| Data Mining: Data | _ |
|--|---|
| Dr. Hui Xiong Rutgers University | |
| THE STATE UNIVERSITY OF NEW JERSEY | |
| Introduction to Data Mining 1/2/2009 1 | |

| Outline | | |
|--|----------|---|
| Attributes and Objects | | |
| Types of Data | | |
| Data Quality | | |
| Data Preprocessing | | |
| | | |
| Introduction to Data Mining | 1/2/2009 | 2 |













| | Attribute Type | Description | Examples | Operations |
|--------------------|-------------------|--|---|--|
| guncar litative | Nominal | Nominal attribute values only distinguish. (=, ≠) | zip codes, employee ID numbers, eye color, sex: { <i>male,</i> <i>female</i> } | mode, entropy, contingency correlation, χ2 test |
| Qua | Ordinal | Ordinal attribute values also order objects. (<, >) | hardness of minerals, {good, better, best}, grades, street numbers | median, percentiles, rank correlation, run tests, sign tests |
| ntitative | Interval | For interval attributes, differences between values are meaningful. (+, -) | calendar dates, temperature in Celsius or Fahrenheit | mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests |
| Quar | Ratio | For ratio variables, both differences and ratios are meaningful. (*, /) | temperature in Kelvin, monetary quantities, counts, age, mass, length, current | geometric mean harmonic mean, percent variation |

| | Attribute Type | Transformation | Comments |
|-------------|-------------------|---|---|
| ve | Nominal | Any permutation of values | If all employee ID numbers were reassigned, would it make any difference? |
| Qualitati | Ordinal | An order preserving change of values, i.e., new_value = f(old_value) 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, 0}. |
| Jantitative | Interval | <i>new_value =a</i> * <i>old_value</i> + <i>b</i> 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). |
| วั | Ratio | new_value = a * old_value | Length can be measured in meters or feet. |

This categorization of attributes is due to S. S. Stevens











| D | Data Matr | ix | | | |
|---|---|--|--|---|---|
| • | If data object attributes, the points in a m dimension re Such data se where there a | s have the sa en the data ol ulti-dimensior presents a di et can be repr | ame fixed se bjects can b nal space, v stinct attribu esented by | et of num be thoug vhere ea ute an m by | neric ht of as ach y n matrix, |
| | columns, one | e for each attr | one for each ibute | object, | and n |
| | columns, one Projection of x Load | e for each attr Projection of y load | ibute | Load | and n Thickness |
| | Columns, one Projection of x Load | Projection of y load | Distance | Load | and n Thickness 1.2 |
| | Columns, one Projection of x Load 10.23 12.65 | Projection of y load 5.27 6.25 | Distance | Load 2.7 2.2 | and n Thickness 1.2 1.1 |





















































































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

| Attribute | Dissimilarity | Similarity | | |
|---|---|--|--|--|
| Type | | | | |
| Nominal | $d = \left\{egin{array}{ll} 0 & 	ext{if } p = q \ 1 & 	ext{if } p eq q \end{array} ight.$ | $s = \left\{egin{array}{cc} 1 & 	ext{if } p = q \ 0 & 	ext{if } p eq q \end{array} ight.$ | | |
| 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 | | | | |
| Introduction | to Data Mining 1/2/2009 | 59 | | |









| | | | L1 | p1 | p2 | р3 | p4 |
|-------|------------|----------------|----------------|-----------|----------|-------|-------|
| | | | p1 | 0 | 4 | 4 | 6 |
| | | | p2 | 4 | 0 | 2 | 4 |
| | | | p3 | 4 | 2 | 0 | 2 |
| | | , | p4 | 6 | 4 | 2 | 0 |
| point | X | y | 10 | | | | |
| p1 | 0 | 2 | L2 | p1 | p2 | p3 | p4 |
| p2 | 2 | 0 | p1 | 0 | 2.828 | 3.162 | 5.099 |
| p3 | 3 | 1 | p2 | 2.828 | 0 | 1.414 | 3.162 |
| p4 | 5 | 1 | p3 | 3.162 | 1.414 | 0 | 2 |
| | | | p4 | 5.099 | 3.162 | 2 | (|
| | | | L _∞ | p1 | p2 | р3 | p4 |
| | | | p1 | 0 | 2 | 3 | 5 |
| | | | p2 | 2 | 0 | 1 | 3 |
| | | | р3 | 3 | 1 | 0 | 2 |
| | | | p4 | 5 | 3 | 2 | 0 |
| | | | | Distanc | e Matrix | | |
| Int | troduction | to Data Mining | | 1/2/200 | ٥ | | 64 |































