# BITS Pilani

**BITS** Pilani
Hyderabad Campus

Prof.Aruna Malapati
Department of CSIS

**Data**

# Today's Learning objective

- **Describe Data**

- **List various Data types**

- **List the issues in Data quality**

- **List and identify the right preprocessing techniques given 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.
  - Other names: variable, filed, characteristic, feature, Predictor, etc.

- A collection of attributes describe an object
  - Other names: record, point, case, sample, entity, or instance

Attributes

Objects

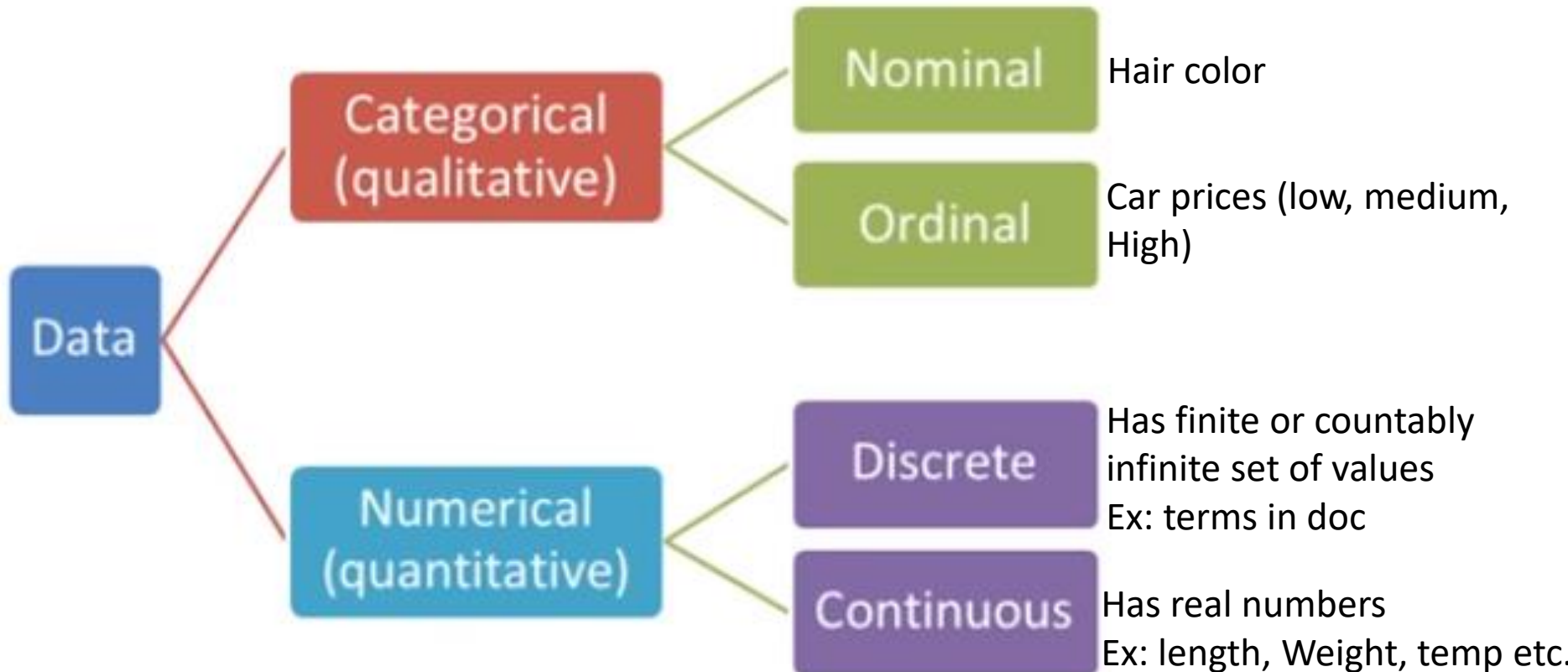| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | 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   |

# Attribute Values

- Each attribute has a set of values object draws from.

- The same attribute can be mapped to different attribute values

  – Example: Temperature can be Celsius in feet or Fahrenheit

- Different attributes can be mapped to the same set of values

  – Example: Attribute values for ID and age are integers

# Types of Attributes

**Data**

**Categorical (qualitative)**
- Nominal — Hair color
- Ordinal — Car prices (low, medium, High)

**Numerical (quantitative)**
- Discrete — Has finite or countably infinite set of values. Ex: terms in doc
- Continuous — Has real numbers. Ex: length, Weight, temp etc.

# 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. (=, ≠) | Zip codes, employee ID numbers, eye color, sex: {*male, female*} | mode, entropy, contingency correlation, $\chi^2$ test |
| Ordinal | The values of an ordinal attribute provide enough information to order objects. (<, >) | hardness of minerals, {*good, better, best*}, 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, *t* and *F* 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 |

# Discrete and Continuous Attributes

- Discrete Attribute
  - Has only a finite or countably infinite set of values
  - Examples: 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 |

- 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 |

# 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



```
<a href="papers/papers.html#bbbb">
Data Mining </a>
<li>
<a href="papers/papers.html#aaaa">
Graph Partitioning </a>
<li>
<a href="papers/papers.html#aaaa">
Parallel Solution of Sparse Linear System of Equations </a>
<li>
<a href="papers/papers.html#ffff">
N-Body Computation and Dense Linear System Solvers
```

# Chemical Data

Benzene Molecule: $C_6H_6$

Sequences of transactions

( A B)    (D)   (C E)
( B D)    (C)   (E)
( C D)    (B)  (A E)

An element of the
sequence

Genomic sequence data

GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

# Ordered Data

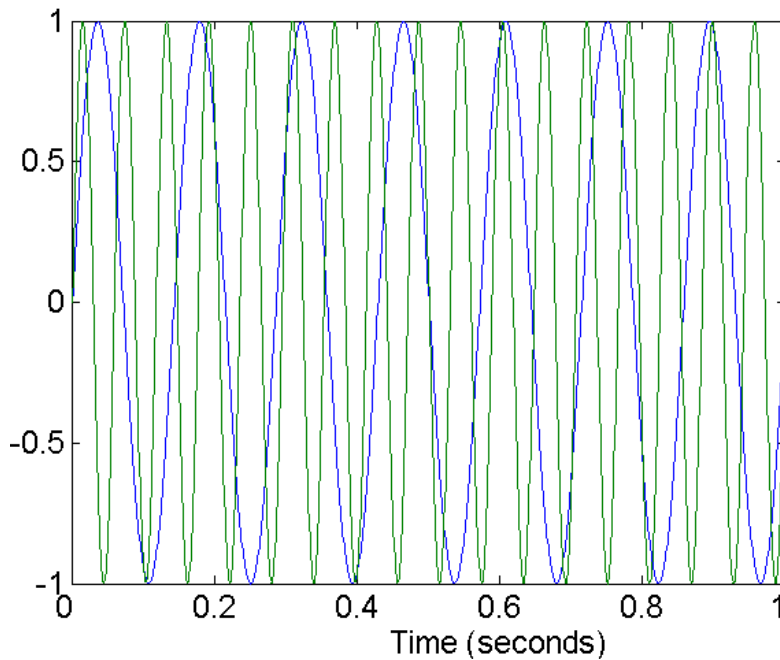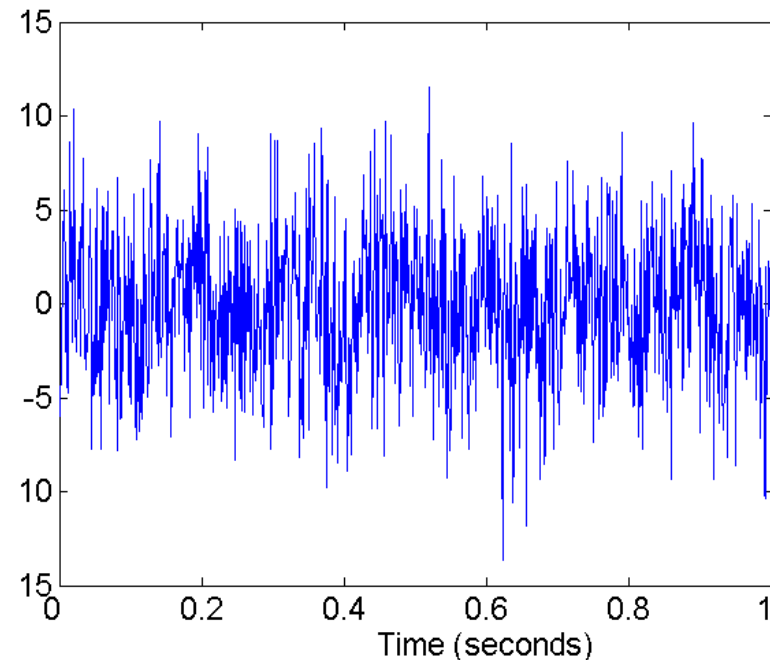## Spatio-Temporal Data

Jan



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: An invalid signal overlapping valid data
  - 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

# 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 heterogeneous 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

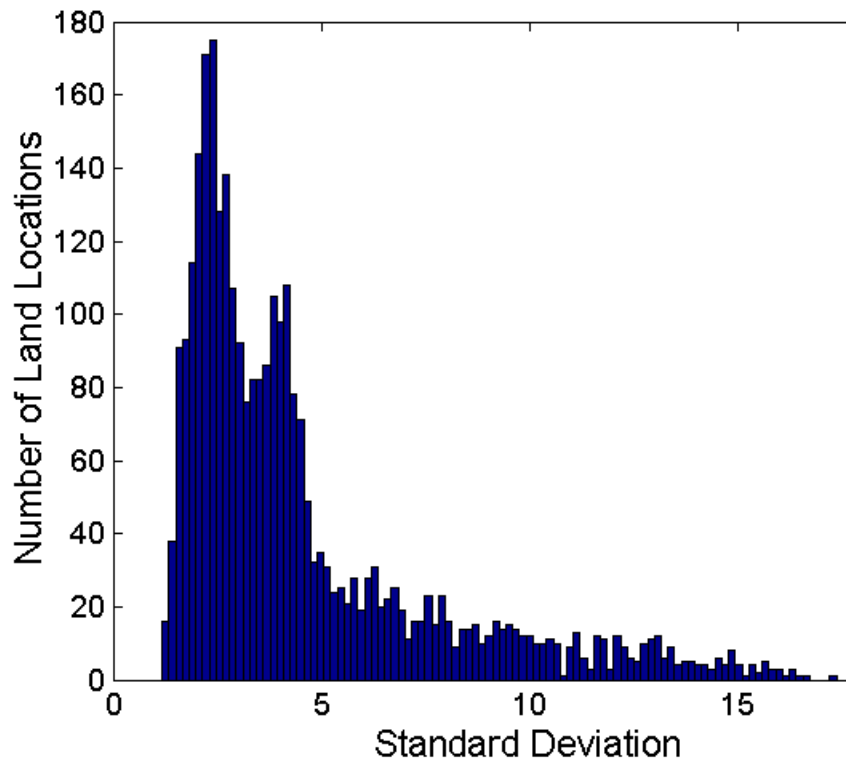- Attribute Transformation

# 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
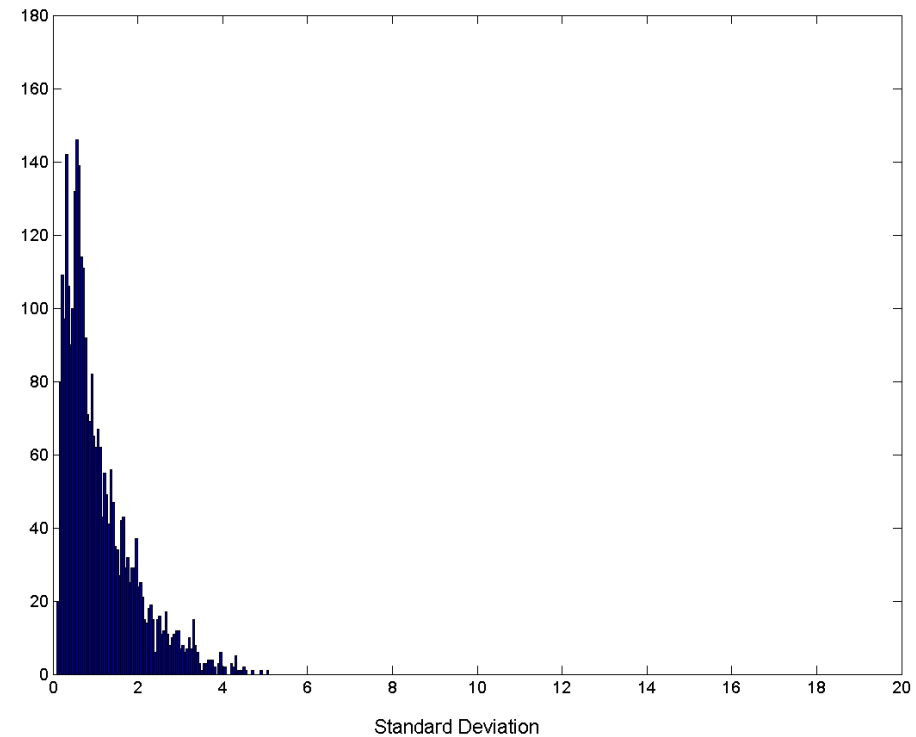    - Aggregated data tends to have less variability

# Aggregation

Variation of Precipitation in Australia



Standard Deviation of Average
Monthly Precipitation

Standard Deviation of
Average Yearly Precipitation

# 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 <span style="color:red">obtaining</span> the entire set of data of interest is too expensive or time consuming.

- Sampling is used in data mining because <span style="color:red">processing</span> the entire set of data of interest is too expensive or time consuming.
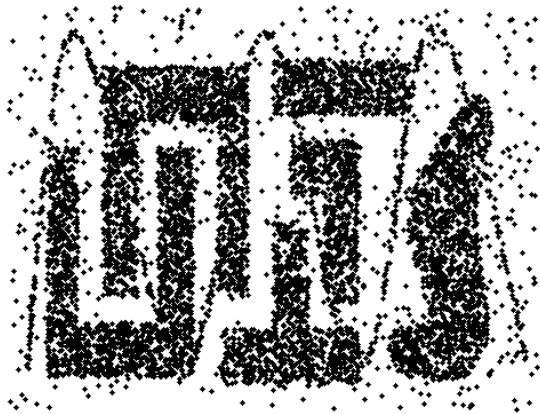
# Sampling

- The key principle for effective sampling is:

- A sample will work almost as well as using the entire data set <span style="color:red">if the sample is representative(different for different data set)</span>.

- Sampling may <span style="color:red">remove outliers</span> and if done improperly can introduce noise.
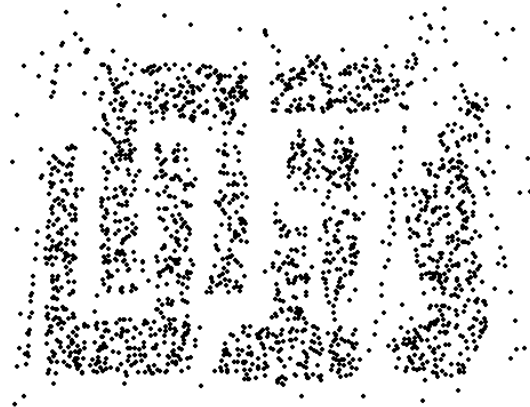
# Types of Sampling

- ## 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
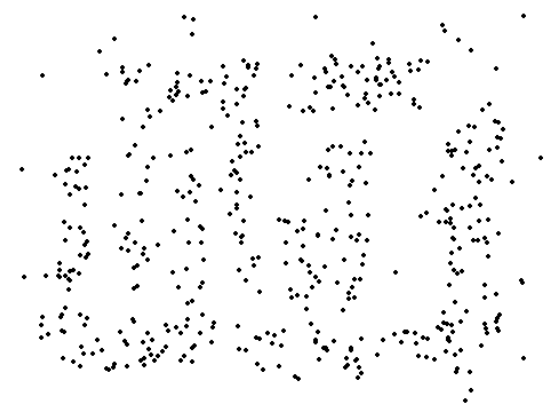
# Sample Size



8000 points                2000 Points                500 Points

# Take home message

- Four different features/attributes/measurements/ independent variables can be of type Nominal, Ordinal, Interval or Ratio type.

- Based on the type of data, the operations vary.

- The data set can be of the record, graph, or ordered type.

- Real-world data is dirty, so preprocessing is a very important step in Data Mining.

- There are several methods for preprocessing, choosing the right method depends on the problem and data obtained.