Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types
- Measuring Data Similarity and Dissimilarity

Types of Data Sets

- Record
 - Relational records
 - Data matrix, e.g., numerical matrix, crosstabs
 - Document data: text documents: termfrequency vector
 - Transaction data
- Graph and network
 - World Wide Web
 - Social or information networks
 - Molecular Structures
- Ordered
 - Video data: sequence of images
 - Temporal data: time-series
 - Sequential Data: transaction sequences
 - Genetic sequence data
- Spatial, image and multimedia:
 - Spatial data: maps
 - Image data:
 - Video data:

Name	FName	City	Age	Salary
Smith	John	3	35	\$280
Doe	Jane	1	28	\$325
Brown	Scott	3	41	\$265
Howard	Shemp	4	<mark>4</mark> 8	\$359
Taylor	Tom	2	22	\$250

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

Data Objects

- Data sets are made up of data objects.
- A data object represents an entity.
- Examples:
 - sales database: customers, store items, sales
 - medical database: patients, treatments
 - university database: students, professors, courses
- Also called samples, examples, instances, data points, objects, tuples.
- Data objects are described by **attributes**.
- Database rows -> data objects; columns ->attributes.

Attributes

- Attribute (dimension, feature, variable): a data field, representing a characteristic or feature of a data object.
 - E.g., customer _ID, name, address
- Attribute types:
 - Non-numeric: symbolic, non-quantitative
 - Nominal (categorical)
 - Binary
 - Ordinal
 - Numeric: quantitative, a measurable quantity, integer or real
 - Interval-scaled
 - Ratio-scaled

Non-numeric Attribute Types

- Nominal (categorical): categories, states, or "names of things"
 - Hair_color = {auburn, black, blond, brown, grey, red, white}
 - marital status, occupation, ID numbers, zip codes
 - Binary
 - Nominal attribute with only 2 states
 - <u>Symmetric binary</u>: both outcomes equally important
 - e.g., gender
 - <u>Asymmetric binary</u>: outcomes not equally important.
 - e.g., medical test (positive vs. negative)
 - Convention: assign 1 to most important outcome (e.g., HIV positive)

Ordinal

- Values have a meaningful order (ranking) but difference between successive values is unknown.
- Size = {small, medium, large}, letter grades, army rankings

Numeric Attribute Types

Interval-scaled: difference between two values is meaningful

- Measured on a scale of equal-sized units
 - E.g., temperature in C°or F°. pH, calendar dates
 - No true zero
 - neither 0C° nor 0F° indicates no heat
 - without zero, we cannot talk of one temperature value as being a multiple of another. We cannot say 10C° is twice as warm as 5C°
- Ratio-scaled: has all the properties of an interval variable, and also has a clear definition of zero (that means none)
 - e.g., temperature in Kelvin (0 kelvin does mean no heat), length, weight, counts, monetary quantities
 - We can speak of a value as being a multiple (or ratio) of another
 - 10K° is twice as high as 5K°

Discrete vs. Continuous Attributes

Another way to categorize data types

Discrete

- Has only a countable (finite or countably infinite) set of values
 - E.g., zip codes, the set of words in a collection of documents
- Sometimes, represented as integer variables
- Note: Binary attributes are a special case of discrete attributes

Continuous

- Has real numbers (which are uncountable) as attribute values
 - E.g., 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

Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types
- Measuring Data Similarity and Dissimilarity

Similarity and Dissimilarity

Similarity

- Numerical measure of how alike two data objects are
- Value is higher when objects are more alike
- Often falls in the range [0,1]
- Dissimilarity (e.g., distance)
 - Numerical measure of how different two data objects are
 - Lower when objects are more alike
 - Minimum dissimilarity is often 0
 - Upper limit varies

Data Matrix and Dissimilarity Matrix

- Data matrix
 - n x p, object-by-variable
 - n data points with p dimensions
 - Two modes stores both objects and attributes
- Dissimilarity matrix
 - n x n, object-by-object
 - n data points, but registers only the distance
 - A triangular matrix
 - Single mode as it only stores dissimilarity values

$$\begin{bmatrix} x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{n1} & \cdots & x_{nf} & \cdots & x_{np} \end{bmatrix}$$

$$\begin{bmatrix} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ \vdots & \vdots & \vdots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

Nominal Attributes

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of binary attributes)
- Method 1: Simple matching

m: # of matches, p: total # of variables

$$d(i,j) = \frac{p-m}{p}$$

- Method 2: Use a large number of binary attributes
 - creating a new binary attribute for each nominal state

Binary Attributes

				Objec	t <i>j</i>	
•	Contingency table for binary data	Object <i>i</i>	1 0 sum	$\begin{array}{c}1\\q\\s\\q+s\end{array}$	$0 \\ r \\ t \\ r+t$	$sum \\ q+r \\ s+t \\ p$
•	Distance for symmetric binary variables:Similarity?	d(i, .)	$j) = -\frac{1}{q}$	r + r + r + r	$\frac{s}{s+t}$	
•	Distance for asymmetric binary variables:	d(i,	j) =	$\frac{r+}{q+r}$	$\frac{s}{+s}$	
	 Jaccard coefficient (<i>similarity</i> measure for <i>asymmetric</i> binary variables): 	sim_{Jac}	ccard(i	, j) =	$\frac{q}{q+r}$	+s

Binary Attributes: example

Name	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	Y	Ν	Р	Ν	Ν	Ν
Mary	Y	Ν	Р	Ν	Р	Ν
Jim	Y	Р	Ν	Ν	Ν	Ν

- All attributes are asymmetric binary
- Let the values Y and P be 1, and the value N be 0

$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$
$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$
$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

Numeric Attributes: Standardization

- Z-score: $z = \frac{x \mu}{\sigma}$
 - x: raw score to be standardized, μ : mean of the population, σ : standard deviation
 - the distance between the raw score and the population mean in units of the standard deviation
 - negative when the raw score is below the mean, positive when above
- An alternative way: Calculate the mean absolute deviation

$$s_{f} = \frac{1}{n} (|x_{1f} - m_{f}| + |x_{2f} - m_{f}| + \dots + |x_{nf} - m_{f}|)$$

- where $m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$ $z_{if} = \frac{x_{if} - m_f}{S_f}$
- standardized measure (*z*-score):
- Using mean absolute deviation is more robust than using standard deviation

Numeric Attributes: Minkowski Distance

Minkowski distance: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h} + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h$$

where $i = (x_{i1}, x_{i2}, ..., x_{ip})$ and $j = (x_{j1}, x_{j2}, ..., x_{jp})$ are two *p*dimensional data objects, and *h* is the order (the distance so defined is also called L-*h* norm)

- Properties
 - d(i, j) > 0 if i ≠ j, and d(i, i) = 0 (Positive definiteness)
 - d(i, j) = d(j, i) (Symmetry)
 - $d(i, j) \le d(i, k) + d(k, j)$ (Triangle Inequality)
- A distance that satisfies these properties is a metric

Special Cases of Minkowski Distance

- h = 1: Manhattan (city block, L₁ norm) distance
 - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i,j) = |x_{i_1} - x_{j_1}| + |x_{i_2} - x_{j_2}| + \dots + |x_{i_p} - x_{j_p}|$$

• h = 2: (L₂ norm) Euclidean distance

$$d(i,j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- $h \rightarrow \infty$. "supremum" (L_{max} norm, L_{∞} norm) distance.
 - This is the maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{h \to \infty} \left(\sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_{f}^{p} |x_{if} - x_{jf}|^h$$

Minkowski Distance: Example

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5

Manhattan (L₁)

L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

Euclidean (L₂)

L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

Supremum (L_{∞})

L_{∞}	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0



Ordinal Variables

Can be treated like interval-scaled

- replace x_{if} by their rank $r_{if} \in \{1, \dots, M_f\}$
- map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by

$$z_{if} = rac{r_{if} - 1}{M_f - 1}$$

- Turned to values with equal intervals
- compute the dissimilarity using methods for interval-scaled variables, e.g., Euclidean distance
- Example: a, b, c, d -> 1, 2, 3, 4 -> 0, 1/3, 2/3, 3/3

Ordinal Variables: Example

- Consider the data in the table:
- Here, the attribute Test has three states: fair, good and excellent, so M_f=3
- For step 1, the 3 attribute values are assigned the ranks 1, 2 and 3 respectively.
- Step 2 normalizes the ranking by mapping rank 1 to 0, rank 2 to 0.5 and rank 3 to 1
- For step 3, using Euclidean distance, a dissimilarity matrix is obtained as shown
- Therefore, students 1 and 2 are most dissimilar, as are students 2 and 4

Student	Test
1	Excellent
2	Fair
3	Good
4	Excellent



Attributes of Mixed Types

- A database may contain multiple attribute types
- use a weighted formula to combine their effects

$$d(i,j) = \frac{\sum_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^{p} \delta_{ij}^{(f)}}$$

• *f* is binary or nominal:

 $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, or $d_{ij}^{(f)} = 1$ otherwise

- *f* is numeric: use the normalized distance
- $d_{ij}^{(f)} = \frac{|x_{if} x_{jf}|}{\max_{ij} x_{hf} \min_{i} x_{hf}}$

- f is ordinal
 - Compute ranks r_{if} and $Z_{if} = \frac{r_{if} 1}{M_{f} 1}$
 - Treat z_{if} as numeric
- The indicator delta is generally set to 1, but
- If f is asymmetric binary and x_{if} = x_{if} = 0, set the indicator to 0
 - recall we removed t from consideration for "Distance for asymmetric binary variables"



 A document can be represented by thousands of attributes, each recording the frequency of a particular word (such as keywords) or phrase in the document.

Document	team	coach	hockey	base ball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

- Other vector objects: gene features in micro-arrays, ...
- Applications: information retrieval, biologic taxonomy, gene feature mapping, ...
- Cosine measure: If d₁ and d₂ are two vectors (e.g., term-frequency vectors), then cos(d₁, d₂) = (d₁ d₂) / ||d₁|| ||d₂|| ,

where \bullet indicates vector dot product, ||d||: the length of vector d

Example: Cosine Similarity

- $\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1|| ||d_2||$, where • indicates vector dot product, ||d|: the length of vector d
- Ex: Find the **similarity** between documents 1 and 2.

 $d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$ $d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$

 $\begin{aligned} &d_1 \bullet d_2 = 5^* 3 + 0^* 0 + 3^* 2 + 0^* 0 + 2^* 1 + 0^* 1 + 0^* 1 + 2^* 1 + 0^* 0 + 0^* 1 = 25 \\ &||d_1|| = (5^* 5 + 0^* 0 + 3^* 3 + 0^* 0 + 2^* 2 + 0^* 0 + 0^* 0 + 2^* 2 + 0^* 0 + 0^* 0)^{0.5} = (42)^{0.5} = 6.481 \\ &||d_2|| = (3^* 3 + 0^* 0 + 2^* 2 + 0^* 0 + 1^* 1 + 1^* 0 + 0 + 1^* 1 + 0^* 0 + 1^* 1)^{0.5} = (17)^{0.5} = 4.12 \\ &\cos(d_1, d_2) = 0.94 \end{aligned}$