Data preprocessing



Elena Baralis, Tania Cerquitelli Politecnico di Torino



- Data types and properties
- Data preparation
- Data preparation for document data
- Similarity and dissimilarity
- Correlation



Data types and properties





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, dimension, or feature
- A collection of attributes describes 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 for a particular object
- 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



Attribute types

- 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
 - 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



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.





More Complicated Examples

- ID numbers
 - Nominal, ordinal, or interval?

- Number of cylinders in an automobile engine
 - Nominal, ordinal, or ratio?



Key Messages for Attribute Types

- The types of operations you choose should be "meaningful" for the type of data you have
 - Distinctness, order, meaningful intervals, and meaningful ratios are only four properties of data
 - The data type you see often numbers or strings may not capture all the properties or may suggest properties that are not there
 - Analysis may depend on these other properties of the data
 - Many statistical analyses depend only on the distribution
 - Many times what is meaningful is measured by statistical significance
 - But in the end, what is meaningful is measured by the domain





Data set types

Record

- Tables
- Document Data
- Transaction Data
- Graph
 - World Wide Web
 - Molecular Structures
- Ordered
 - Spatial Data
 - Temporal Data
 - Sequential Data
 - Genetic Sequence Data





Tabular Data

A collection of records

 Each record is characterized by a fixed set of attributes
Tid Refund Marital Taxable

I Id	Refund	Marital Status	l axable 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





Document data

- It includes textual data that can be semistructured or unstructured
 - Plain text can be organized in sentences, paragraphs, sections, documents
- Text acquired in different contexts may have a structure and/or a semantics
 - Web pages are enriched with tags
 - Documents in digital libraries are enriched with metadata
 - E-learning documents can be annotated or partly highglihted





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





Examples: Generic graph, a molecule, and webpages





Sequences of transactions



An element of the sequence





Ordered Data

Genomic sequence data

GGTTCCGCCTTCAGCCCGCGCCC CGCAGGGCCCGCCCGCGCGCGCG GAGAAGGGCCCGCCTGGCGGGGCG GGGGGAGGCGGGGGCCGCCCGAGC CCAACCGAGTCCGACCAGGTGCC CCCTCTGCTCGGCCTAGACCTGA GCTCATTAGGCGGCAGCGGACAG GCCAAGTAGAACACGCGAAGCGC





Ordered Data

Spatio-Temporal Data

Jan

Average Monthly Temperature of land and ocean





Data Quality

Poor data quality negatively affects many data processing efforts

"The most important point is that poor data quality is an unfolding disaster. Poor data quality costs the typical company at least ten percent (10%) of revenue; twenty percent (20%) is probably a better estimate." Thomas C. Redman, DM Review, August 2004

- Data mining example: a classification model for detecting people who are loan risks is built using poor data
 - Some credit-worthy candidates are denied loans
 - More loans are given to individuals that default





- 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
 - Wrong data



From: Tan, Steinbach, Karpatne, Kumar, Introduction to Data Mining 2nd Ed., McGraw Hill 2018



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

- Case 1: Outliers are noise that interferes with data analysis
- Case 2: Outliers are the goal of our analysis
 - Credit card fraud
 - Intrusion detection





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 or variables
- Estimate missing values
 - Example: time series of temperature
- Ignore the missing value during analysis





- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources
- Examples
 - Different words/abbreviations for the same concept (e.g., Street, St.)
- Data cleaning
 - Process of dealing with duplicate data issues
- When should duplicate data not be removed?



From: Tan, Steinbach, Karpatne, Kumar, Introduction to Data Mining 2nd Ed., McGraw Hill 2018

Data preparation





Data Preprocessing

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





- 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



Example: Precipitation in Australia

- The next slide shows precipitation in Australia from the period 1982 to 1993
 - A histogram for the standard deviation of average monthly precipitation for 3,030 0.5° by 0.5° grid cells in Australia
 - A histogram for the standard deviation of the average yearly precipitation for the same locations.
- The average yearly precipitation has less variability than the average monthly precipitation.
- All precipitation measurements (and their standard deviations) are in centimeters.





Variation of Precipitation in Australia





Data reduction

- Generates a reduced representation of the dataset
- This representation is smaller in volume, but it can provide similar analytical results
 - sampling
 - reduces the cardinality of the set
 - feature selection
 - reduces the number of attributes
 - discretization
 - reduces the cardinality of the attribute domain





- 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.





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





Sample Size: examples



8000 points

2000 Points

500 Points



From: Tan, Steinbach, Karpatne, Kumar, Introduction to Data Mining 2nd Ed., McGraw Hill 2018



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





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

- Principal Component Analysis (PCA)
- Singular Value Decomposition
- Others: supervised and non-linear techniques




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





From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006



Dimensionality Reduction: PCA

256







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 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 transformWavelet transform



From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006



Discretization

- Discretization is the process of converting a continuous attribute into an ordinal attribute
 - A potentially infinite number of values are mapped into a small number of categories
 - Discretization is commonly used in classification
 - Many classification algorithms work best if both the independent and dependent variables have only a few values



Iris Sample Data Set

- Iris Plant data set
 - Can be obtained from the UCI Machine Learning Repository <u>http://www.ics.uci.edu/~mlearn/MLRepository.html</u>
 - From the statistician Douglas Fisher
 - Three flower types (classes)
 - Setosa
 - Versicolour
 - Virginica
 - Four (non-class) attributes
 - Sepal width and length
 - Petal width and length



Virginica. Robert H. Mohlenbrock. USDA NRCS. 1995. Northeast wetland flora: Field office guide to plant species. Northeast National Technical Center, Chester, PA. Courtesy of USDA NRCS Wetland Science Institute.



Discretization: Iris Example



Petal width low or petal length low implies Setosa. Petal width medium or petal length medium implies Versicolour. Petal width high or petal length high implies Virginica.

Discretization: Iris Example ...

How can we tell what the best discretization is?

- Unsupervised discretization: find breaks in the data values
 - Example: Petal Length



Supervised discretization: Use class labels to find breaks



Discretization

Examples of unsupervised discretization techniques

- N intervals with the same width $W = (v_{max} v_{min})/N$
 - Easy to implement
 - It can be badly affected by outliers and sparse data
 - Incremental approach
- N intervals with (approximately) the same cardinality
 - It better fits sparse data and outliers
 - Non incremental approach
- clustering
 - It fits well sparse data and outliers
- analysis of data distribution
 - e.g., 4 intervals, one for each quartile







Data consists of four groups of points and two outliers. Data is onedimensional, but a random y component is added to reduce overlap.







Equal interval width approach used to obtain 4 values.







Equal frequency approach used to obtain 4 values.







K-means approach to obtain 4 values.



Binarization

- Binarization maps an attribute into one or more binary variables
- Continuous attribute: first map the attribute to a categorical one
 - Example: height measured as {low, medium, high}
- Categorical attribute
 - Mapping to a set of binary attributes
 - Example: Low, medium, high as 1 0 0, 0 1 0, 0 0 1
 - One-hot encoding
 - Only 1 bit takes value 1
 - It represents the specific value taken by the attribute



Attribute Transformation

- An attribute transform is 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|
- Normalization
 - Refers to various techniques to adjust to differences among attributes in terms of frequency of occurrence, mean, variance, range
 - Take out unwanted, common signal, e.g., seasonality
- In statistics, standardization refers to subtracting off the means and dividing by the standard deviation





Normalization

It is a type of data transformation

- The values of an attribute are scaled so as to fall within a small specified range, typically [-1,+1] or [0,+1]
- Techniques
 - min-max normalization

$$v' = \frac{v - min_{A}}{max_{A} - min_{A}} (new max_{A} - new min_{A}) + new min_{A}$$

z-score normalization

$$v' = \frac{v - mean_A}{stand _ dev_A}$$

decimal scaling

 $v' = \frac{v}{10^{j}}$ j is the smallest integer such that max(|v'|) < 1



Example: Sample Time Series of Plant Growth



Net Primary Production (NPP) is a measure of plant growth used by ecosystem scientists.

	Minneapolis	Atlanta	Sao Paolo		
Minneapolis	1.0000	0.7591	-0.7581		
Atlanta	0.7591	1.0000	-0.5739		
Sao Paolo	-0.7581	-0.5739	1.0000		



Example: Sample Time Series of Plant Growth



Normalized using monthly Z Score:

Subtract off monthly mean and divide by monthly standard deviation

Correlations between time series

	Minneapolis	Atlanta	Sao Paolo
Minneapolis	1.0000	0.0492	0.0906
Atlanta	0.0492	1.0000	-0.0154
Sao Paolo	0.0906	-0.0154	1.0000



Data preparation for document data





Document representation

- A document might be modeled in different ways
 - The choice heavily affects the quality of the mining result
- The most common representation models a document as a set of features
 - Each feature might represent a set of characters, a word, a term, a concept





Document processing

- It is the activity to generate a structured data representation of document data
- It includes five sequential steps
 - Document splitting
 - Tokenisation
 - Case normalisation
 - Stopword removal
 - Stemming





Document splitting

- Based on the data analytics goal, documents can be split into
 - sentences, paragraphs, or analyzed in their entire content
- Short documents are typically not split
 - e.g., emails or social posts
- Long documents can be
 - broken up into sections or paragraphs
 - analyzed as a whole





Tokenization

- It is the process of breaking text into sentences or text into tokens (i.e., words)
 - Identify sentence boundaries based on punctuation, capitalization
 - Separate words in sentences
 - Language-dependent



Case normalization

- This step converts each token to completely upper-case or lower-case characters
 - Capitalisation helps human readers differentiate, for example, between nouns and proper nouns and can be useful for automated algorithms as well
 - However, an upper-case word at the beginning of the sentence should be treated no differently than the same word in lower case appearing elsewhere in a document



Stemming

Reduce a word to its root form (i.e., the **stem**)

- It includes the identification and removal of prefixes, suffixes, and pluralisation
- It operates on a single word without knowledge of the context
 - It cannot discriminate between words which have different meanings depending on the part of speech

Stemmers are

- Easy to implement
- Available for most spoken languages
- Run significantly faster than lemmatization and POS tagging algorithms



Stopword elimination

- "Stop words" refers to the most common words in a language
 - E.g., prepositions, articles, conjunctions in English
- Stop words are filtered out before or after processing of textual data
 - They are likely to have little semantic meaning



Stopword elimination

- There is no single universal list of stop words used by all natural language processing tools
- Any group of words can be chosen as the stop words for a given purpose
 - different search engines use different stop word lists
 - Some of them remove lexical words, such as "want", from a query in order to improve performance
- Some tools specifically avoid removing these stop words to support phrase search



Weighted document representation



Text representation: feature vectors

- Most data mining algorithms are unable to directly process textual data in their original form
 - documents are transformed into a more manageable representation
- Documents are represented by feature vectors
- A feature is simply an entity without internal structure
 - A dimension of the feature space
- A document is represented as a vector in this space
 - a collection of features and their weights





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	w <u>i</u> 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



From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006

Bag-of-word representation

- All words in a document are considered as separate features
 - the dimension of the feature space is equal to the number of different words in the entire document collection
- The feature vector of a document consists of a set of weights, one for each distinct word
- The methods for giving weights to the features may vary



Weighting schemes

- Binary
 - One, if the corresponding word is present in the document
 - Zero, otherwise
 - Occurrences of all words have the same importance
- Simple document frequency
 - The number of times in which the corresponding word occurs in the document
 - Most frequent words are not always representative of the document content



Weighting schemes

- More complex weighting schemes are possible to take into account the frequency of the word
 - in the document
 - in the section/paragraph
 - in the category (for indexed documents)
 - in the collection of documents



Weighting schemes

Term frequency inverse document frequency (tf-idf)

Tf-idf of term t in document d of collection D (consisting of m documents)

tf-idf(t) = freq(t, d) * log(m/freq(t, D))

 Terms occurring frequently in a single document but rarely in the whole collection are preferred

Suitable for

- A single document consisting of many sections or subsections
- A collection of *heterogeneous* documents


Weighting schemes

- Term frequency document frequency (tf-df)
 - Tf-df of term t in document d of collection D

tf-df(t) = freq(t, d) * log(freq(t, D))

 Terms occurring frequently both in a single document and in the whole collection are preferred

Suitable for

- Single documents or parts of a document with homogeneous content
- A collection of documents ranging over the same topic





Tf-idf matrix example

												Ļ									
major	malform	materi	matric	matrix	mean	measur	mechan	medicin	medium	medlin	method	methodolog	g micro	microarch	. migrat	mo	model	molecular	morbid	moreov	mortal
0	0	0.153	0.051	0.021	0	0	0	0	0	0	0.051	0.069	0.072	0	0.020	0	0.034	0.072	0	0.072	0.063
0.032	0.032	0.048	0.032	0.020	0.032	0.032	0.032	0.064	0.032	0.032	0.048	0.043	0.023	0.032	0.018	0.032	0.022	0.023	0.095	0.023	0.033
0	0	0	0	0.016	0	0.077	0.077	0	0	0	0.039	0.026	0	0.077	0.007	0.077	0	0	0	0	0.016
0.085	0.171	0	0	0	0	0	0	0	0.171	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0.153	0.051	0.021	0	0	0	0	0	0	0.051	0.069	0.072	0	0.020	0	0.034	0.072	0	0.072	0.063
0	0	0	0.052	0	0.105	0	0	0.052	0	0.052	0	0.035	0	0	0.020	0	0.035	0	0	0	0.022
0.093	0	0	0	0.039	0	0	0	0.093	0	0.093	0	0	0	0	0.018	0	0	0	0	0	0
0.077	0	0.154	0	0.032	0	0	0	0.077	0	0.077	0	0	0	0	0.030	0	0.052	0	0	0	0.032

- Most common words (e.g., "model") have low values
- Peculiar words ("e.g., "medlin", "micro", "methodolog") have high values





Document-Term matrix X

- Local weight I_{ii}
- Global weight g_j
- $\mathbf{X}_{ij} = \mathbf{I}_{ij} * \mathbf{g}_{j}$



Similarity and dissimilarity



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

The following table shows the similarity and dissimilarity between two objects, *x* and *y*, with respect to a single, simple attribute.

Attribute	Dissimilarity	Similarity				
Type						
Nominal	$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$				
Ordinal	d = x - y /(n - 1) (values mapped to integers 0 to $n-1$, where n is the number of values)	s = 1 - d				
Interval or Ratio	d = x - y	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min d}{\max d - \min d}$				





Euclidean Distance

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^{n} (x_k - y_k)^2}$$

where *n* is the number of dimensions (attributes) and x_k and y_k are, respectively, the k^{th} attributes (components) of data objects **x** and **y**.

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	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 is a generalization of Euclidean Distance

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{k=1}^{n} |x_k - y_k|^r\right)^{1/r}$$

Where *r* is a parameter, *n* is the number of dimensions (attributes) and x_k and y_k are, respectively, the k^{th} attributes (components) or data objects *x* and *y*.





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
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
· · · · · · · · · · · · · · · · · · ·	I			
\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	0

Distance Matrix



Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well-known properties.
 - *I*. $d(\mathbf{x}, \mathbf{y}) \ge 0$ for all \mathbf{x} and \mathbf{y} and $d(\mathbf{x}, \mathbf{y}) = 0$ only if $\mathbf{x} = \mathbf{y}$. (Positive definiteness)
 - 2. $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all \mathbf{x} and \mathbf{y} . (Symmetry)
 - *d*(\mathbf{x}, \mathbf{z}) $\leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all points \mathbf{x}, \mathbf{y} , and \mathbf{z} . (Triangle Inequality)

where $d(\mathbf{x}, \mathbf{y})$ is the distance (dissimilarity) between points (data objects) \mathbf{x} and \mathbf{y} .

A distance that satisfies these properties is a metric





- Similarities also have some well known properties
 - 1. $s(\mathbf{x}, \mathbf{y}) = 1$ (or maximum similarity) only if $\mathbf{x} = \mathbf{y}$

$$s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}) \text{ for all } \mathbf{x} \text{ and } \mathbf{y} \text{ (symmetry)}$$

where $s(\mathbf{x}, \mathbf{y})$ is the similarity between points (data objects) \mathbf{x} and \mathbf{y}





Mahalanobis Distance





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



Mahalanobis Distance



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
 = (M₁₁ + M₀₀) / (M₀₁ + M₁₀ + M₁₁ + M₀₀)
 - J = number of 11 matches / number of not-both-zero attributes values = $(M_{11}) / (M_{01} + M_{10} + M_{11})$





SMC versus Jaccard: Example

 $M_{01} = 2$ (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 d₁ and d₂ are two document vectors, then
cos(d₁, d₂) = (d₁ • d₂) / ||d₁|| ||d₂||,
where • indicates vector dot product and || d || is the norm of vector d

• Example:

 $d_1 = 3205000200$ $d_2 = 100000102$

 $d_1 \bullet d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$ $||d_1|| = (3*3+2*2+0*0+5*5+0*0+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5} = 6.481$

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

 $\cos(d_{1'}, d_2) = .3150$



From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006

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, $s_k(\mathbf{x}, \mathbf{y})$, in the range [0, 1].
- 2: Define an indicator variable, δ_k , for the k^{th} attribute as follows:
 - $\delta_k = 0$ if the k^{th} attribute is an asymmetric attribute and both objects have a value of 0, or if one of the objects has a missing value for the kth attribute
 - $\delta_k = 1$ otherwise

3. Compute similarity
$$(\mathbf{x}, \mathbf{y}) = \frac{\sum_{k=1}^{n} \delta_k s_k(\mathbf{x}, \mathbf{y})}{\sum_{k=1}^{n} \delta_k}$$



Using Weights to Combine Similarities

May not want to treat all attributes the same. Use non-negative weights ω_k

• similarity(
$$\mathbf{x}, \mathbf{y}$$
) = $\frac{\sum_{k=1}^{n} \omega_k \delta_k s_k(\mathbf{x}, \mathbf{y})}{\sum_{k=1}^{n} \omega_k \delta_k}$

Can also define a weighted form of distance

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{k=1}^{n} w_k |x_k - y_k|^r\right)^{1/r}$$



From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006

Correlation





Data correlation

- Measure of the linear relationship between two data objects
 - having binary or continuous variables
- Useful during the data exploration phase
 - To be better aware of data properties
- Analysis of feature correlation
 - Correlated features should be removed
 - simplifying the next analytics steps
 - improving the performance of the data-driven algorithms





$$\operatorname{corr}(\mathbf{x}, \mathbf{y}) = \frac{\operatorname{covariance}(\mathbf{x}, \mathbf{y})}{\operatorname{standard_deviation}(\mathbf{x}) * \operatorname{standard_deviation}(\mathbf{y})} = \frac{s_{xy}}{s_x \ s_y},$$

where we are using the following standard statistical notation and definitions

covariance(
$$\mathbf{x}, \mathbf{y}$$
) = $s_{xy} = \frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})(y_k - \overline{y})$

standard_deviation(
$$\mathbf{x}$$
) = $s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \overline{x})^2}$
standard_deviation(\mathbf{y}) = $s_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \overline{y})^2}$

$$\overline{x} = \frac{1}{n} \sum_{k=1}^{n} x_k \text{ is the mean of } \mathbf{x}$$
$$\overline{y} = \frac{1}{n} \sum_{k=1}^{n} y_k \text{ is the mean of } \mathbf{y}$$



Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

Perfect linear correlation when value is 1 or -1



Drawback of Correlation

•
$$\mathbf{x} = (-3, -2, -1, 0, 1, 2, 3)$$

• $\mathbf{y} = (9, 4, 1, 0, 1, 4, 9)$
 $y_i = x_i^2$
• mean(\mathbf{x}) = 0, mean(\mathbf{y}) = 4
• std(\mathbf{x}) = 2.16, std(\mathbf{y}) = 3.74

corr = (-3)(5)+(-2)(0)+(-1)(-3)+(0)(-4)+(1)(-3)+(2)(0)+3(5) / (6 * 2.16 * 3.74))= 0

