

# INTRODUCTION TO DATA MINING: DATA PREPROCESSING

Chiara Renso  
KDD-LAB  
ISTI- CNR, Pisa, Italy  
[chiara.renso@isti.cnr.it](mailto:chiara.renso@isti.cnr.it)



1

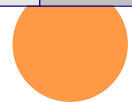
# WHAT IS DATA?

- 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, or feature
- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

## Attributes

<i>Tid</i>	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

## Objects



# TYPES OF ATTRIBUTES

- 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, temperatures in Celsius or Fahrenheit.
  - **Ratio**
    - ◆ Examples: temperature in Kelvin, length, time, counts

# 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.

# 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

<i>Tid</i>	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



# 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

# 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.

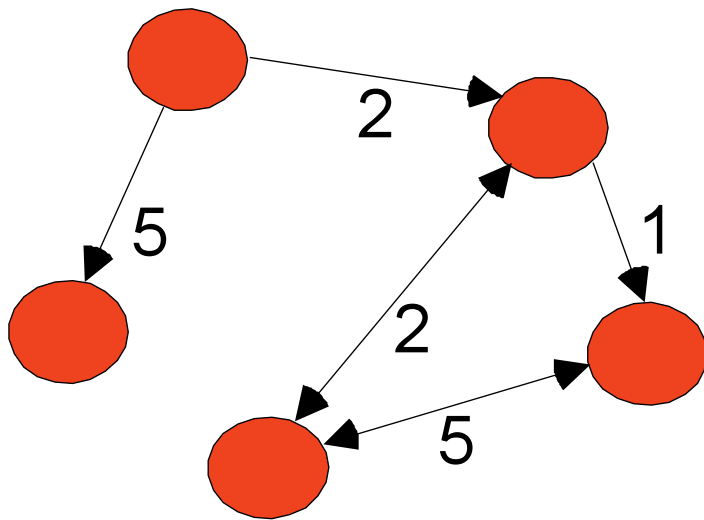
<i>TI</i> <i>D</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

item

transaction

# GRAPH DATA

- Examples: Generic graph and HTML Links



`<a href="papers/papers.html#bbbb"> Data Mining </a>`

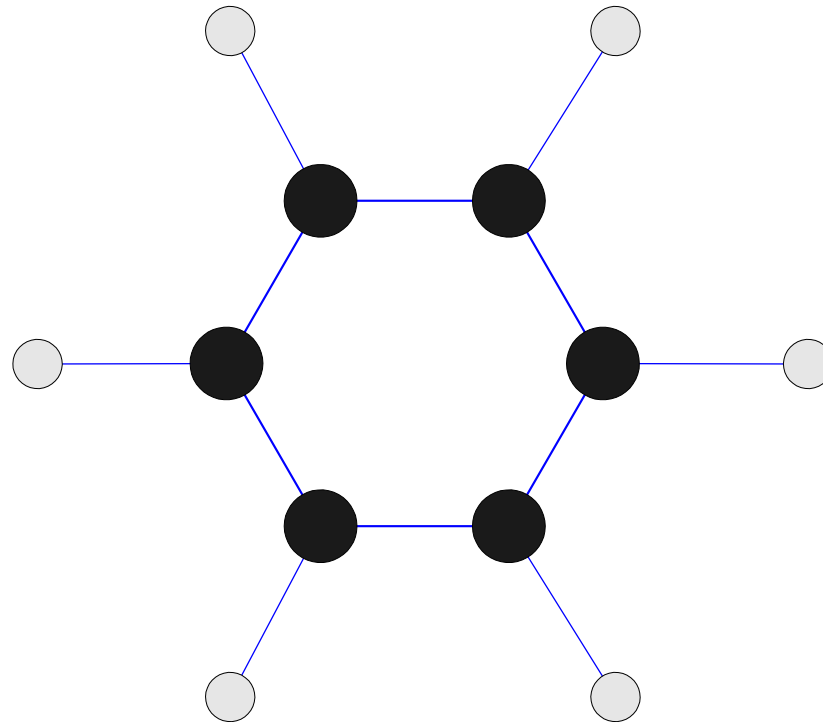
`<li> <a href="papers/papers.html#aaaa"> Graph Partitioning </a>`

`<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

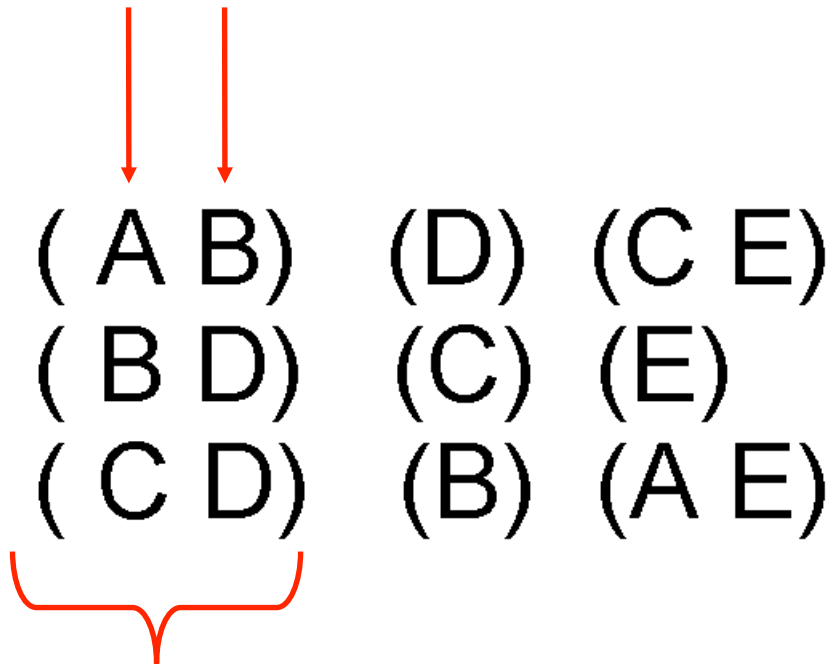
- Benzene Molecule:  $C_6H_6$



# ORDERED DATA

- Sequences of transactions

Items/Events



An element of  
the sequence

# ORDERED DATA

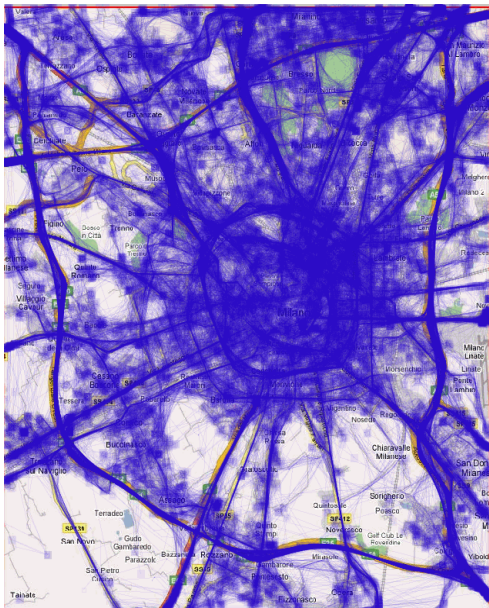
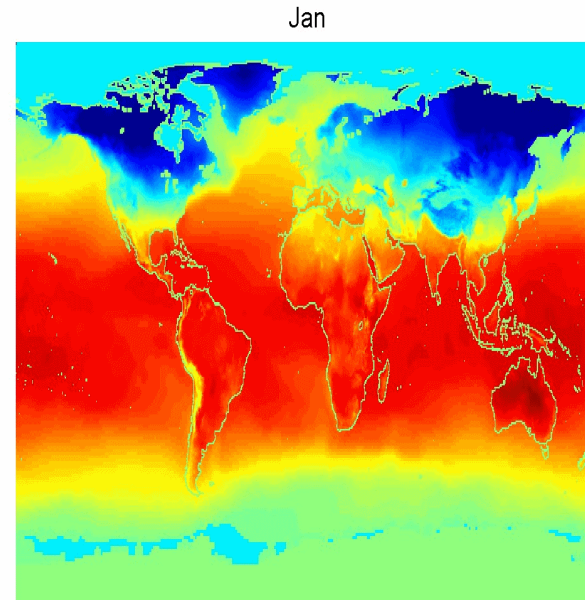
- Genomic sequence data

```
GGTTC CGCCTTCAGCCCCGCGCC  
CGCAGGGCCCGCCCCGCGCCGTC  
GAGAAGGGCCCGCCTGGCGGGCG  
GGGGGAGGCGGGGCCGCCCGAGC  
CCAACCGAGTCCGACCAGGTGCC  
CCCTCTGCTCGGCCTAGACCTGA  
GCTCATTAGGCGGCAGCGGACAG  
GCCAAGTAGAACACGCGAAGCGC  
TGGGCTGCCTGCTGCGACCAGGG
```

# ORDERED DATA

- Spatio-Temporal Data

**Average Monthly  
Temperature of  
land and ocean**



**Trajectories of  
Moving Objects**

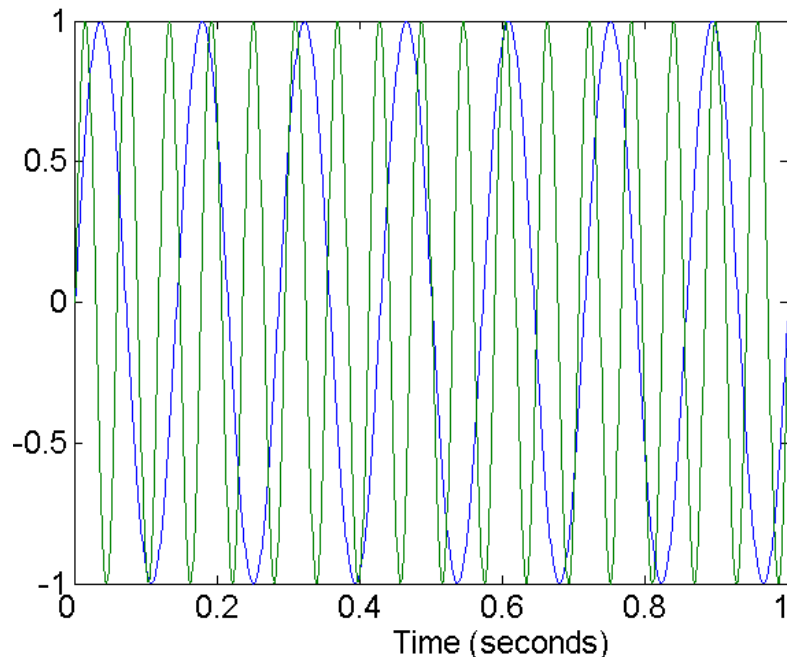
# DATA QUALITY

- 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

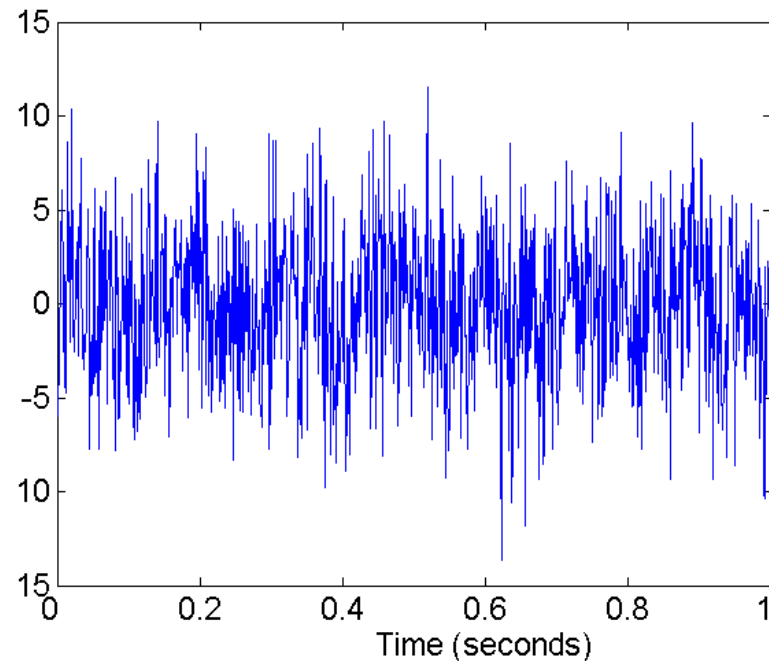


# NOISE

- Noise refers to modification of original values
  - Examples: distortion of a person's voice when talking on a poor phone and “snow” on television screen



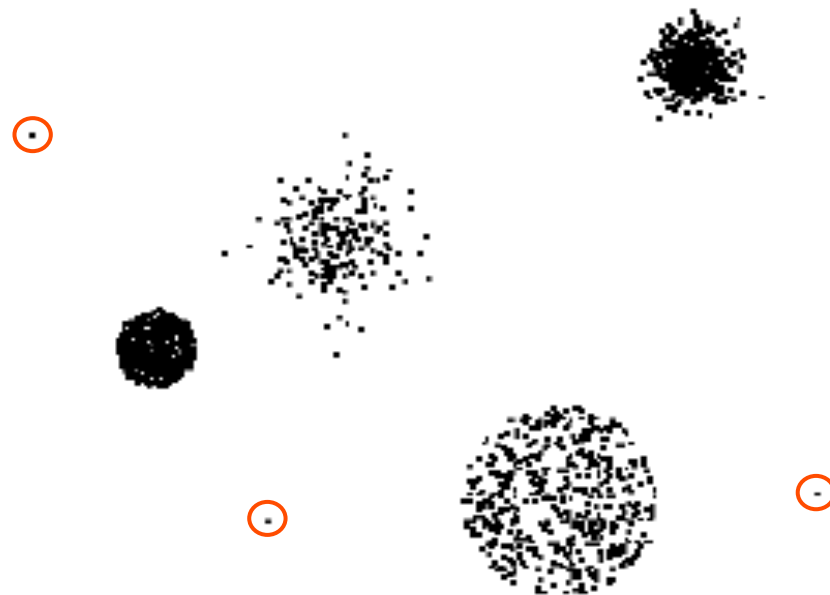
**Two Sine Waves**



**Two Sine Waves + Noise**

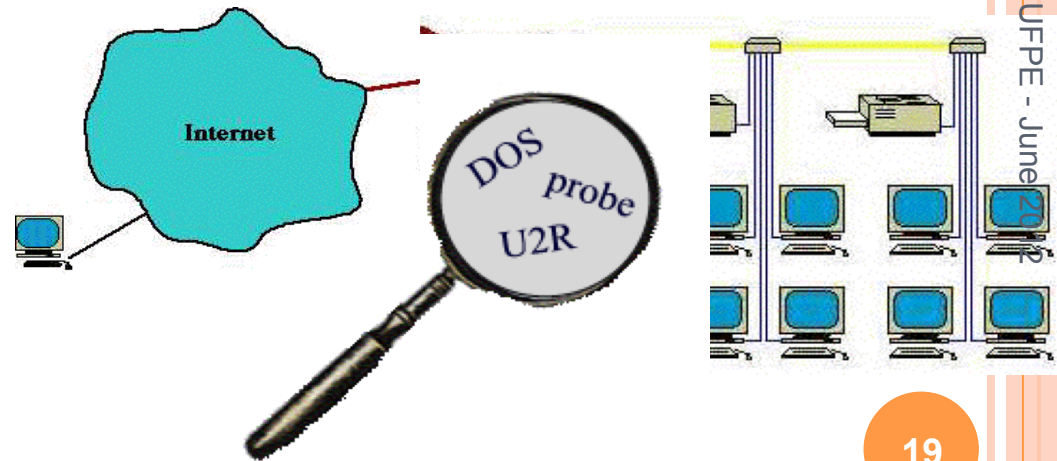
# OUTLIERS

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



# DEVIATION/ANOMALY DETECTION

- Outliers are useful when we need to detect significant deviations from normal behavior
- Applications:
  - Credit Card Fraud Detection
  - Network Intrusion Detection



# MISSING VALUES

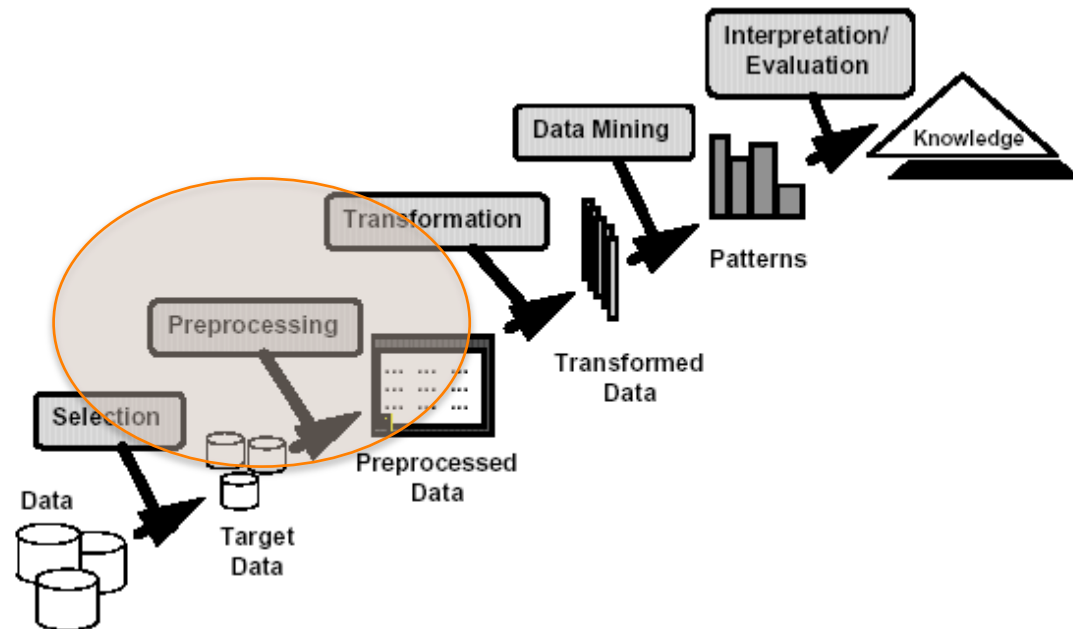
- Reasons for missing values
  - Information is not collected (e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases (e.g., annual income is not applicable to children)
- Handling missing values
  - Eliminate Data Objects
  - Estimate Missing Values
  - 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

# SAMPLING

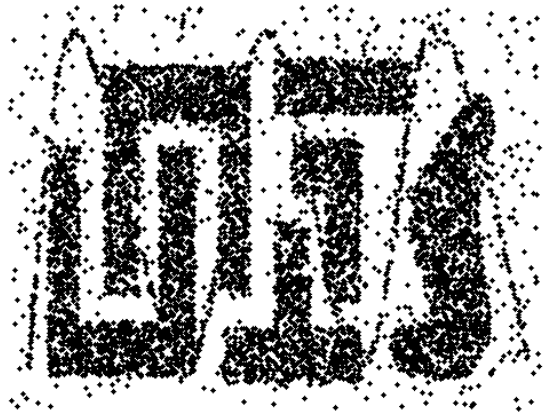
- Sampling is the main technique employed for data selection.
  - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming.
- Sampling is used in data mining because **processing the entire set of data of interest is too expensive** or time consuming.



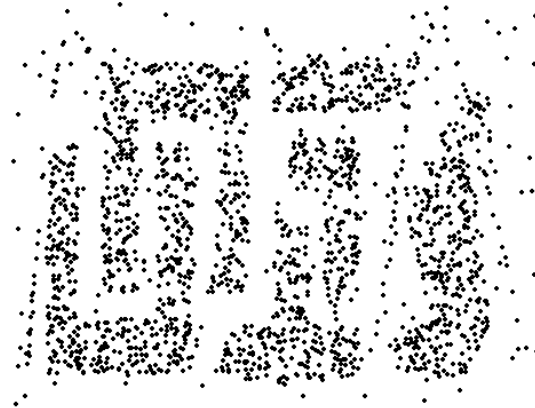
## SAMPLING ...

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

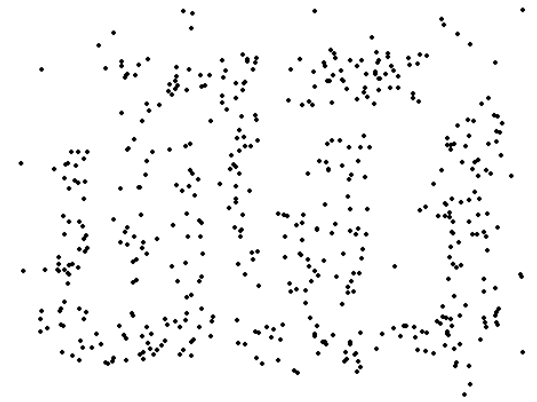
# SAMPLE SIZE



**8000 points**



**2000 Points**



**500 Points**

# 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**.



# 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

# FEATURE SUBSET SELECTION

- Reduce dimensionality of data

Remove:

- **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 often 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