Why Mine Data? Commercial Viewpoint

- Lots of data is being collected and warehoused
 - Web data, e-commerce
 - purchases at department/ grocery stores
 - Bank/Credit Card transactions



- Computers have become cheaper and more powerful
- Competitive Pressure is Strong
 - Provide better, customized services for an *edge* (e.g. in Customer Relationship Management)

Why Mine Data? Scientific Viewpoint

- Data collected and stored at enormous speeds (GB/hour)
 - remote sensors on a satellite
 - telescopes scanning the skies
 - microarrays generating gene expression data
 - scientific simulations
 generating terabytes of data
- Traditional techniques infeasible for raw data
- Data mining may help scientists
 - in classifying and segmenting data
 - in Hypothesis Formation



Mining Large Data Sets - Motivation

- There is often information "hidden" in the data that is not readily evident
- Human analysts may take weeks to discover useful information
- Much of the data is never analyzed at all.

What is Data Mining?

Many Definitions

 Non-trivial extraction of implicit, previously unknown and potentially useful information from data



What is (not) Data Mining?

What is not Data Mining?

Look up phone
 number in phone
 directory

Query a Web
search engine for
information about
"Amazon"

• What is Data Mining?

Certain names are more prevalent in certain locations in Particular area)

 Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest, Amazon.com,)

Origins of Data Mining

- Draws ideas from machine learning/Al, pattern recognition, statistics, and database systems
- Traditional Techniques may be unsuitable due to
 - Enormity of data
 - High dimensionality of data
 - Heterogeneous, distributed nature of data



Data Mining Tasks

Prediction Methods

 Use some variables to predict unknown or future values of other variables.

- Description Methods
 - Find human-interpretable patterns that describe the data.

Data Mining Tasks...

- Classification [Predictive]
- Clustering [Descriptive]
- Association Rule Discovery [Descriptive]
- Sequential Pattern Discovery [Descriptive]
- Regression [Predictive]
- Deviation Detection /anomaly[Predictive]

Classification: Definition

- Given a collection of records (training set)
 - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: <u>previously unseen</u> records should be assigned a class as accurately as possible.
 - A test set is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

Classification Example



Refund	Marital Status	Taxable Income	Cheat		
No	Single	75K	?		
Yes	Married	50K	?		
No	Married	150K	?	Δ.	
Yes	Divorced	90K	?		
No	Single	40K	?		
No	Married	80K	?		Test

Direct Marketing

- Goal: Reduce cost of mailing by *targeting* a set of consumers likely to buy a new cell-phone product.
- Approach:
 - Use the data for a similar product introduced before.
 - We know which customers decided to buy and which decided otherwise. This {buy, don't buy} decision forms the class attribute.
 - Collect various demographic, lifestyle, and companyinteraction related information about all such customers.
 - Type of business, where they stay, how much they earn, etc.
 - Use this information as input attributes to learn a classifier model.

Fraud Detection

- Goal: Predict fraudulent cases in credit card transactions.
- Approach:
 - Use credit card transactions and the information on its account-holder as attributes.
 - When does a customer buy, what does he buy, how often he pays on time, etc
 - Label past transactions as fraud or fair transactions. This forms the class attribute.
 - Learn a model for the class of the transactions.
 - Use this model to detect fraud by observing credit card transactions on an account.

- Customer Attrition/Churn:
 - Goal: To predict whether a customer is likely to be lost to a competitor.
 - Approach:
 - Use detailed record of transactions with each of the past and present customers, to find attributes.
 - How often the customer calls, where he calls, what time-of-the day he calls most, his financial status, marital status, etc.
 - Label the customers as loyal or disloyal.
 - Find a model for loyalty.

Sky Survey Cataloging

 Goal: To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images (from Palomar Observatory).

- 3000 images with 23,040 x 23,040 pixels per image.

- Approach:
 - Segment the image.
 - Measure image attributes (features) 40 of them per object.
 - Model the class based on these features.
 - Success Story: Could find 16 new high red-shift quasars, some of the farthest objects that are difficult to find!

Classifying Galaxies

Early



Class:

• Stages of Formation

Attributes:

- Image features,
- Characteristics of light waves received, etc.



Data Size:

- 72 million stars, 20 million galaxies
- Object Catalog: 9 GB
- Image Database: 150 GB

Late



Clustering Definition

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that
 - Data points in one cluster are more similar to one another.
 - Data points in separate clusters are less similar to one another.
- Similarity Measures:
 - Euclidean Distance if attributes are continuous.
 - Other Problem-specific Measures.

Illustrating Clustering

⊠Euclidean Distance Based Clustering in 3-D space.

Intracluster distances Intercluster distances are minimized are maximized

Clustering: Application 1

- Market Segmentation:
 - Goal: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with a distinct marketing mix.
 - Approach:
 - Collect different attributes of customers based on their geographical and lifestyle related information.
 - Find clusters of similar customers.
 - Measure the clustering quality by observing buying patterns of customers in same cluster vs. those from different clusters.

Clustering: Application 2

- Document Clustering:
 - Goal: To find groups of documents that are similar to each other based on the important terms appearing in them.
 - Approach: To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.
 - Gain: Information Retrieval can utilize the clusters to relate a new document or search term to clustered documents.

Illustrating Document Clustering

- Clustering Points: 3204 Articles of Los Angeles Times.
- Similarity Measure: How many words are common in these documents (after some word filtering).

Category	Total Articles	Correctly Placed
Financial	555	364
Foreign	341	260
National	273	36
Metro	943	746
Sports	738	573
Entertainment	354	278

Clustering of S&P 500 Stock Data

- **#** Observe Stock Movements every day.
- **Clustering points:** Stock-{UP/DOWN}
- Similarity Measure: Two points are more similar if the events described by them frequently happen together on the same day.

We used association rules to quantify a similarity measure.

	Discovered Clusters	Industry Group
1	Applied-Matl-DOW N, Bay-Net work-Down, 3-COM-DOW N, Cabletron-Sys-DOW N, CISCO-DOW N, HP-DOW N, DSC-Comm-DOW N, INTEL-DOW N, LSI-Logic-DOW N, Micron-Tech-DOW N, Te xas-Inst-Down, Te llabs-Inc-Down, Natl-Se miconduct-DOW N, Orac l-DOW N, SGI-DOW N, Sun-DOW N	Technology1-DOWN
2	Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN	Technology2-DOWN
3	Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN	Financial-DOWN
4	Baker-Hughes-UP, Dresser-Inds-UP, Halliburton-HLD-UP, Louisiana-Land-UP, Phillips-Petro-UP, Unocal-UP, Schlumberger-UP	Oil-UP

Association Rule Discovery: Definition

- Given a set of records each of which contain some number of items from a given collection;
 - Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

TID	Items	
1	Bread, Coke, Milk	Rules Discovered:
2	Beer, Bread	{Milk}> {Coke}
3	Beer, Coke, Diaper, Milk	{Diaper, Milk}> {Beer}
4	Beer, Bread, Diaper, Milk	
5	Coke, Diaper, Milk	

Association Rule Discovery: Application 1

- Marketing and Sales Promotion:
 - Let the rule discovered be

{Bagels, ... } --> {Potato Chips}

- <u>Potato Chips as consequent</u> => Can be used to determine what should be done to boost its sales.
- <u>Bagels in the antecedent</u> => Can be used to see which products would be affected if the store discontinues selling bagels.
- Bagels in antecedent and Potato chips in consequent
 Scan be used to see what products should be sold with Bagels to promote sale of Potato chips!

Association Rule Discovery: Application 2

- Supermarket shelf management.
 - Goal: To identify items that are bought together by sufficiently many customers.
 - Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.
 - A classic rule --
 - If a customer buys diaper and milk, then he is very likely to buy beer.
 - So, don't be surprised if you find six-packs stacked next to diapers!

Association Rule Discovery: Application 3

- Inventory Management:
 - Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with right parts to reduce on number of visits to consumer households.
 - Approach: Process the data on tools and parts required in previous repairs at different consumer locations and discover the co-occurrence patterns.

Sequential Pattern Discovery: Definition

Given is a set of *objects*, with each object associated with its own *timeline of events*, find rules that predict strong sequential dependencies among different events.

$$(\mathbf{A} \ \mathbf{B}) \quad (\mathbf{C}) \longrightarrow (\mathbf{D} \ \mathbf{E})$$

 Rules are formed by first disovering patterns. Event occurrences in the patterns are governed by timing constraints.



Sequential Pattern Discovery: Examples

- In telecommunications alarm logs,
 - (Inverter_Problem Excessive_Line_Current)

(Rectifier_Alarm) --> (Fire_Alarm)

- In point-of-sale transaction sequences,
 - Computer Bookstore:

(Intro_To_Visual_C) (C++_Primer) -->

(Perl_for_dummies,Tcl_Tk)

- Athletic Apparel Store:

(Shoes) (Racket, Racketball) --> (Sports_Jacket)

Regression

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Greatly studied in statistics, neural network fields.
- Examples:
 - Predicting sales amounts of new product based on advetising expenditure.
 - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
 - Time series prediction of stock market indices.

Challenges of Data Mining

- Scalability
- Dimensionality
- Complex and Heterogeneous Data
- Data Quality
- Data Ownership and Distribution
- Privacy Preservation
- Streaming Data

Data : 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 an object
 - Object is also known as record, point, case, sample, entity, or instance



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

Measurement of Length

 The way you measure an attribute is somewhat may not match the attributes properties.



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: = \neq
 - 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. $(=, \neq)$	<pre>zip codes, employee ID numbers, eye color, sex: {male, female}</pre>	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, better, 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

Attribute Level	Transformation	Comments
Nominal	Any permutation of values	If all employee ID numbers were reassigned, would it make any difference?
Ordinal	An order preserving change of values, i.e., <i>new_value</i> = <i>f</i> (<i>old_value</i>) where <i>f</i> is a monotonic function.	An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by { 0.5, 1, 10}.
Interval	$new_value = a * old_value + b$ where a and b are constants	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	<i>new_value</i> = <i>a</i> * <i>old_value</i>	Length can be measured in meters or feet.
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

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

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	pla y	ball	score	game	wi n	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



 Data Mining Graph Partitioning Parallel Solution of Sparse Linear System of Equations N-Body Computation and Dense Linear System Solvers

Chemical Data

• Benzene Molecule: C₆H₆



Ordered Data

Sequences of transactions

Items/Events (AB) (D) (CE) (BD) (C) (E) (CD) (B) (AE) An element of the sequence

Ordered Data

Genomic sequence data

GGTTCCGCCTTCAGCCCGCGCCC CGCAGGGCCCGCCCGCGCGCGCG GAGAAGGGCCCGCCTGGCGGGGCG GGGGGAGGCGGGGGCCGCCCGAGC CCAACCGAGTCCGACCAGGTGCC CCCTCTGCTCGGCCTAGACCTGA GCTCATTAGGCGGCAGCGGACAG GCCAAGTAGAACACGCGAAGCGC

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



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 heterogeous 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



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

Dimensionality Reduction: PCA



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 approch:
 - 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



3 categories for both x and y

5 categories for both x and y
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	Dissimilarity	Similarity	
Type			
Nominal	$d = \left\{ egin{array}{cc} 0 & ext{if} \; p = q \ 1 & ext{if} \; p eq q \end{array} ight.$	$s = \left\{ egin{array}{ccc} 1 & ext{if} \; p = q \ 0 & ext{if} \; p eq q \end{array} ight.$	
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	$s = -d, \ s = rac{1}{1+d} \ ext{or}$	
		$s = 1 - rac{d-min_d}{max_d-min_d}$	

 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 kth attributes (components) or data objects *p* and *q*.

• Standardization is necessary, if scales differ.

Euclidean Distance



point	X	У
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	р3	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

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

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

Minkowski Distance: Examples

- r = 1. City block (Manhattan, taxicab, L₁ 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
- $r \to \infty$. "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

L1	p1	p2	p3	p4		
p1	0	4	4	6		
p2	4	0	2	4		
р3	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		
\mathbf{L}_{∞}	p1	p2	p3	p4		
p1	0	2	3	5		
p2	2	0	1	3		
p3	3	1	0	2		
p4	5	3	2			

point	X	У
p1	0	2
p2	2	0
p3	3	1
p4	5	1

Distance Matrix

Mahalanobis Distance

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



 Σ is the covariance matrix of the input data *X*

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

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

Mahalanobis Distance



Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.
 - 1. $d(p, q) \ge 0$ 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. $d(p, r) \le d(p, q) + d(q, r)$ 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
 M₀₁ = the number of attributes where p was 0 and q was 1
 M₁₀ = the number of attributes where p was 1 and q was 0
 M₀₀ = the number of attributes where p was 0 and q was 0
 M₁₁ = 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 = 1000000000q = 0000001001

 $M_{01} = 2$ (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:

 $d_1 = 3205000200$ $d_2 = 100000102$

 $\begin{aligned} &d_1 \bullet d_2 = 3^*1 + 2^*0 + 0^*0 + 5^*0 + 0^*0 + 0^*0 + 0^*0 + 2^*1 + 0^*0 + 0^*2 = 5 \\ &||d_1|| = (3^*3 + 2^*2 + 0^*0 + 5^*5 + 0^*0 + 0^*0 + 0^*0 + 2^*2 + 0^*0 + 0^*0)^{0.5} = (42)^{0.5} = 6.481 \\ &||d_2|| = (1^*1 + 0^*0 + 0^*0 + 0^*0 + 0^*0 + 0^*0 + 1^*1 + 0^*0 + 2^*2)^{0.5} = (6)^{0.5} = 2.245 \end{aligned}$

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

Extended Jaccard Coefficient (Tanimoto)

- Variation of Jaccard for continuous or count attributes
 - Reduces to Jaccard for binary attributes

$$T(p,q) = \frac{p \bullet q}{\|p\|^2 + \|q\|^2 - p \bullet q}$$

Correlation

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

$$p'_k = (p_k - mean(p)) / std(p)$$

$$q'_k = (q_k - mean(q)) / std(q)$$

 $correlation(p,q) = p' \bullet q'$

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:

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

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

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

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.

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

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

Density

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



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.