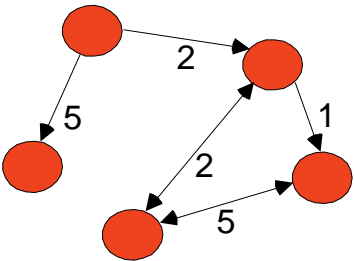
- ▶ A special type of record data, where
 - ▶ each record (transaction) involves a set of items.
 - ▶ For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk



Graph Data

- ▶ Examples: Generic graph and HTML Links



```
<a href="papers/papers.html#bbb">
Data Mining </a>
<li>
<a href="papers/papers.html#aaa">
Graph Partitioning </a>
<li>
<a href="papers/papers.html#aaa">
Parallel Solution of Sparse Linear System of Equations </a>
<li>
<a href="papers/papers.html#fff">
N-Body Computation and Dense Linear System Solvers
```



Ordered Data

► Genomic sequence data

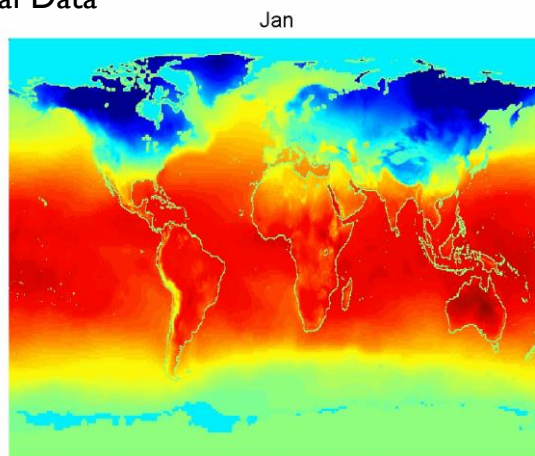
```
GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCCGCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG
```



Ordered Data

► Spatio-Temporal Data

Average Monthly
Temperature of land
and ocean



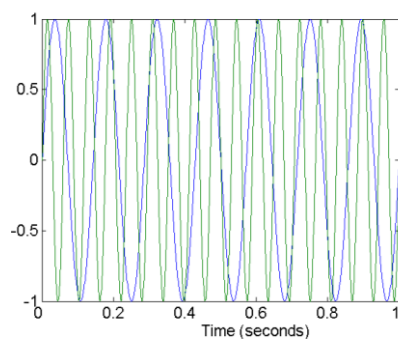
Data Quality

- ▶ What kinds of data quality problems?
- ▶ How can we detect problems with the data?
- ▶ What can we do about these problems?
- ▶ Examples of data quality problems:
 - ▶ Noise and outliers
 - ▶ missing values
 - ▶ duplicate data

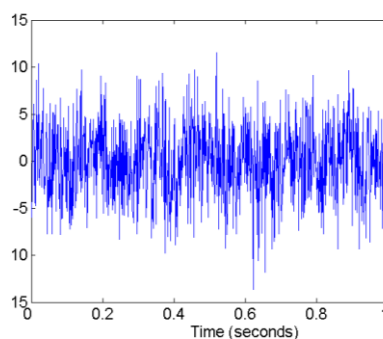


Noise

- ▶ Noise refers to modification of original values
 - ▶ Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen



Two Sine Waves

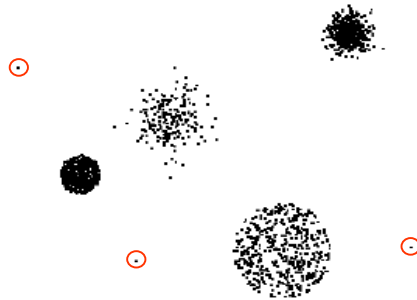


Two Sine Waves + Noise



Outliers

- ▶ Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



Missing Values

- ▶ **Reasons for missing values**
 - ▶ Information is not collected (e.g., people decline to give their age and weight)
 - ▶ Attributes may not be applicable to all cases (e.g., annual income is not applicable to children)
- ▶ **Handling missing values**
 - ▶ Eliminate Data Objects
 - ▶ Estimate Missing Values
 - ▶ Ignore the Missing Value During Analysis
 - ▶ Replace with all possible values (weighted by their probabilities)

Duplicate Data

- ▶ Data set may include data objects that are duplicates, or almost duplicates of one another
 - ▶ Major issue when merging data from heterogeneous sources
- ▶ Examples:
 - ▶ Same person with multiple email addresses
- ▶ Data cleaning
 - ▶ Process of dealing with duplicate data issues



Data Preprocessing

- ▶ Aggregation
- ▶ Sampling
- ▶ Dimensionality Reduction
- ▶ Feature subset selection
- ▶ Feature creation
- ▶ Discretization and Binarization
- ▶ Attribute Transformation



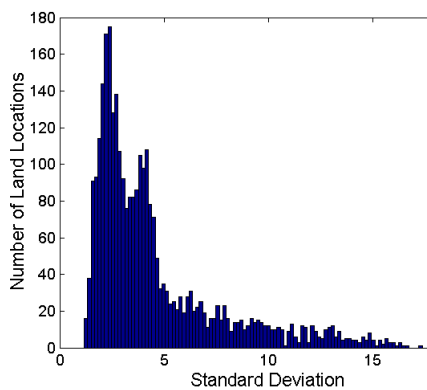
Aggregation

- ▶ Combining two or more attributes (or objects) into a single attribute (or object)
- ▶ Purpose
 - ▶ Data reduction
 - ▶ Reduce the number of attributes or objects
 - ▶ Change of scale
 - ▶ Cities aggregated into regions, states, countries, etc
 - ▶ More “stable” data
 - ▶ Aggregated data tends to have less variability

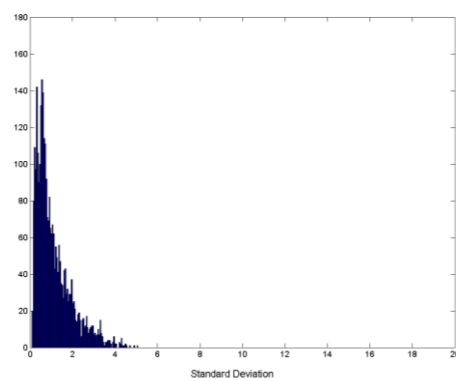


Aggregation

Variation of Precipitation in Australia



Standard Deviation of
Average Monthly
Precipitation



Standard Deviation of
Average Yearly Precipitation



Sampling

- ▶ **Sampling is the main technique employed for data selection.**
 - ▶ It is often used for both the preliminary investigation of the data and the final data analysis.
- ▶ **Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming.**
- ▶ **Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.**



Sampling ...

- ▶ **The key principle for effective sampling is the following:**
 - ▶ using a sample will work almost as well as using the entire data sets, if the sample is representative
 - ▶ A sample is representative if it has approximately the same property (of interest) as the original set of data



Types of Sampling

- ▶ **Simple Random Sampling**
 - ▶ There is an equal probability of selecting any particular item
- ▶ **Sampling without replacement**
 - ▶ As each item is selected, it is removed from the population
- ▶ **Sampling with replacement**
 - ▶ Objects are not removed from the population as they are selected for the sample.
 - ▶ In sampling with replacement, the same object can be picked up more than once
- ▶ **Stratified sampling**
 - ▶ Split the data into several partitions; then draw random samples from each partition



Sample Size



8000 points

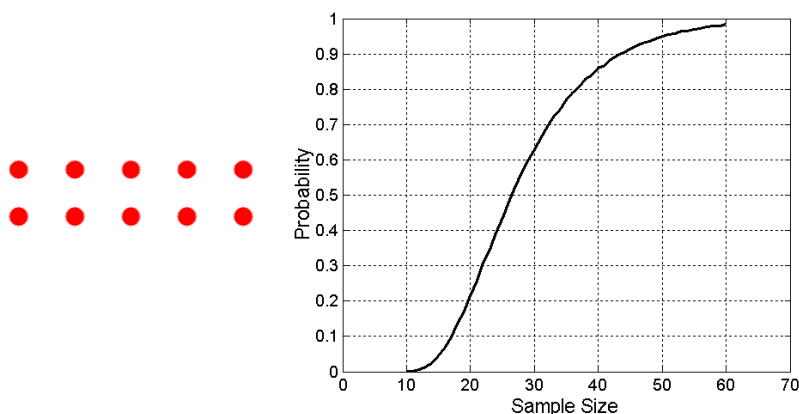
2000 Points

500 Points



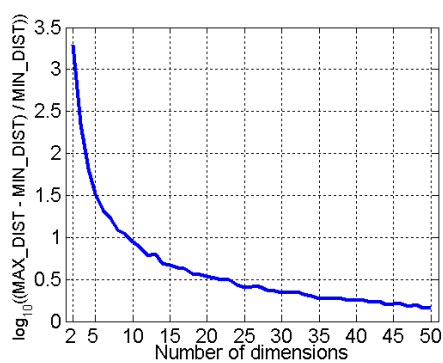
Sample Size

- ▶ What sample size is necessary to get at least one object from each of 10 groups.



Curse of Dimensionality

- ▶ When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- ▶ Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



- Randomly generate 500 points
- Compute difference between max and min distance between any pair of points

Dimensionality Reduction

► Purpose:

- Avoid curse of dimensionality
- Reduce amount of time and memory required by data mining algorithms
- Allow data to be more easily visualized
- May help to eliminate irrelevant features or reduce noise

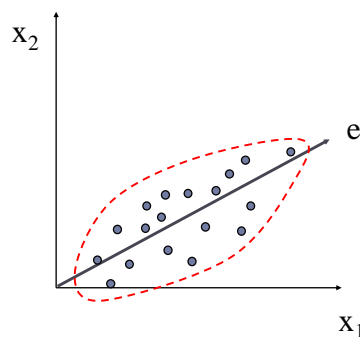
► Techniques

- Principle Component Analysis
- Singular Value Decomposition
- Others: supervised and non-linear techniques



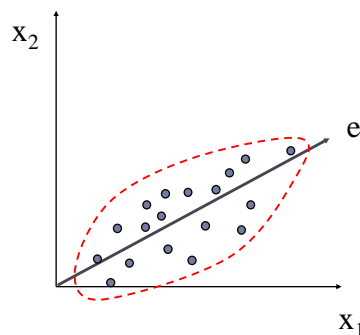
Dimensionality Reduction: PCA

- Goal is to find a projection that captures the largest amount of variation in data



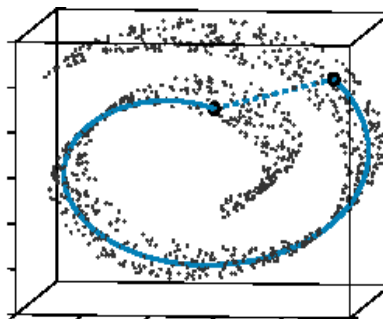
Dimensionality Reduction: PCA

- ▶ Find the eigenvectors of the covariance matrix
- ▶ The eigenvectors define the new space



Dimensionality Reduction: ISOMAP

By: Tenenbaum, de Silva,
Langford (2000)



- ▶ Construct a neighbourhood graph
- ▶ For each pair of points in the graph, compute the shortest path distances – geodesic distances



Feature Subset Selection

- ▶ Another way to reduce dimensionality of data
- ▶ Redundant features
 - ▶ duplicate much or all of the information contained in one or more other attributes
 - ▶ Example: purchase price of a product and the amount of sales tax paid
- ▶ Irrelevant features
 - ▶ contain no information that is useful for the data mining task at hand
 - ▶ Example: students' ID is often irrelevant to the task of predicting students' GPA



Feature Subset Selection

- ▶ Techniques:
 - ▶ Brute-force approach:
 - ▶ Try all possible feature subsets as input to data mining algorithm
 - ▶ Embedded approaches:
 - ▶ Feature selection occurs naturally as part of the data mining algorithm
 - ▶ Filter approaches:
 - ▶ Features are selected before data mining algorithm is run
 - ▶ Wrapper approaches:
 - ▶ Use the data mining algorithm as a black box to find best subset of attributes



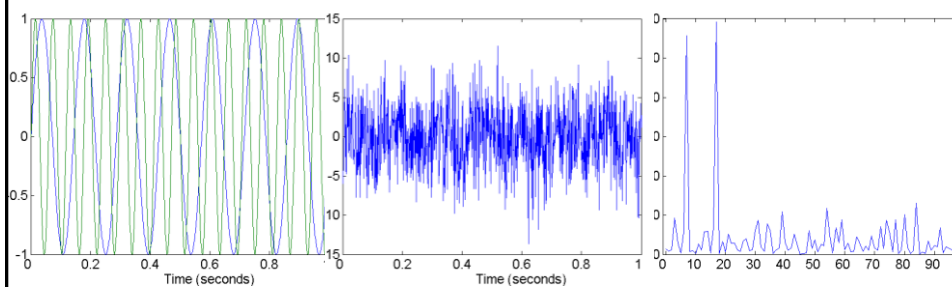
Feature Creation

- ▶ Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- ▶ Three general methodologies:
 - ▶ Feature Extraction
 - ▶ domain-specific
 - ▶ Mapping Data to New Space
 - ▶ Feature Construction
 - ▶ combining features



Mapping Data to a New Space

- Fourier transform
- Wavelet transform



Two Sine Waves

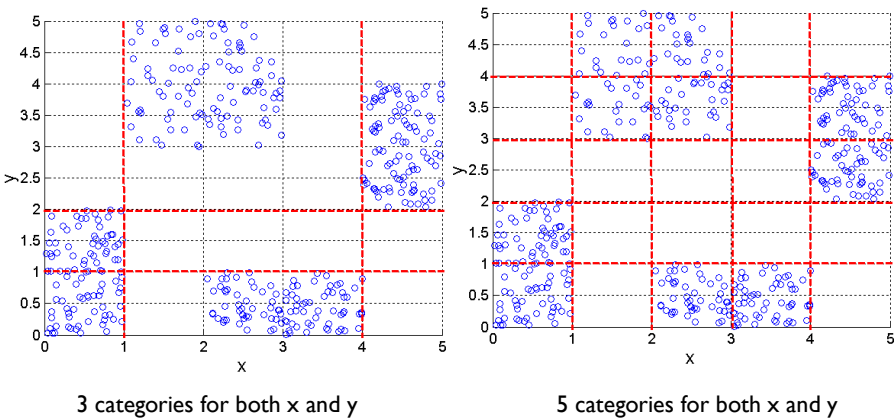
Two Sine Waves + Noise

Frequency

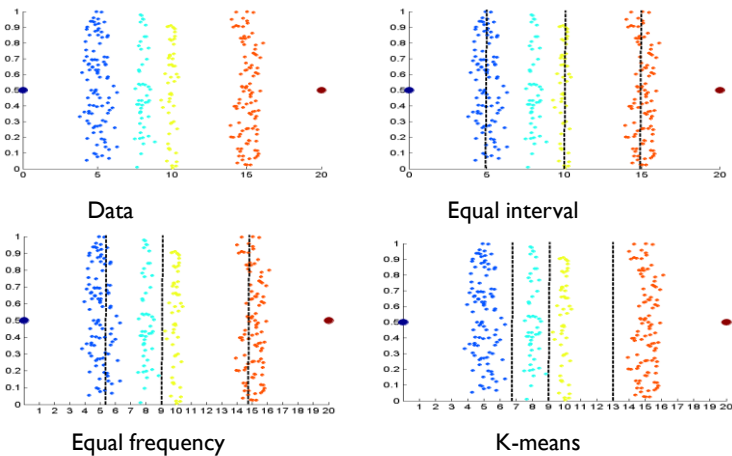


Discretization Using Class Labels

► **Entropy based approach**

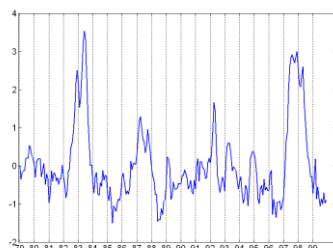


Discretization Without Using Class Labels



Attribute Transformation

- ▶ A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values
 - ▶ Simple functions: x^k , $\log(x)$, e^x , $|x|$
 - ▶ Standardization and Normalization



Similarity and Dissimilarity

- ▶ **Similarity**
 - ▶ Numerical measure of how alike two data objects are.
 - ▶ Is higher when objects are more alike.
 - ▶ Often falls in the range $[0,1]$
- ▶ **Dissimilarity**
 - ▶ Numerical measure of how different are two data objects
 - ▶ Lower when objects are more alike
 - ▶ Minimum dissimilarity is often 0
 - ▶ Upper limit varies
- ▶ Proximity refers to a similarity or dissimilarity



Similarity/Dissimilarity for Simple Attributes

p and q are the attribute values for two data objects.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$, where n is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d = p - q $	$s = -d, s = \frac{1}{1+d}$ or $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

Table 5.1. Similarity and dissimilarity for simple attributes



Euclidean Distance

- ▶ Euclidean Distance

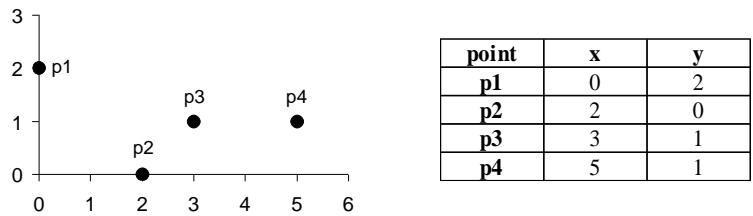
$$dist = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Where n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k^{th} attributes (components) or data objects p and q .

- ▶ Standardization is necessary, if scales differ.



Euclidean Distance



	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix



Minkowski Distance

- ▶ Minkowski Distance is a generalization of Euclidean Distance

$$dist = (\sum_{k=1}^n |p_k - q_k|^r)^{\frac{1}{r}}$$

Where r is a parameter, n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k th attributes (components) or data objects p and q .



Minkowski Distance: Examples

- ▶ $r = 1$. City block (Manhattan, taxicab, L_1 norm) distance.
 - ▶ A common example of this is the Hamming distance, which is just the number of bits that are different between two binary vectors
- ▶ $r = 2$. Euclidean distance
- ▶ $r \rightarrow \infty$. “supremum” (L_{\max} norm, L_{∞} norm) distance.
 - ▶ This is the maximum difference between any component of the vectors
- ▶ Do not confuse r with n , i.e., all these distances are defined for all numbers of dimensions.



Minkowski Distance

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

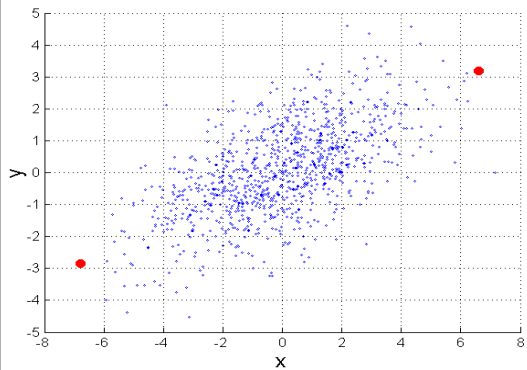
L_{∞}	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

Distance Matrix



Mahalanobis Distance

$$mahalanobis(p,q) = (p - q) \Sigma^{-1} (p - q)^T$$

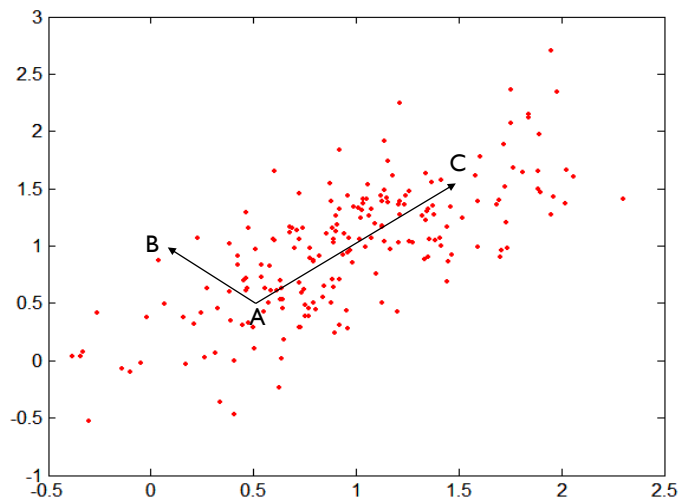


Σ is the covariance matrix of the input data X

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.

Mahalanobis Distance



Covariance Matrix:

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

Mahal(A,B) = 5

Mahal(A,C) = 4

Common Properties of a Distance

- ▶ Distances, such as the Euclidean distance, have some well known properties.

1. $d(p, q) \geq 0$ for all p and q and $d(p, q) = 0$ only if $p = q$. (Positive definiteness)
2. $d(p, q) = d(q, p)$ for all p and q . (Symmetry)
3. $d(p, r) \leq d(p, q) + d(q, r)$ for all points p, q , and r . (Triangle Inequality)

where $d(p, q)$ is the distance (dissimilarity) between points (data objects), p and q .

- ▶ A distance that satisfies these properties is a **metric**

Common Properties of a Similarity

- ▶ Similarities, also have some well known properties.

1. $s(p, q) = 1$ (or maximum similarity) only if $p = q$.
2. $s(p, q) = s(q, p)$ for all p and q . (Symmetry)

where $s(p, q)$ is the similarity between points (data objects), p and q .

Similarity Between Binary Vectors

- ▶ Common situation is that objects, p and q , have only binary attributes
- ▶ Compute similarities using the following quantities
 - M_{01} = the number of attributes where p was 0 and q was 1
 - M_{10} = the number of attributes where p was 1 and q was 0
 - M_{00} = the number of attributes where p was 0 and q was 0
 - M_{11} = the number of attributes where p was 1 and q was 1
- ▶ Simple Matching and Jaccard Coefficients
 - SMC = number of matches / number of attributes

$$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$
 - J = number of 11 matches / number of not-both-zero attributes values

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$



SMC versus Jaccard: Example

$p = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$
 $q = 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$

$M_{01} = 2$ (the number of attributes where p was 0 and q was 1)
 $M_{10} = 1$ (the number of attributes where p was 1 and q was 0)
 $M_{00} = 7$ (the number of attributes where p was 0 and q was 0)
 $M_{11} = 0$ (the number of attributes where p was 1 and q was 1)

$$SMC = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0+7) / (2+1+0+7) = 0.7$$

$$J = (M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$



Cosine Similarity

- ▶ If d_1 and d_2 are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\| ,$$

where \bullet indicates vector dot product and $\|d\|$ is the length of vector d .

- ▶ Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$d_1 \bullet d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$\|d_1\| = (3*3 + 2*2 + 0*0 + 5*5 + 0*0 + 0*0 + 0*0 + 2*2 + 0*0 + 0*0)^{0.5} = (42)^{0.5} = 6.481$$

$$\|d_2\| = (1*1 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 1*1 + 0*0 + 2*2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$



Correlation

- ▶ Correlation measures the linear relationship between objects
- ▶ To compute correlation, we standardize data objects, p and q , and then take their dot product

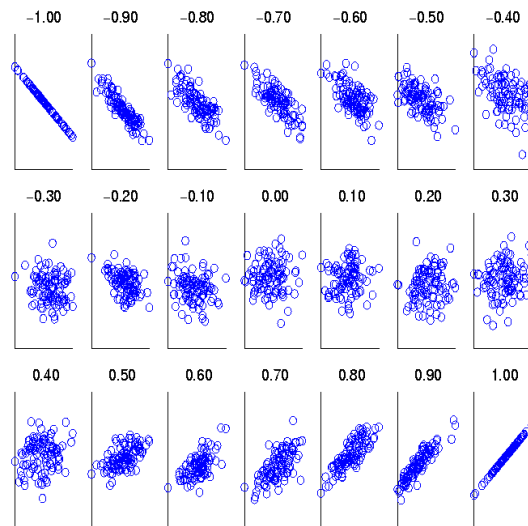
$$p'_k = (p_k - \text{mean}(p)) / \text{std}(p)$$

$$q'_k = (q_k - \text{mean}(q)) / \text{std}(q)$$

$$\text{correlation}(p, q) = p' \bullet q'$$



Visually Evaluating Correlation



Scatter plots showing the similarity from -1 to 1.

General Approach for Combining Similarities

- Sometimes attributes are of many different types, but an overall similarity is needed.

1. For the k^{th} attribute, compute a similarity, s_k , in the range $[0, 1]$.
2. Define an indicator variable, δ_k , for the k^{th} attribute as follows:

$$\delta_k = \begin{cases} 0 & \text{if the } k^{th} \text{ attribute is a binary asymmetric attribute and both objects have a value of 0, or if one of the objects has a missing values for the } k^{th} \text{ attribute} \\ 1 & \text{otherwise} \end{cases}$$

3. Compute the overall similarity between the two objects using the following formula:

$$similarity(p, q) = \frac{\sum_{k=1}^n \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

Using Weights to Combine Similarities

- ▶ May not want to treat all attributes the same.
 - ▶ Use weights w_k which are between 0 and 1 and sum to 1.

$$similarity(p, q) = \frac{\sum_{k=1}^n w_k \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

$$distance(p, q) = \left(\sum_{k=1}^n w_k |p_k - q_k|^r \right)^{1/r}.$$



Density

- ▶ Density-based clustering require a notion of density
- ▶ Examples:
 - ▶ Euclidean density
 - ▶ Euclidean density = number of points per unit volume
 - ▶ Probability density



Euclidean Density – Cell-based

- Simplest approach is to divide region into a number of rectangular cells of equal volume and define density as # of points the cell contains

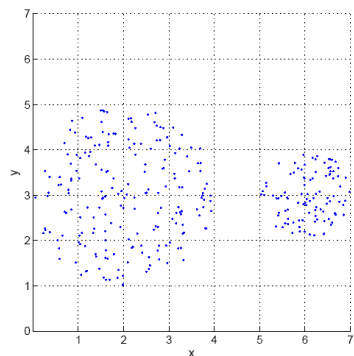


Figure 7.13. Cell-based density.

0	0	0	0	0	0	0
0	0	0	0	0	0	0
4	17	18	6	0	0	0
14	14	13	13	0	18	27
11	18	10	21	0	24	31
3	20	14	4	0	0	0
0	0	0	0	0	0	0

Table 7.6. Point counts for each grid cell.

Euclidean Density – Center-based

- Euclidean density is the number of points within a specified radius of the point

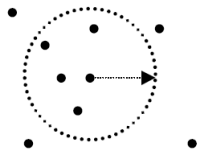


Figure 7.14. Illustration of center-based density.