□ 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 {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/operations it possesses: *Nominal*

- Distinctness:  $= \neq$  *ordinal*
- Order:  $< >$
- Differences are meaningful  $+ -$  *Interval*
- Ratios are meaningful  $* /$  *Ratio Attribute*

- Nominal attribute: distinctness
- Ordinal attribute: distinctness & order
- Interval attribute: distinctness, order & meaningful differences
- Ratio attribute: all 4 properties/operations

| | Attribute Type | Description | Examples | Operations |
|---|---|---|---|---|
| Categorical Qualitative | Nominal | Nominal attribute values only distinguish. (=, ≠) | zip codes, employee ID numbers, eye color, sex: { *male, female* } | mode, entropy, contingency correlation, χ2 test |
| | Ordin al | Ordinal attribute values also order objects. (<, >) | hardness of minerals, {*good, better, best* }, grades, street numbers | median, percentiles, rank correlation, run tests, sign tests |
| Numeric Quant | Interval | For interval attributes, differences bet ween values are meaningful. (+, - ) | calendar dates, temperature in Celsius or Fahrenheit | mean, standard deviation, Pearson's correlation, *t* and *F* tests |
| | 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 |

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

# 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.

# Asymmetric Attributes

- Only presence (a non-zero attribute value) is regarded as important
  - Words present in documents
  - Items present in customer transactions

- If we met a friend in the grocery store would we ever say the following? *We don't care about the attributes not in consideration*
  *"I see our purchases are very similar since we didn't buy most of the same things."* eg: we don't care about stuff not present, we don't even regard them as attributes

- We need two asymmetric binary attributes to represent one ordinary binary attribute
  - Association analysis uses asymmetric attributes

- Asymmetric attributes typically arise from objects that are sets

---

Symmetric : Yes or No  (we also need to know if
                          *present?*           not present)

Asymmetric attributes : Only if they are present.
                          absence is not noted

## Record Data

# 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

# 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 |

# 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 |

# 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.

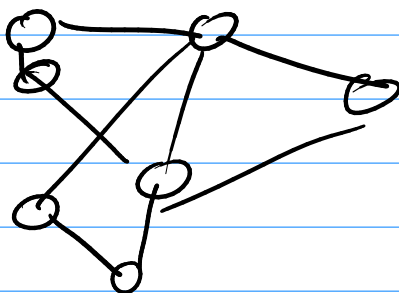| | team | coach | play | ball | score | game | win | lost | timeout | season |
|-----------|------|-------|------|------|-------|------|-----|------|---------|--------|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1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

□ 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 |

graphs → molecular structures generic graphs etc.

What kind of Common structures do these compounds have that make them have same properties

Ordered data
↳ genomic sequence data
Spatio-temporal data

## 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
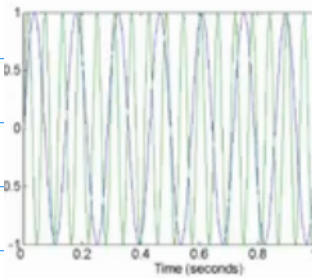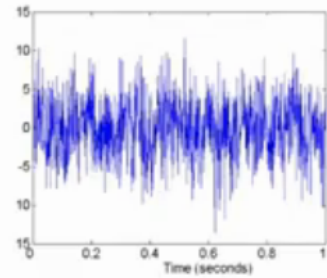  - More loans are given to individuals that default

## Noise

- For objects, noise is an extraneous object
- For attributes, 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**

## 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 ✓
  - Wrong data ✓