

Data Mining : Data

Marco Loog



Outline

- Data, objects, attributes
- Data properties and types
- Data quality
- Preprocessing of data
- The curse of dimensionality

“What is Data?”

- Collection of 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, 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

Attribute Values

- Attribute values are numbers or symbols assigned to an attribute
- Distinction between attributes and attribute values
 - Same attribute can be mapped to different attribute values
 - Example: height can be measured in feet or meters
 - Different attributes can be mapped to the same set of values
 - Example: Attribute values for ID and age are integers

Properties of Attribute Values

- Mathematical properties / operations:
 - Distinctness $= \neq$
 - Order $< >$
 - Addition $+ -$
 - Multiplication $* /$
- The type of an attribute depends on which of these apply
 - Nominal attribute : distinctness
 - Ordinal attribute : distinctness & order
 - Interval attribute : distinctness, order, & addition
 - Ratio attribute : all 4 properties

Examples of the Four Types

- 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
- Which properties are allowed is largely about making sense
 - Formally any operation could be applied to any data type
 - Depends on context, objective, etc.

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
 - Continuous attributes are typically represented as floating-point variables
- Practically, any value can only be measured and represented in a finite way

Examples

- For each of the following attributes, consider whether it is binary, discrete, or continuous and whether it's nominal, ordinal, interval, or ratio.
 - Age in years
 - Time in terms of AM or PM
 - Brightness as measured by a light meter
 - Brightness as measured by people's judgments
 - Bronze, silver, and gold medals as awarded at the Olympics
 - Height above sea level
 - Number of patients in a hospital
 - ISBN numbers for books
 - Military rank
 - Distance from the center of campus
 - Temperature in degrees Kelvin
 - Temperature in degrees Celsius
 - Coat check number

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



Record Data

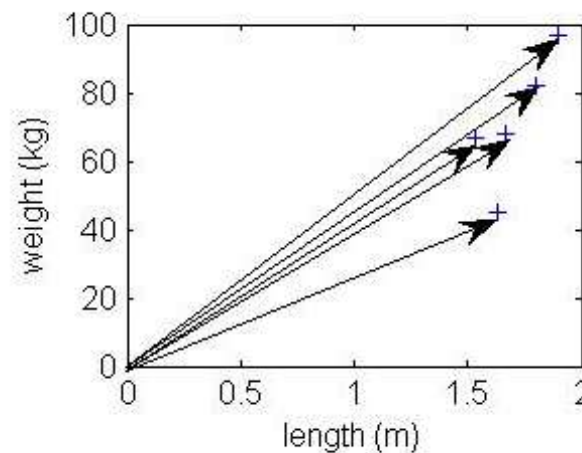
- Data that consists of a collection of records, each of which consists of a fixed set of attributes

<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 and Vector Space

- 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

<i>Tid</i>	Length	Weight
1	1.80	82
2	1.53	67
3	1.67	68
4	1.90	97
5	1.63	45



Document Data

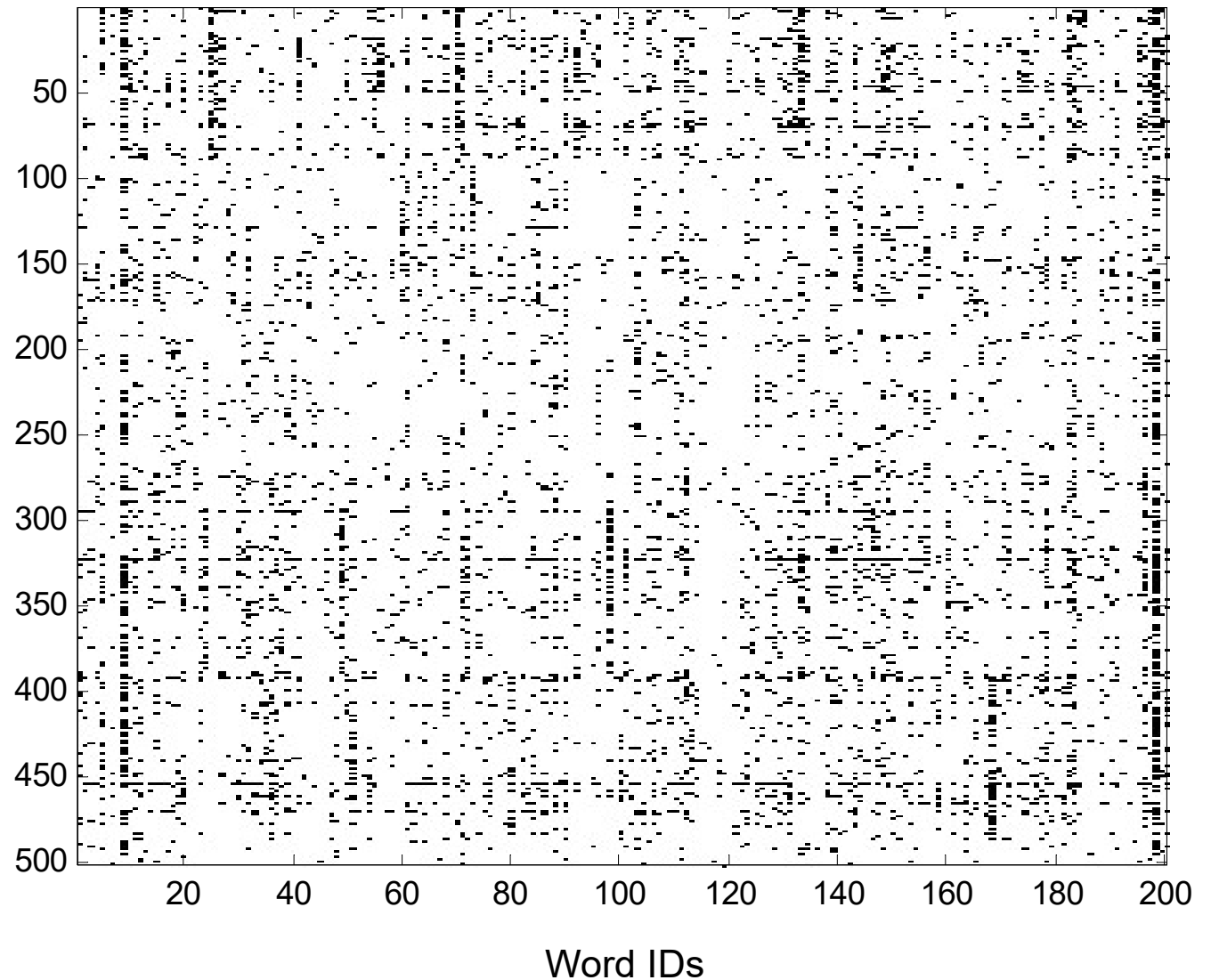
- Each document becomes a “term” vector
 - Each term is a component [attribute] of the vector
 - The value of each component is
the number of times the corresponding term occurs in the document

Document	team	coach	play	ball	score	game	win	lost	timeout	season
1	3	0	5	0	2	6	0	2	0	2
2	0	7	0	2	1	0	0	3	0	0
3	0	1	0	0	1	2	2	0	3	0
4	1	4	0	2	3	0	1	6	2	1
5	2	3	3	1	6	1	3	0	0	4

Sparse Document Matrix



Text
documents



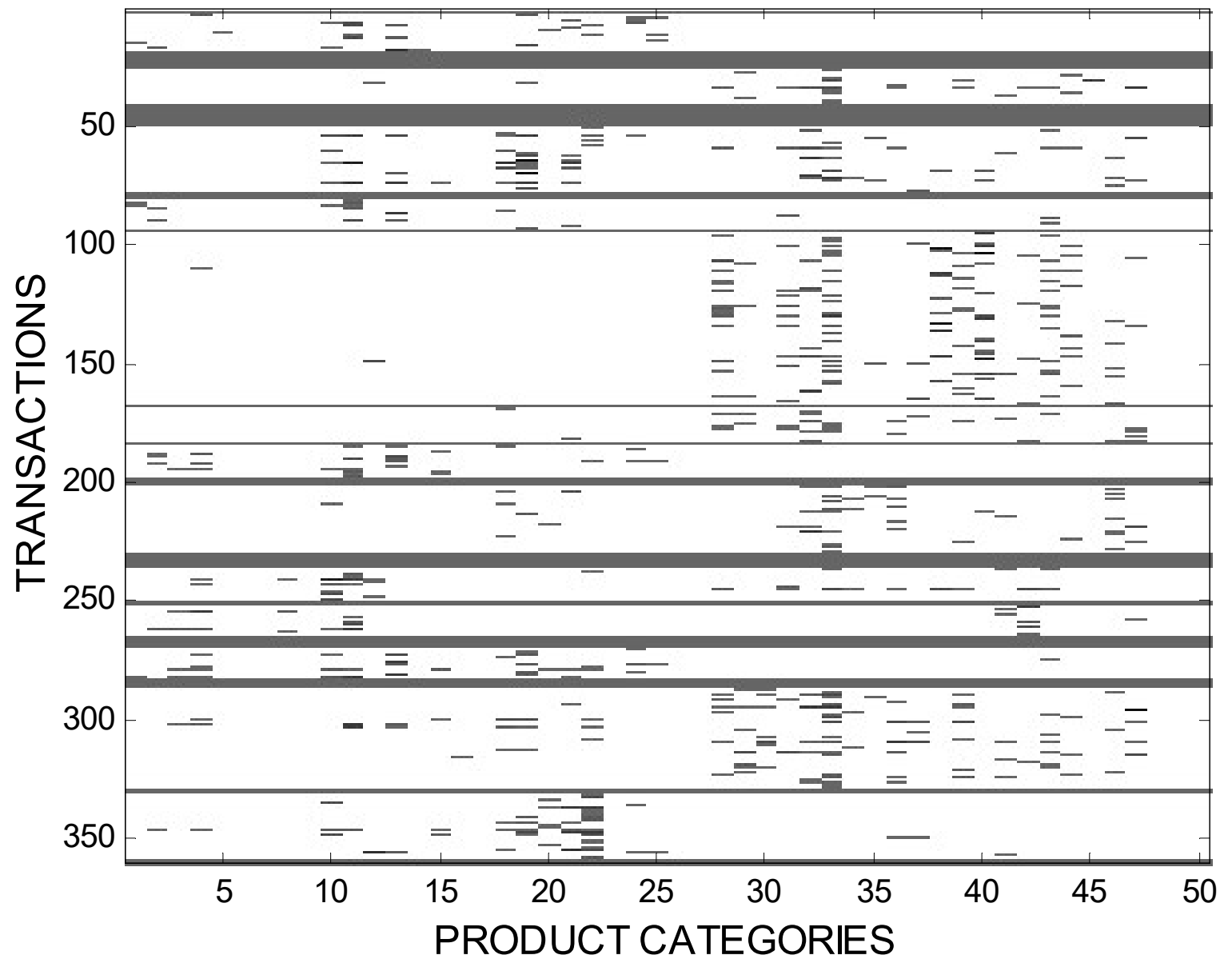
Transaction Data

- 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 constitutes a transaction, while the individual products that were purchased are the items

<i>TID</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

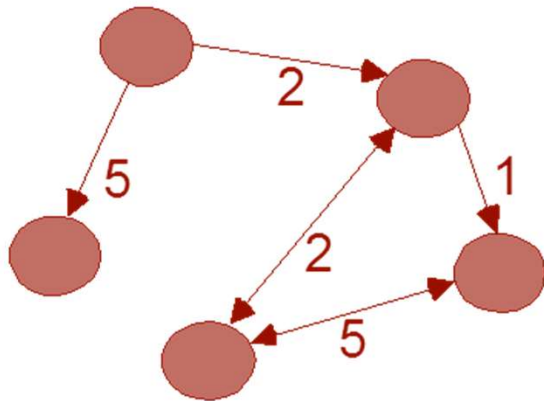


Market Basket Data



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>
```

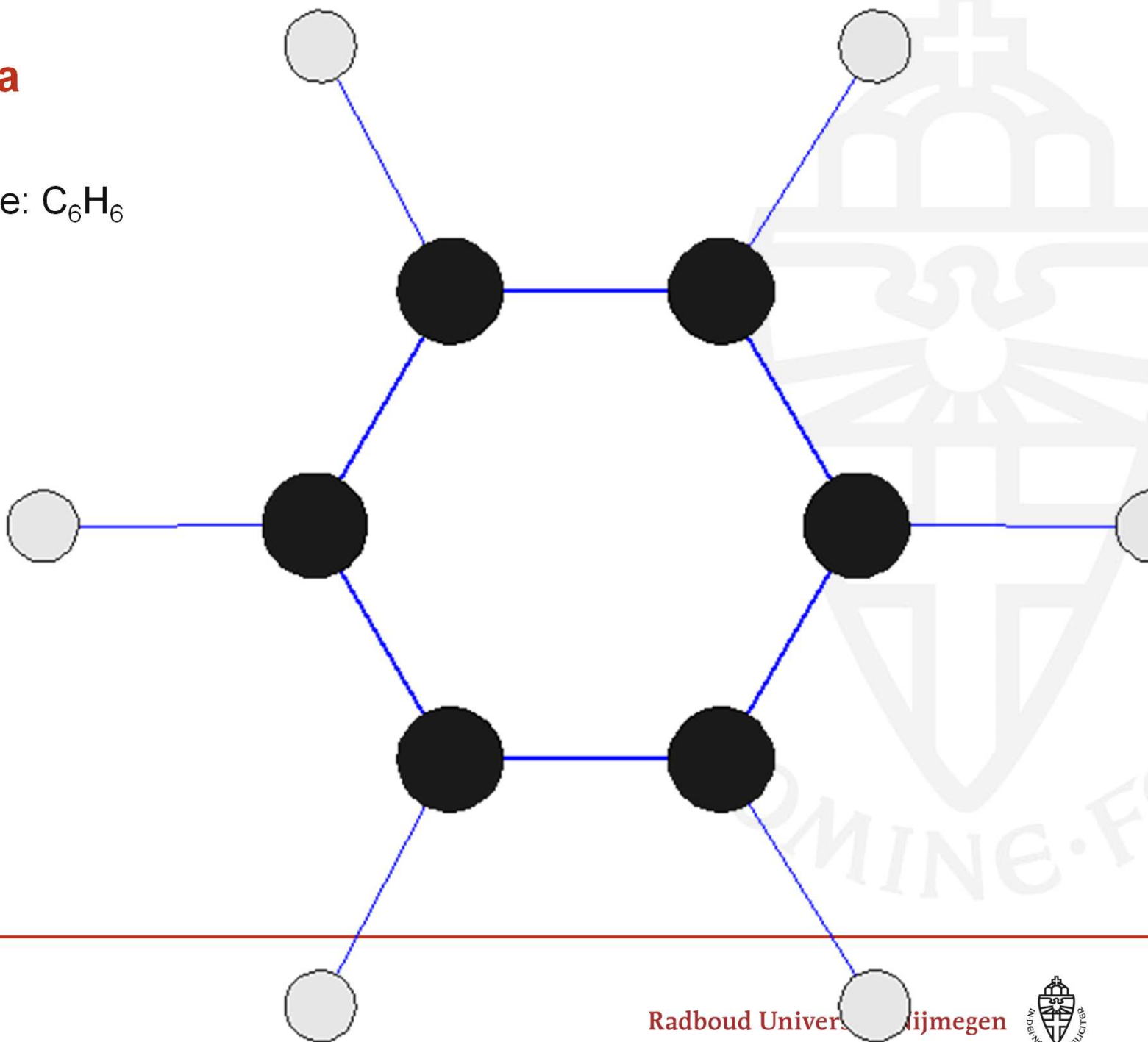
```
<li>
```

```
<a href="papers/papers.html#ffff">
```

```
N-Body Computation and Dense Linear System Solvers
```


Chemical Data

Benzene molecule: C_6H_6



Genomic Sequence Data

ADACABDABAABBDDBCADDDDDBCDDBCCBBCCDADADAADA
BDBBDABABBCDDDCDDABDCBBDBDBCBBABBBBCBBABCBB
ACBBDBAACCADDADBDBBCBBCCBBBDCABDDBBADDBBBB
CCACDABBABDDCDDBBABDBDDBDDBCACDBBCCBBACDCA
DCBACCADCCCACCDDADCBCADADBAACCDDDCBDBDCCCC
ACACACCDABDDBCADADBCBDDADABCCABDAACABCABAC
BDDDCBADCBADDDDDCDDCADCCBBADABBAADAABCCB
CABDBAADCBCDACBCABABCCBACBDABDDDADAABADCDC
CDBBCDBDADDCCBBCDBAADADBCAAAADBDCADBDBBBBCD
CCBCCCDCCADAADACABDABAABBDDBCADDDDDBCDDBCCB
BCCDADADACCCDABAABBCBDBDBADBBBBBCDADABABBDA
CDCDDDBBCDBBCBBCCDABCADDADBACBBBCCDBAAADDD
BDDCABACBCADCDCBAAADCADDADAABBACCBB

Genomic Sequence Data

ADACABDABAABBDDBCADDDDDBCDDBC**CBBC**DADADAADA
BDBBDABABBCDDDCDDABDCBBDBDBCBBABBBBCBBABCBB
ACBBDBAACCADDADBDBB**CBBC**BBBBDCABDDBBADDBBBB
CCACDABBABDDCDDBBABDBDDDBDDBCACDBBCCBBACDCA
DCBACCADCCCACCDDADCBCADADBAACCDDDCBDBDCCCC
ACACACCDABDDBCADADBCBDDADABCCABDAACABCABAC
BDDDCBADCBADDDDDCDDCADCCBBADABBAAADAAABCCB
CABDBAADCBCDACBCABABCCBACBDABDDDADAABADCDC
CDBBCDBDADDCCBBCDBAADADBCAAAADBDCADBDBBBBCD
CCBCCCDCCADAADACABDABAABBDDBCADDDDDBCDDBC**CB**
BCCDADADACCCDABAABBCBDBDBADBBBBBCDADABABBDA
CDCDDDBBCDBB**CBBC**DABCADDADBACBBBCCDBAAADDD
BDDCABACBCADCDCBAAADCADDADAABBACCBB

Spatio-Temporal Data

- Average monthly temperature of land and ocean

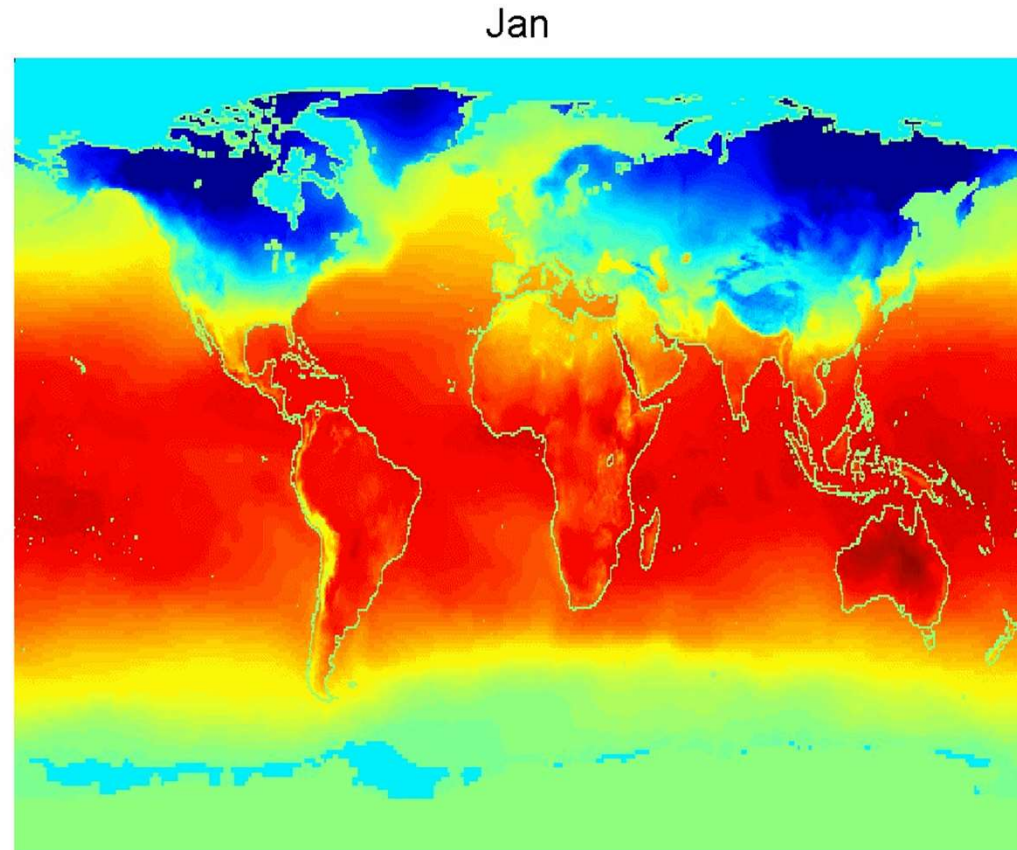


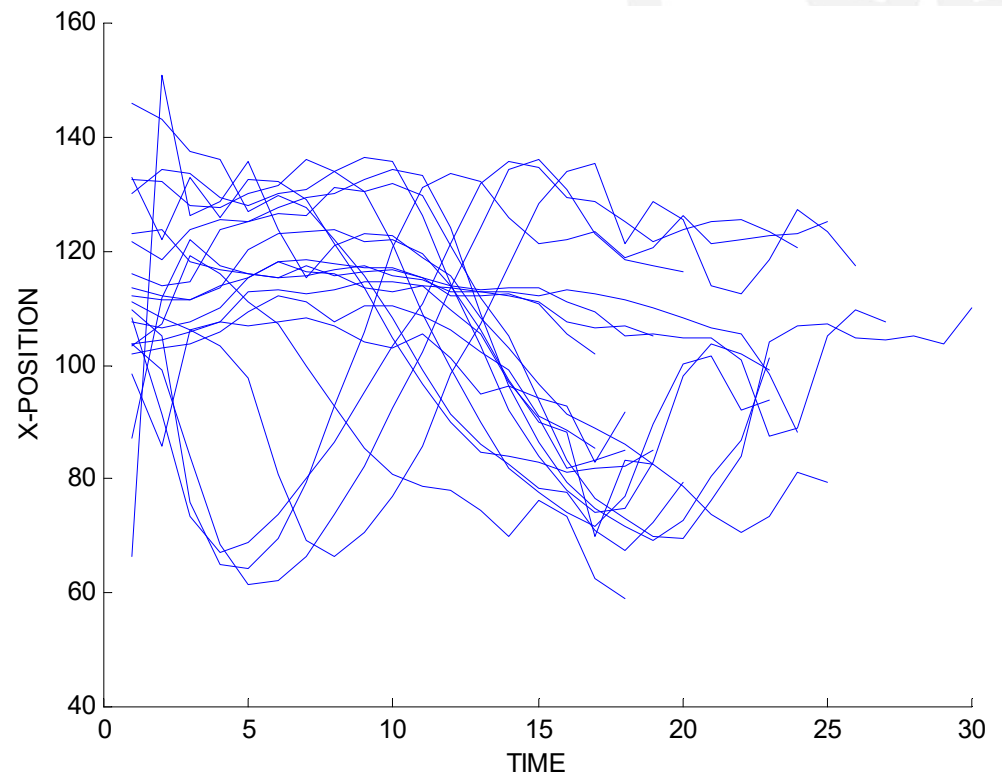
Image Data



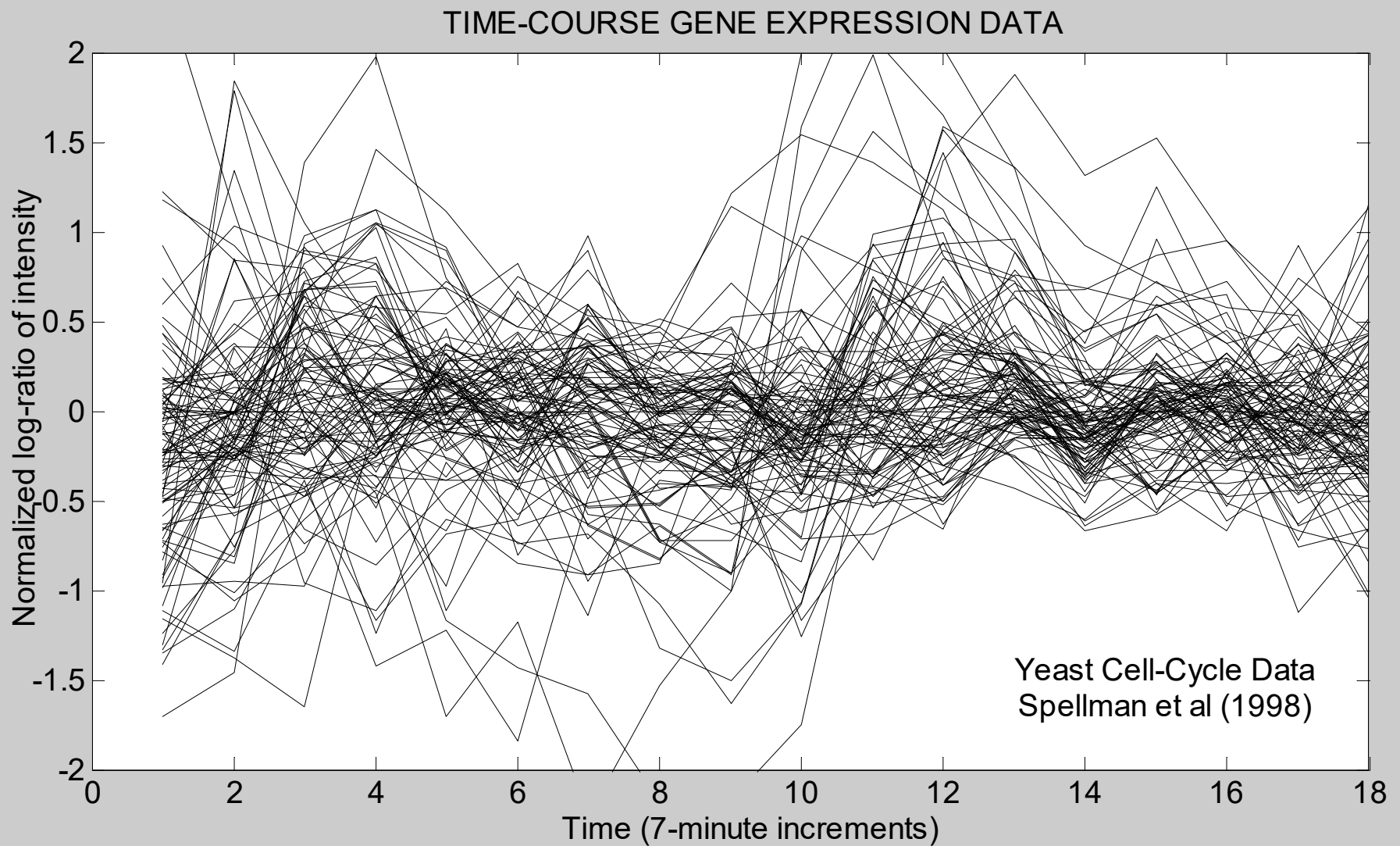
- Wikipedia

Time Series Data

Trajectories of centroids of moving hand in video streams

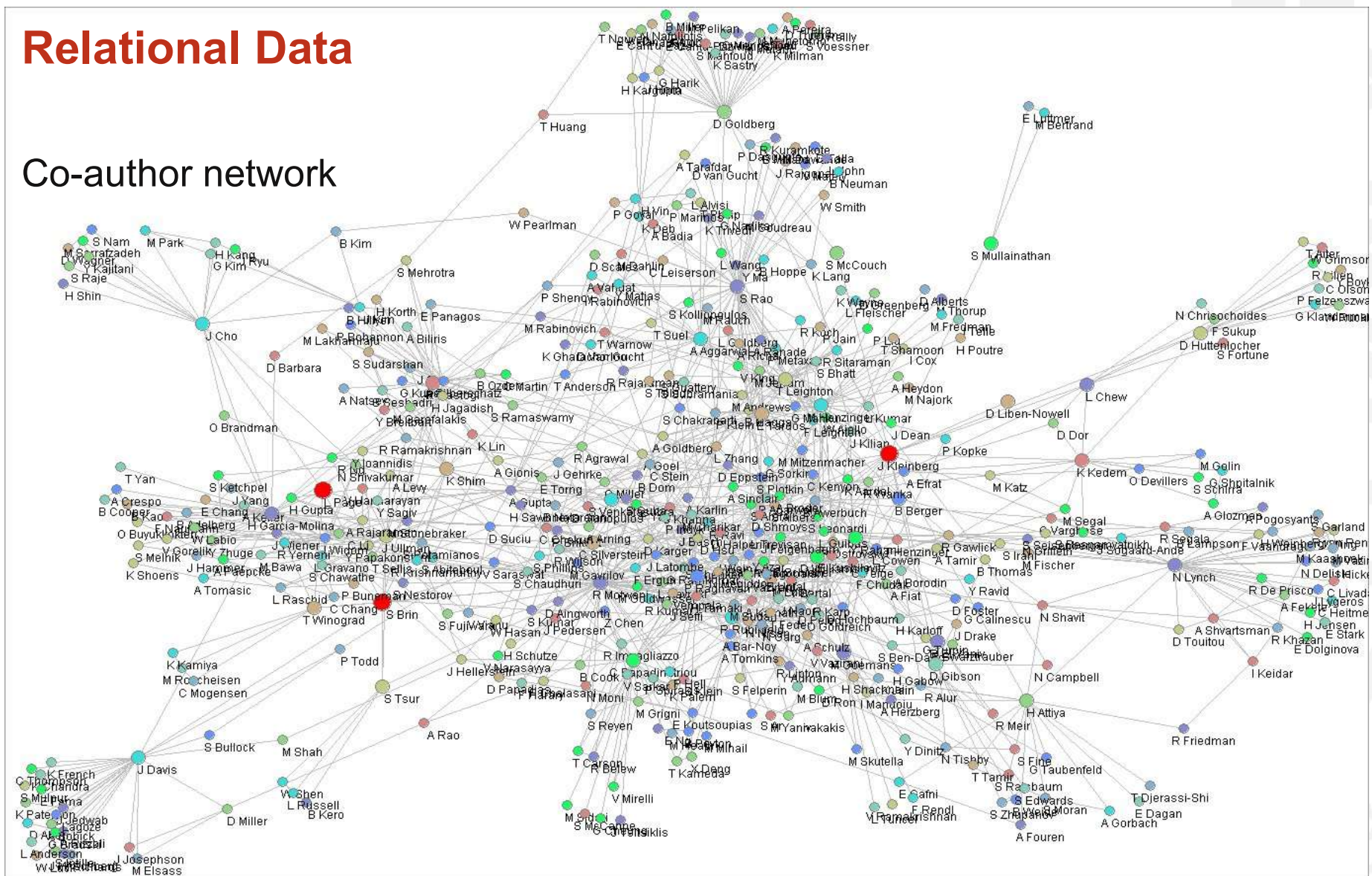


Biological Time Series

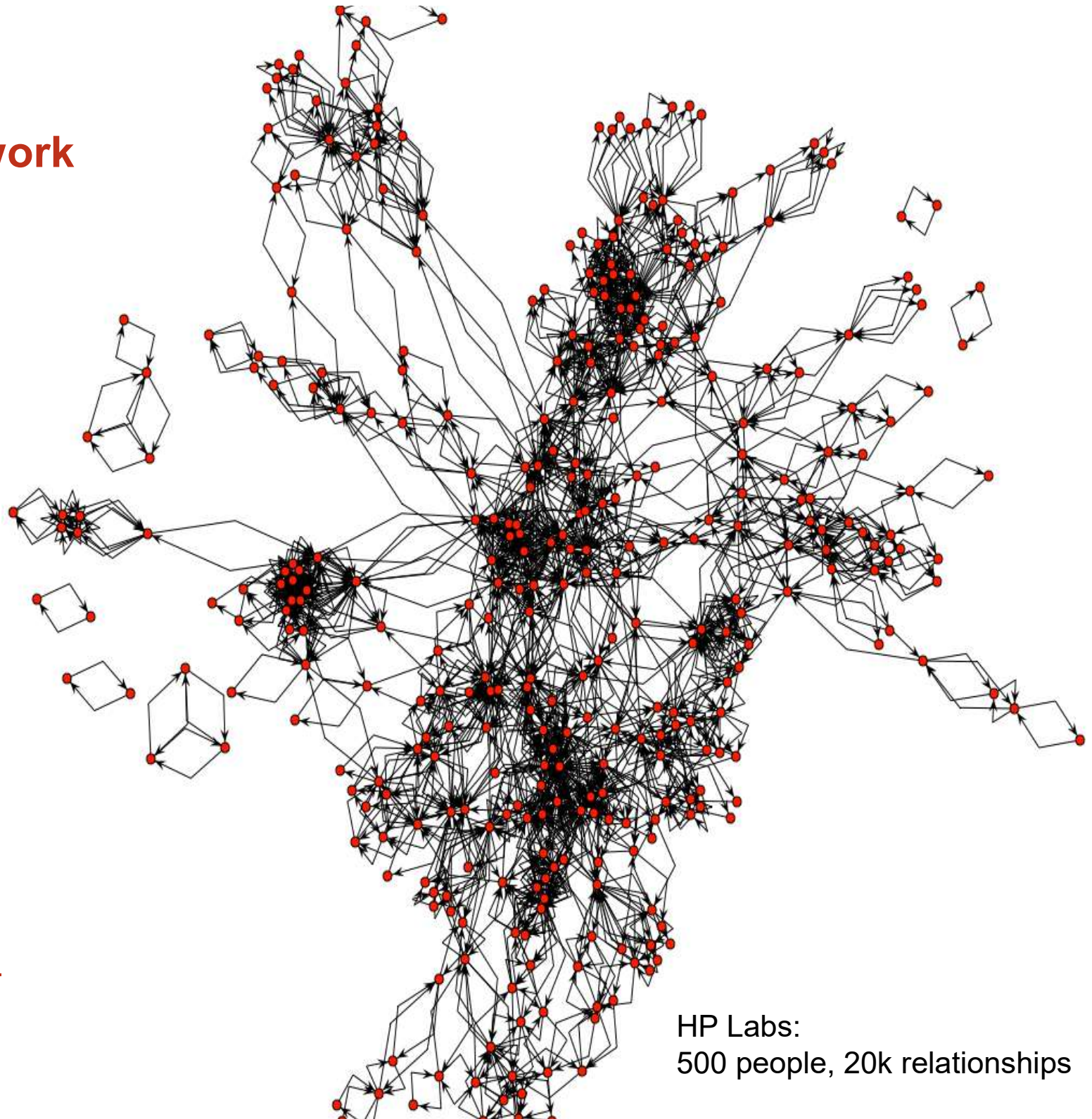


Relational Data

Co-author network



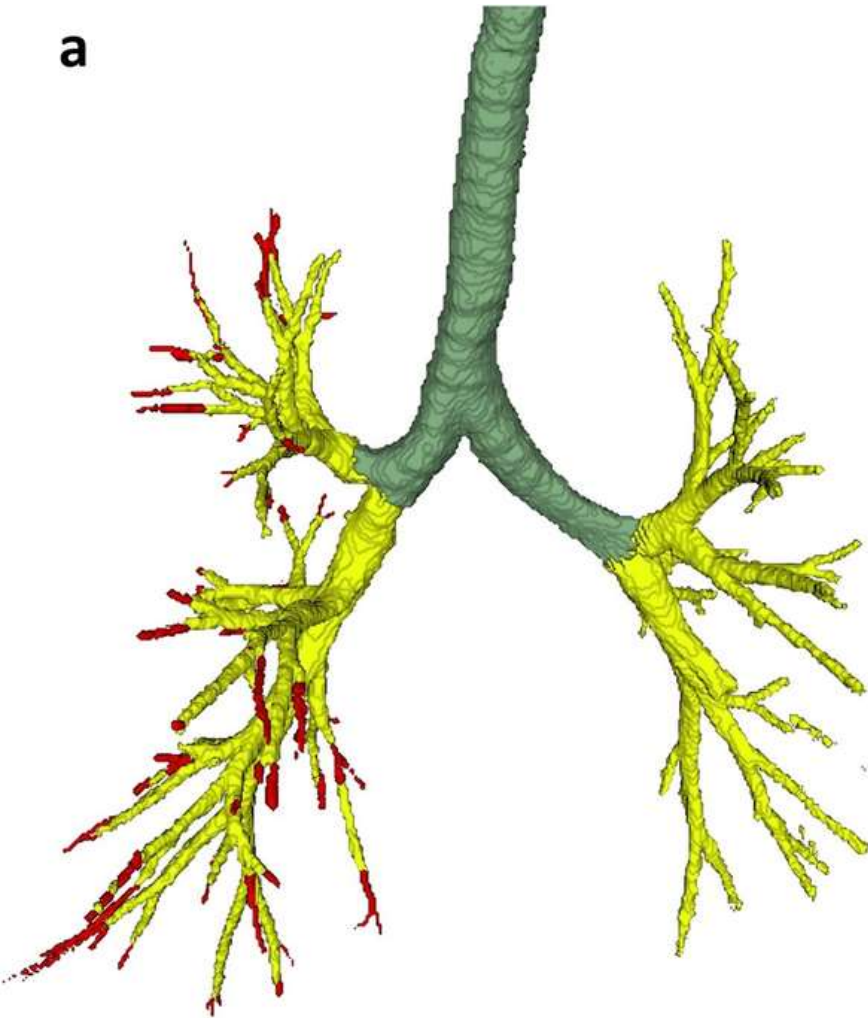
Email Network



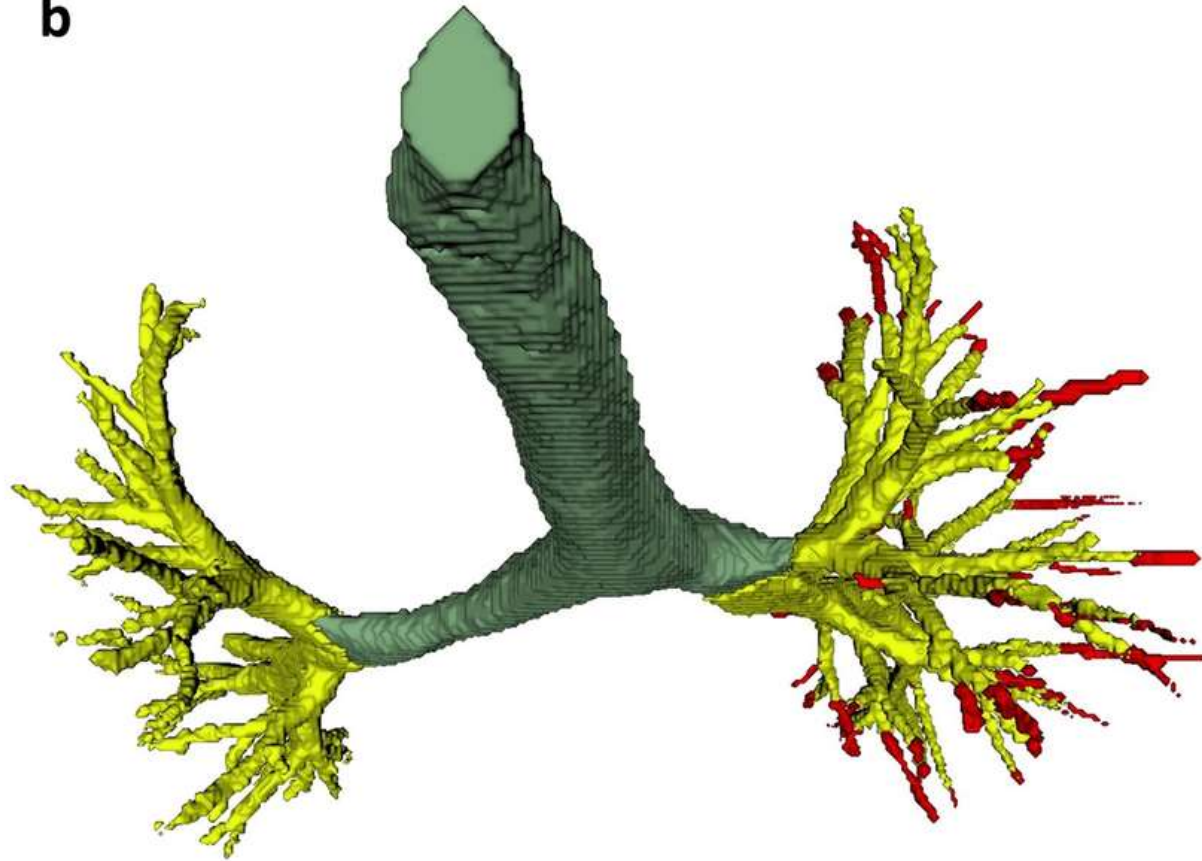
Airways



a



b



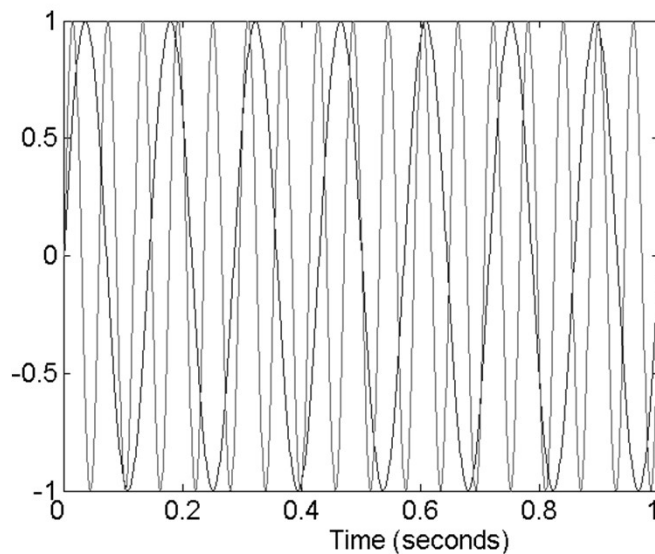
- Dudurych et al. (2021)

Data Quality

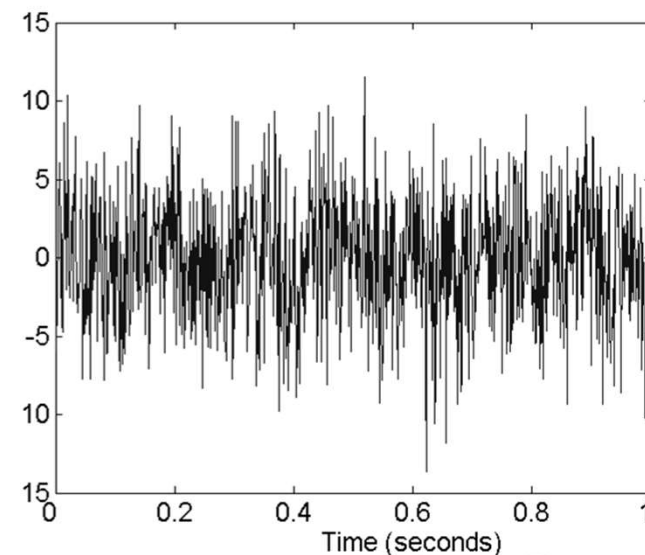
- Data is of high quality if they
 - Are fit for their intended use
 - Correctly represent the phenomena they correspond to
- Examples of data quality problems :
 - Noise and outliers
 - Missing values
 - Duplicate data

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



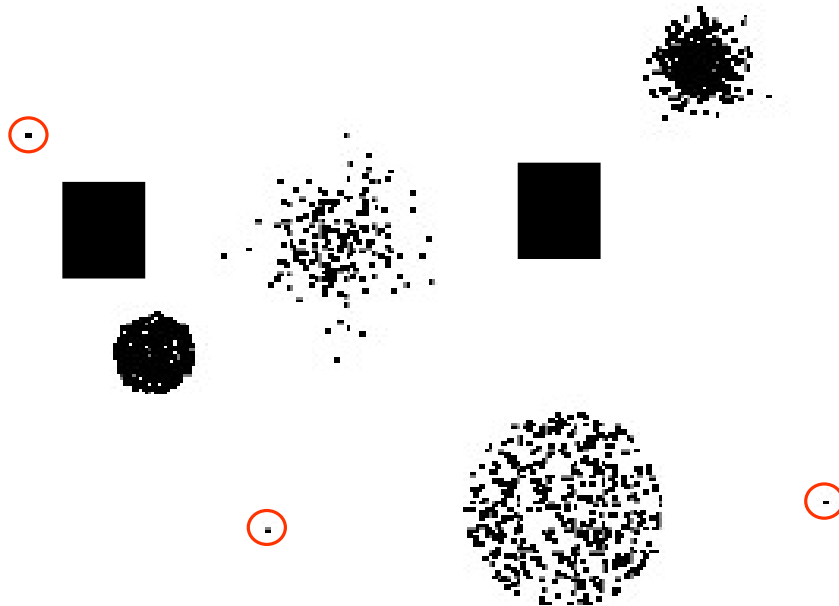
Two sine waves



Two sine waves + noise

Outliers

Outliers are data objects with characteristics that are considerably different from most [or even any?] of the other data objects in the data set



Missing Values

- Some 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., no annual income for children
- Handling missing values [some suggestions]
 - 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
- Example : same person with multiple email addresses
- This results in the need for data cleaning

Data Preprocessing

- Aggregation
- Sampling
- Dimensionality reduction
- Feature subset selection
- Feature creation
- Discretization
- Attribute transformation

Aggregation

- Combining multiple 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

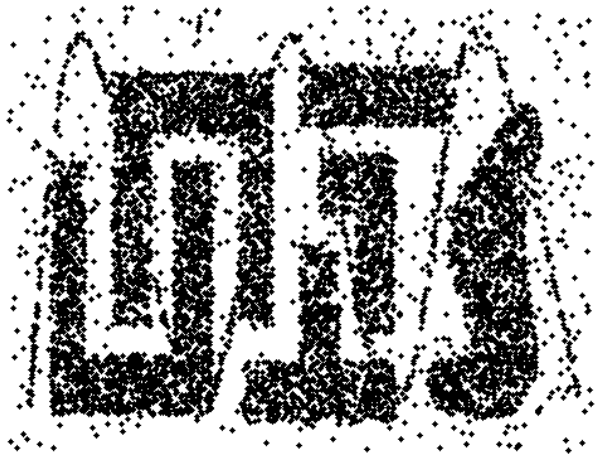
Sampling

- Sampling is the main technique employed for data selection
 - Used for both the preliminary investigation of the data and the final analysis
- Statisticians sample because
obtaining all data of interest is too expensive or time consuming
- Sampling is used in data mining because
processing all data of interest is too expensive or time consuming.
- Key principle for effective sampling
 - 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 data

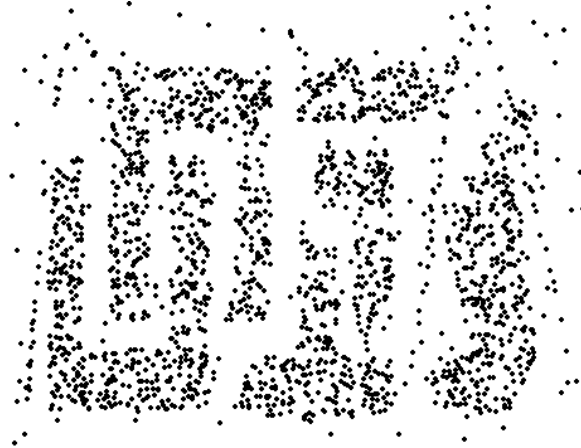
Types of Sampling

- 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
 - 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



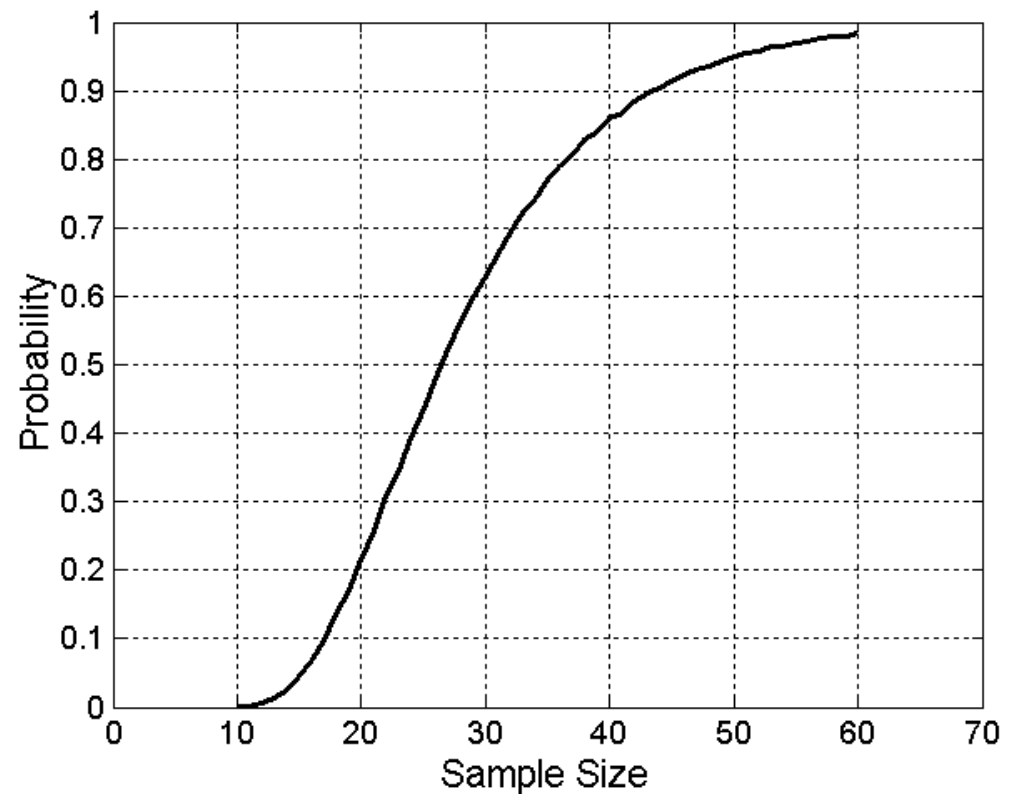
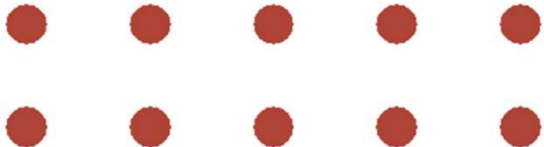
2000 Points



500 Points

Sample Size

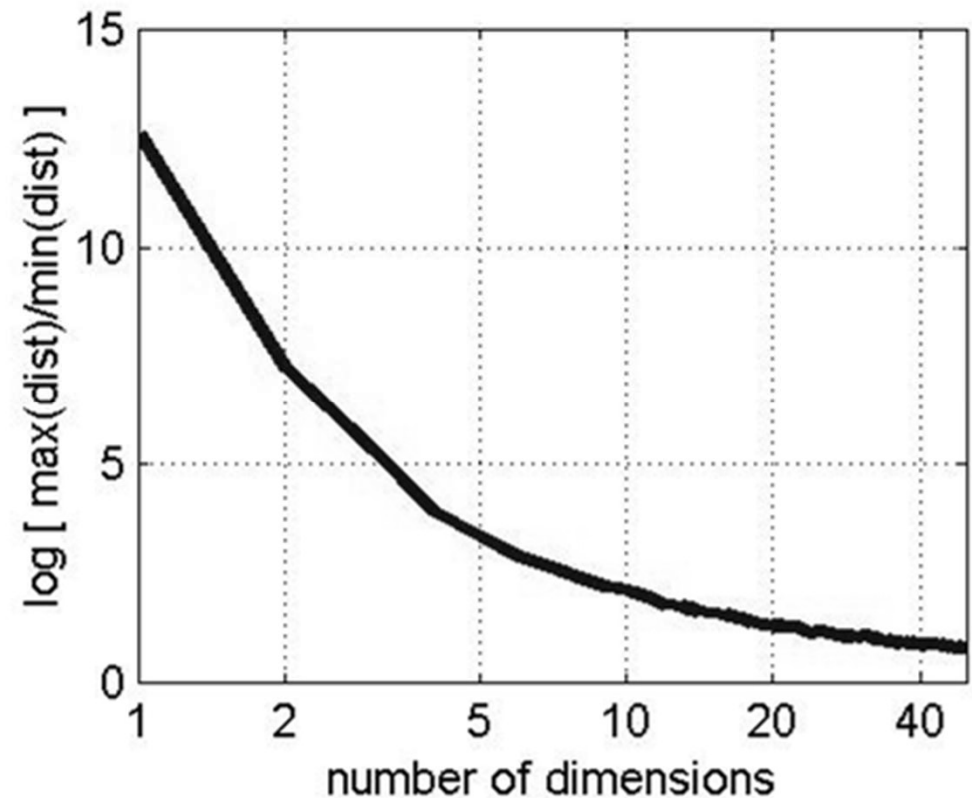
What sample size is necessary to get at least one object from each of 10 groups?



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, may become less meaningful

- Randomly generate 500 points
- Compute [log] ratio of max and min distance between any pair of points



Dimensionality Reduction

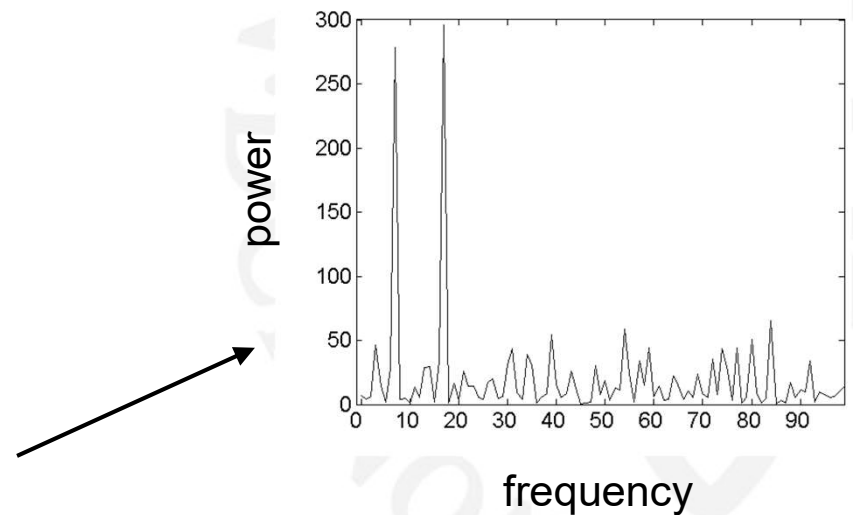
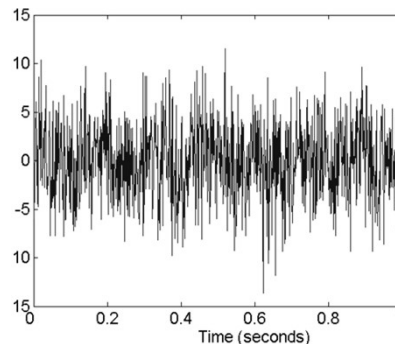
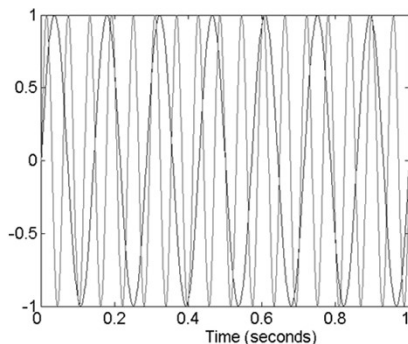
- Purpose
 - Avoid curse of dimensionality
 - Reduce amount of time and memory required by data mining algorithms
 - Allow data to be more easily visualized
 - May help to eliminate irrelevant features or reduce noise
- Techniques
 - Principal Component Analysis
 - Singular Value Decomposition
 - Others: supervised and non-linear techniques

Feature Subset Selection

- Some 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
 - Wrapper approaches :
use the data mining algorithm as a black box to find best subset of attributes

Feature Creation

- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Combining features
 - For example, BMI instead of length and weight separately
 - Particularly relevant for restricted [e.g., linear] models
- Mapping data to a new space
 - For example, Fourier transform



Attribute Transformation

- Function that maps a set of attribute values to a new set of values such that each old value can be identified with one of the new values
 - Simple functions: $x^k, \log x, e^x, |x|$
 - Standardization and normalization
- For many data mining algorithms, continuous features may preferably be more or less normally distributed
- Log transformation often useful for positive features such as income, height, etc.
- Sometimes the sign of the feature doesn't matter, sometimes the magnitude doesn't

All in All...

- Think about the objects you are dealing with
 - What do you know about them?
- Think about what your attributes describe your objects?
 - What are they? What do they mean?
- Think about the right kind of representation
 - Which scale? Which unit? Log, exp, etc.?
 - Downstream task?
 - Use prior knowledge and experience!
 - Use [your] common sense!
- Final note : curse of dimensionality and related issues will return!