



1



2

---

# Introduction

- Data is growing at a phenomenal rate
- Users expect more sophisticated information
- How?

## UNCOVER HIDDEN INFORMATION *DATA MINING*

---

# Data Mining Definition

- Finding hidden information in a database
- Fit data to a model
- Similar terms
  - Exploratory data analysis
  - Data driven discovery
  - Deductive learning

# Data Mining Algorithm

- Objective: Fit Data to a Model
  - Descriptive
  - Predictive
- Preference – Technique to choose the best model
- Search – Technique to search the data
  - “Query”



# Database Processing vs. Data Mining Processing

- |   |   |
|---|---|
| <ul style="list-style-type: none"> <li>• Query               <ul style="list-style-type: none"> <li>• Well defined</li> <li>• SQL</li> </ul> </li> <li>■ Data               <ul style="list-style-type: none"> <li>– Operational data</li> </ul> </li> <li>■ Output               <ul style="list-style-type: none"> <li>– Precise</li> <li>– Subset of database</li> </ul> </li> </ul> | <ul style="list-style-type: none"> <li>• Query               <ul style="list-style-type: none"> <li>• Poorly defined</li> <li>• No precise query language</li> </ul> </li> <li>■ Data               <ul style="list-style-type: none"> <li>– Not operational data</li> </ul> </li> <li>■ Output               <ul style="list-style-type: none"> <li>– Fuzzy</li> <li>– Not a subset of database</li> </ul> </li> </ul> |
|---|---|



# Query Examples

- Database

- Find all credit applicants with last name of Smith.
- Identify customers who have purchased more than \$10,000 in the last month.
- Find all customers who have purchased milk

- Data Mining

- Find all credit applicants who are poor credit risks. (classification)
- Identify customers with similar buying habits. (Clustering)
- Find all items which are frequently purchased with milk. (association rules)

# Basic Data Mining Tasks

- Classification maps data into predefined groups or classes
  - Supervised learning
  - Prediction
  - Regression
- Clustering groups similar data together into clusters.
  - Unsupervised learning
  - Segmentation
  - Partitioning

---

## Basic Data Mining Tasks (cont'd)

- Link Analysis uncovers relationships among data.
  - Affinity Analysis
  - Association Rules
  - Sequential Analysis determines sequential patterns.

---

## CLASSIFICATION

- Assign data into predefined groups or classes.



## But it isn't Magic

- You must know what you are looking for
- You must know how to look for you

Suppose you knew that a specific cave had gold:

What would you look for?

How would you look for it?

Might need an expert miner



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

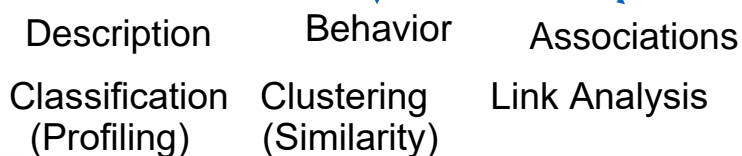
11

11

“If it looks like a duck,  
walks like a duck, and  
quacks like a duck, then  
it's a duck.”



“If it looks like a terrorist,  
walks like a terrorist, and  
quacks like a terrorist, then  
it's a terrorist.”

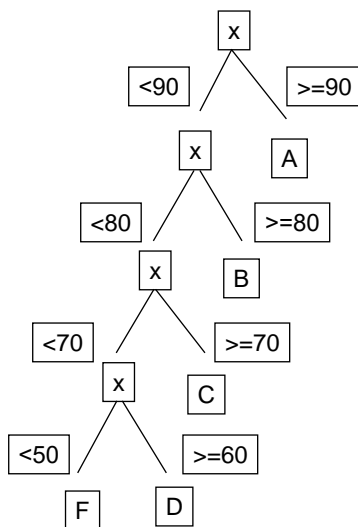


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

12

12

## Classification Ex: Grading

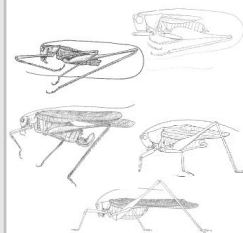


13

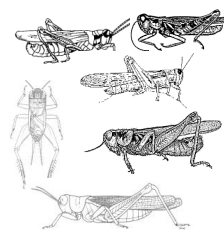
Given a collection of annotated data. (in this case 5 instances of Katydid and five of Grasshoppers), decide what type of insect the unlabeled example is.



### Katydid



### Grasshoppers



14

The classification problem can now be expressed as:

Given a training database predict the class label of a previously unseen instance

Insect ID	Abdomen Length	Antennae Length	Insect Class
1	2.7	5.5	Grasshopper
2	8.0	9.1	Katydid
3	0.9	4.7	Grasshopper
4	1.1	3.1	Grasshopper
5	5.4	8.5	Katydid
6	2.9	1.9	Grasshopper
7	6.1	6.6	Katydid
8	0.5	1.0	Grasshopper
9	8.3	6.6	Katydid
10	8.1	4.7	Katydid

previously unseen instance =

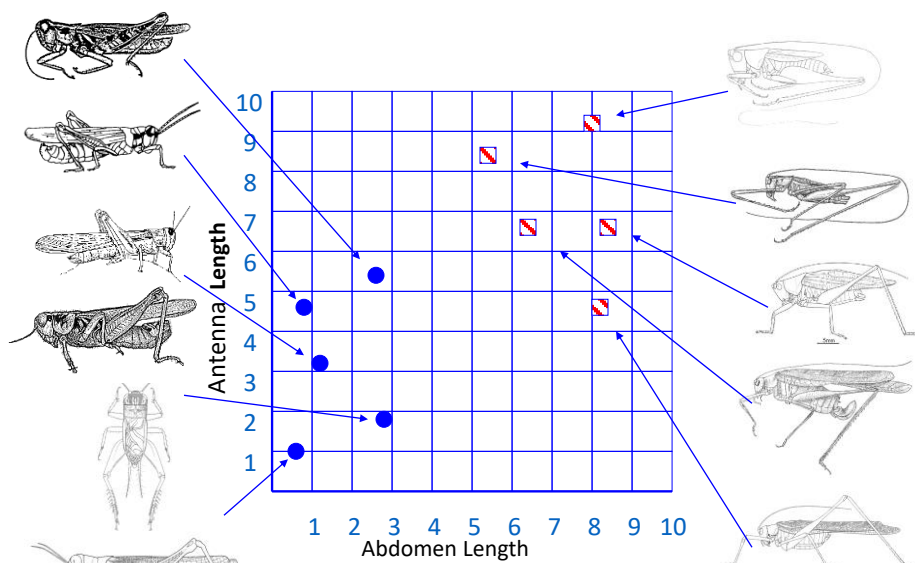
11	5.1	7.0	???????
----	-----	-----	---------



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

15

15



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

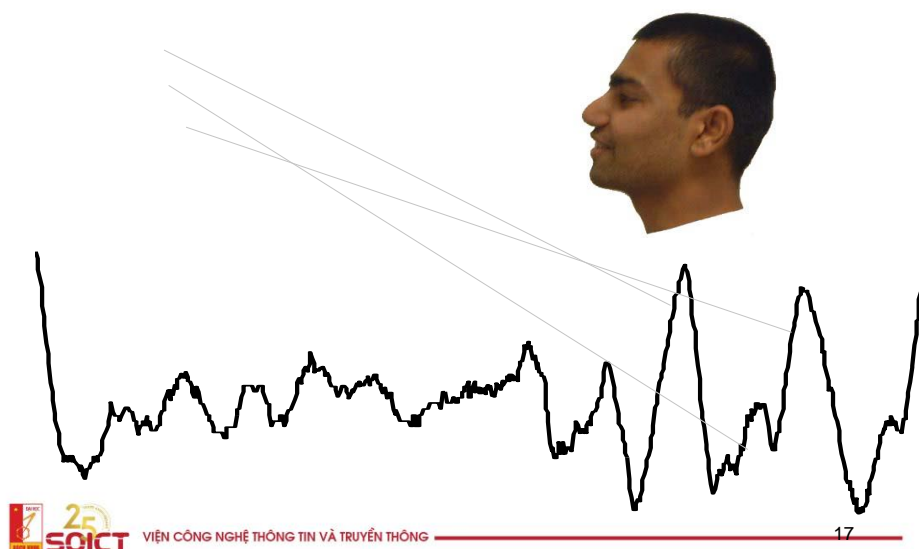
Grasshoppers

Katydid

16

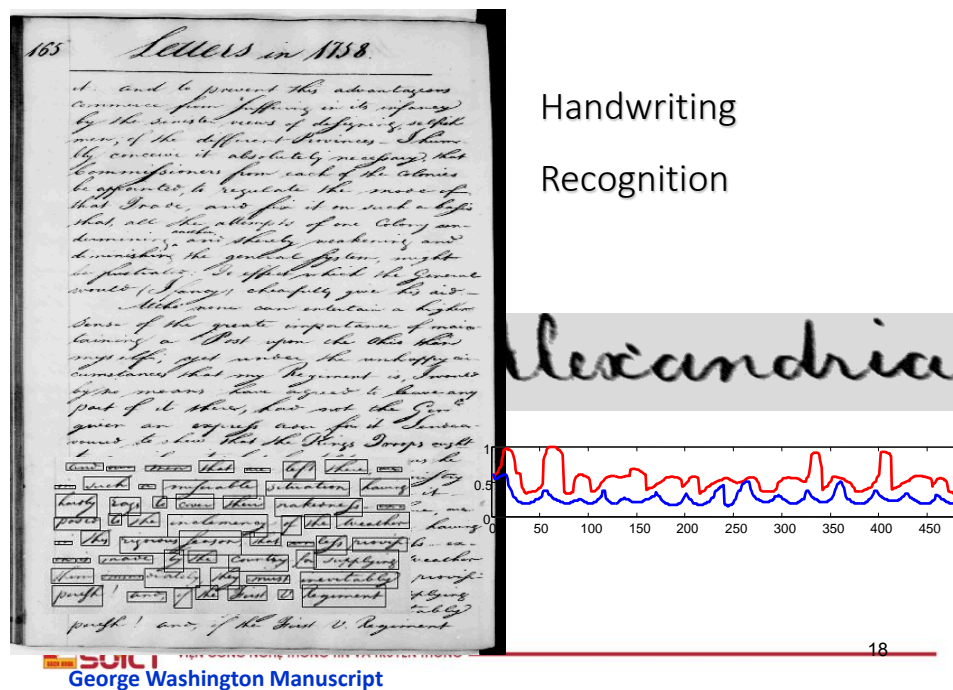


## Facial Recognition



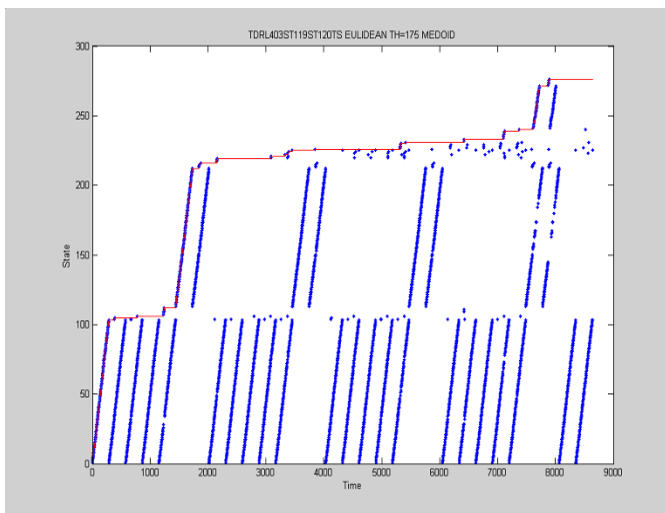
17

## Handwriting Recognition



18

# Anomaly Detection



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

19

19



But if the largest manhunt in the Washington area since Abraham Lincoln's assassination shows anything, it serves as a reminder that criminal profiling is more art than science. And that multiple killers don't always fit neatly into statistical profiling categories. Profiling the characteristics of a criminal, once dismissed as conjecture, is widely used today to help investigators solve hard-to-crack cases. "We are going to see new crimes and new techniques," says James Alan Fox, a criminologist at Northeastern University. "It's basically educated hunches."

"What people should know about profiling is that there's no magic to it," says James Alan Fox, a criminologist at Northeastern University. "It's basically educated hunches."



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

20

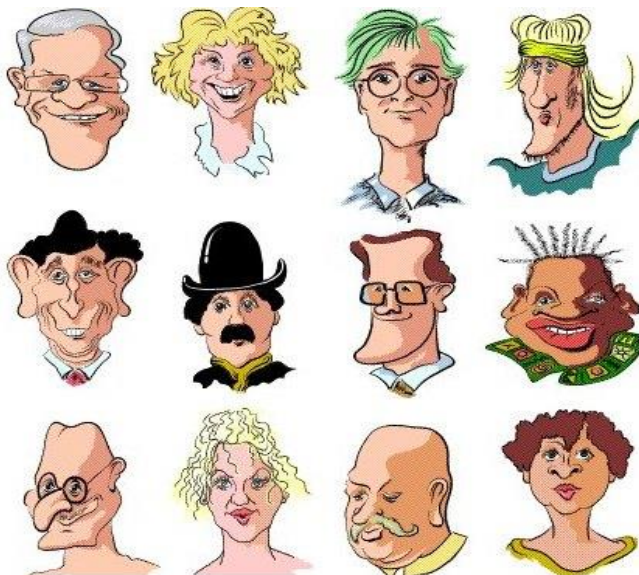
20

# CLUSTERING

- Partition data into previously undefined groups.



21



22

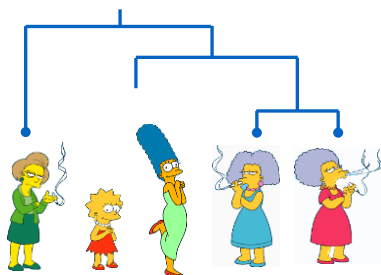
## What is Similarity?



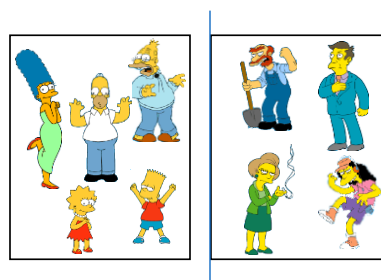
23

## Two Types of Clustering

Hierarchical



Partitional



24

# Hierarchical Clustering Example

## Iris Data Set



Sentosa



Versicolor



Virginica



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

25

25



<http://www.time.com/time/magazine/article/0,9171,1541283,00.html>



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

26

26



# Microarray Data Analysis

- Each probe location associated with gene
- Color indicates degree of gene expression
- Compare different samples (normal/disease)
- Track same sample over time
- Questions
  - Which genes are related to this disease?
  - Which genes behave in a similar manner?
  - What is the function of a gene?
- Clustering
  - Hierarchical
  - K-means



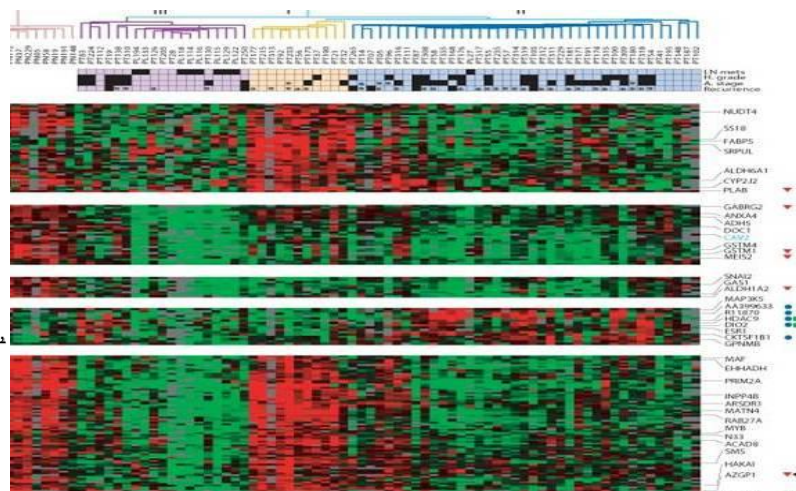
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

27

27

## Microarray Data - Clustering

"Gene expression profiling identifies clinically relevant subtypes of prostate cancer"  
[Proc. Natl. Acad. Sci. USA](#), Vol. 101, Issue 3, 811-816, January 20, 2004



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

28

28

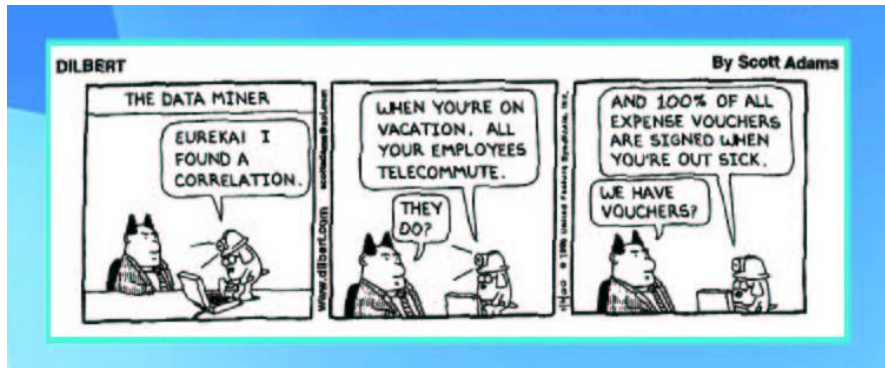
# ASSOCIATION RULES/ LINK ANALYSIS

- Find relationships between data



# ASSOCIATION RULES EXAMPLES

- People who buy diapers also buy beer
- If gene A is highly expressed in this disease then gene A is also expressed
- Relationships between people
- Book Stores
- Department Stores
- Advertising
- Product Placement



*Data Mining Introductory and Advanced Topics*, by Margaret H. Dunham, Prentice Hall, 2003.  
DILBERT reprinted by permission of United Feature Syndicate, Inc.



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

31

31

## The Dallas Morning News

Sing Newspaper

\$1.50

Dallas, Texas, Sunday, June 3, 2007

dallasnews.com

Joshua Benton and Holly K. Hacker, "At Charters, Cheating's off the Charts," *Dallas Morning News*, June 4, 2007.

Student	Shaded rows: Flagged for cheating (similar wrong answers)	Period: correct answer	Letter: incorrect answer	Dash: Skipped question
1	...C.A./C3.PD...B.AGD.C.D...JA...BG.J.PC.B...GB.B...G.			
2	...A.CGCC.C.F...A.B.AGD.C.DJ...A.B...G...A.B...GB.B...G.			
3	...A.CGCC.C.B...F...A.B...GD.C.DJ...A.B...G...A.B...GB.B...G.			
4	...A.CGCC.C.B...A...D...H...B...G...A...F...D...B...G...G.			
5	...G...F...A...F...B...J...H...B...G...A...F...D...B...G...G.			
6	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
7	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
8	...A.CGCC.C.F...B...C...D...H...A...A...D...B...G...A...G.			
9	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
10	...F...A...A...D...H...B...G...A...F...D...B...G...G.			
11	...A.CGCC.C.F...B...C...D...H...A...A...D...B...G...A...G.			
12	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
13	...B...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
14	...A.CGCC.C.F...B...C...D...H...A...A...D...B...G...A...G.			
15	...A...B...F...A...C...D...G...B...C...B...B...A...F...B...A...D...G...B...A...B			
16	...A...B...F...A...C...D...G...B...C...B...B...A...F...B...A...D...G...B...A...B			
17	...A...B...F...A...C...D...G...B...C...B...B...A...F...B...A...D...G...B...A...B			
18	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
19	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
20	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
21	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
22	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
23	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
24	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
25	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
26	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
27	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
28	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
29	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			
30	...G...F...C...B...B...F...B...G...C...D...J...G...F...B...B...D.			

\* The odds indicate the chances of having an associated pair of students in the school with such similar answers.  
SOURCE: Dallas Morning News research

TROP OXFORD/Staff Artist



VIỆN CÔNG NGHỆ TH

32



## No/Little Cheating



Joshua Benton and Holly K. Hacker, "At Charters, Cheating's off the Charts:", *Dallas Morning News*, June 4, 2007.



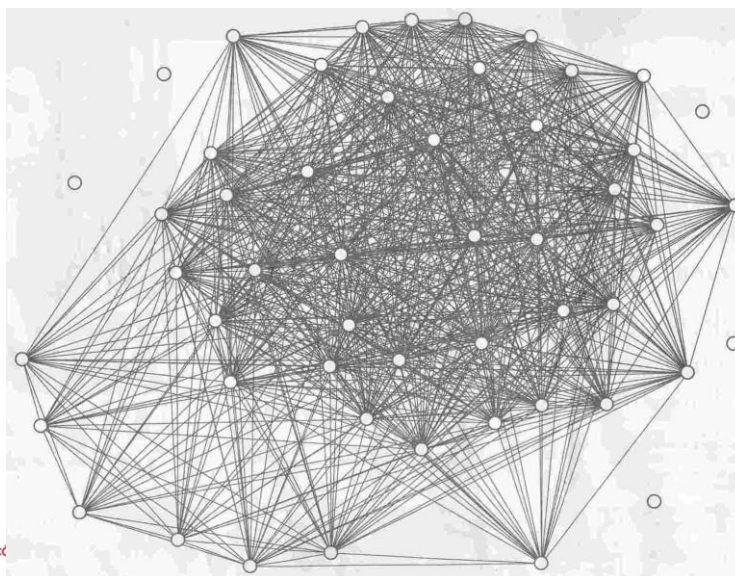
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

33

33

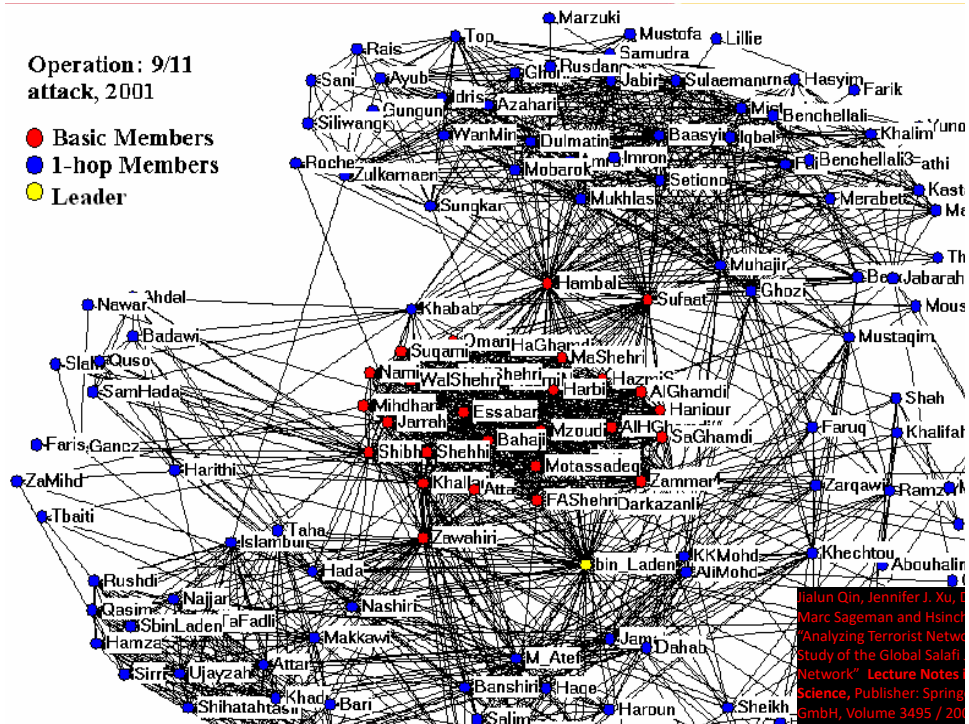
## Rampant Cheating

Joshua Benton and Holly K. Hacker, "At Charters, Cheating's off the Charts:", *Dallas Morning News*, June 4, 2007.



VIỆN C

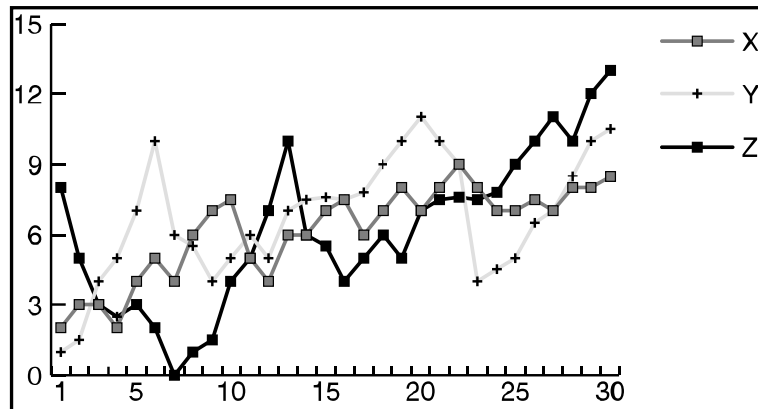
34



## Ex: Stock Market Analysis

- Example: Stock Market
- Predict future values
- Determine similar patterns over time
- Classify behavior

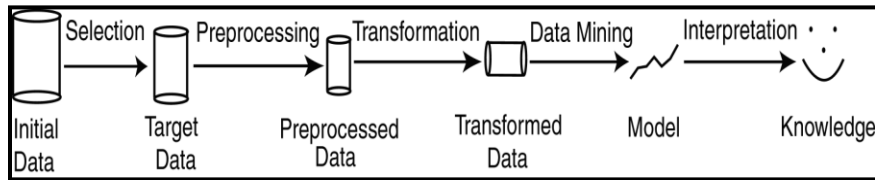
## Ex: Stock Market Analysis



## Data Mining vs. KDD

- Knowledge Discovery in Databases (KDD): process of finding useful information and patterns in data.
- Data Mining: Use of algorithms to extract the information and patterns derived by the KDD process.

## KDD Process



Modified from [FPSS96C]

- Selection: Obtain data from various sources.
- Preprocessing: Cleanse data.
- Transformation: Convert to common format.  
Transform to new format.
- Data Mining: Obtain desired results.
- Interpretation/Evaluation: Present results to user in meaningful manner.



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

39

39

## KDD Process Ex: Web Log

- Selection:
  - Select log data (dates and locations) to use
- Preprocessing:
  - Remove identifying URLs; Remove error logs
- Transformation:
  - Sessionize (sort and group)
- Data Mining:
  - Identify and count patterns; Construct data structure
- Interpretation/Evaluation:
  - Identify and display frequently accessed sequences.
- Potential User Applications:
  - Cache prediction
  - Personalization



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

40

40

## Related Topics

- Databases
- OLTP
- OLAP
- Information Retrieval



## DB & OLTP Systems

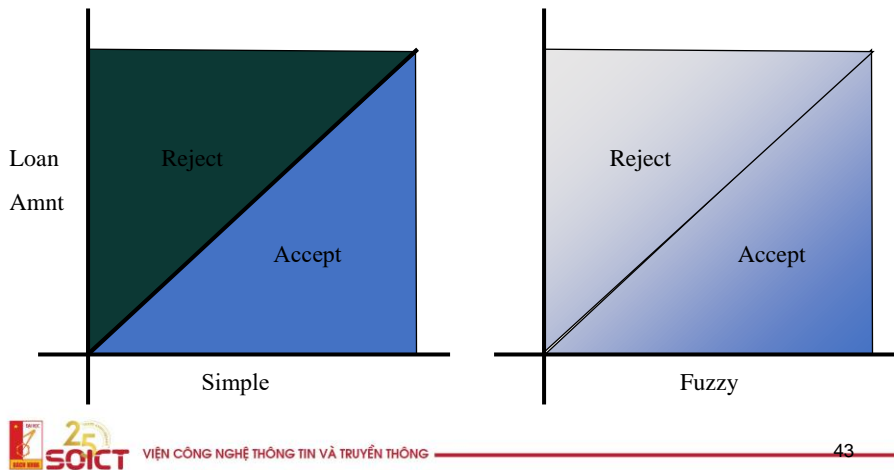
- Schema
  - (ID,Name,Address,Salary,JobNo)
- Data Model
  - ER
  - Relational
- Transaction
- Query:
 

```
SELECT Name
FROM T
WHERE Salary > 100000
```

*DM: Only imprecise queries*



## Classification/Prediction is Fuzzy



43

## Information Retrieval

- **Information Retrieval (IR):** retrieving desired information from textual data.
- Library Science
- Digital Libraries
- Web Search Engines
- Traditionally keyword based
- Sample query:  
Find all documents about “data mining”.

**DM: Similarity measures;  
Mine text/Web data.**

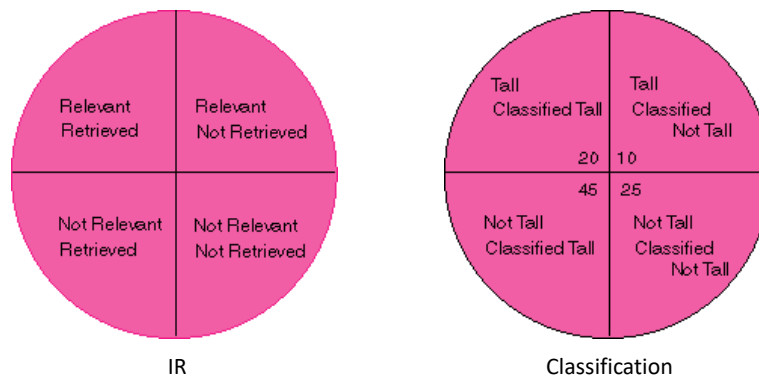
44

## Information Retrieval (cont'd)

- **Similarity:** measure of how close a query is to a document.
- Documents which are “close enough” are retrieved.
- Metrics:
  - **Precision** =  $\frac{|\text{Relevant and Retrieved}|}{|\text{Retrieved}|}$
  - **Recall** =  $\frac{|\text{Relevant and Retrieved}|}{|\text{Relevant}|}$



## IR Query Result Measures and Classification



# OLAP

- **Online Analytic Processing (OLAP):** provides more complex queries than OLTP.
- **OnLine Transaction Processing (OLTP):** traditional database/transaction processing.
- Dimensional data; cube view
- Visualization of operations:
  - *Slice:* examine sub-cube.
  - *Dice:* rotate cube to look at another dimension.
  - *Roll Up/Drill Down*

*DM: May use OLAP queries.*



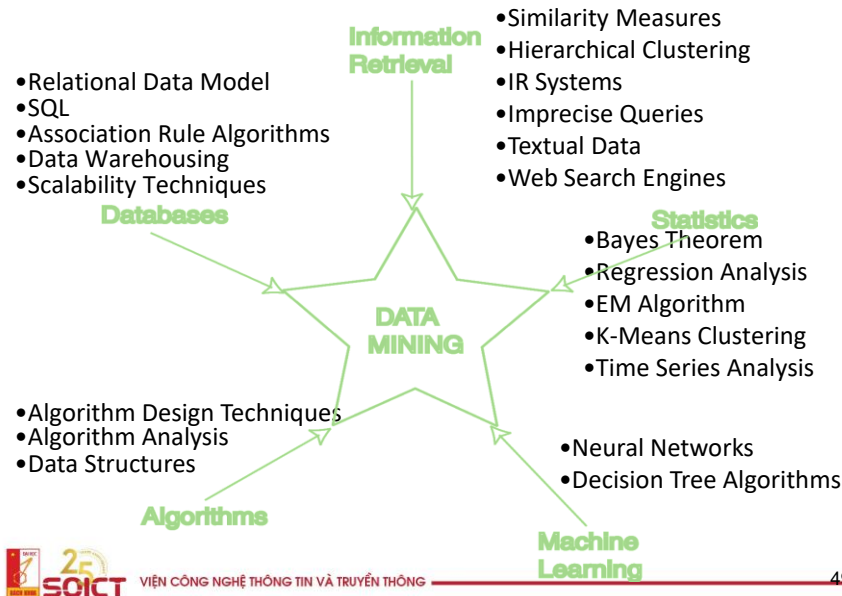
## DM vs. Related Topics

Area	Query	Data	Results	Output
DB/OLTP	Precise	Database	Precise	DB Objects or Aggregation
IR	Precise	Documents	Vague	Documents
OLAP	Analysis	Multidimensional	Precise	DB Objects or Aggregation
DM	Vague	Preprocessed	Vague	KDD Objects





# Data Mining Development



49

## KDD Issues

- Human Interaction
- Overfitting
- Outliers
- Interpretation
- Visualization
- Large Datasets
- High Dimensionality

50

# Overfitting

- Suppose we want to predict whether an individual is short, medium, or tall. What is wrong with this data?

Name	Gender	Height	Output
Mary	F	1.6	Short
Maggie	F	1.9	Medium
Martha	F	1.88	Medium
Stephanie	F	1.7	Short
Bob	M	1.85	Medium
Kathy	F	1.6	Short
George	M	1.7	Short
Debbie	F	1.8	Medium
Todd	M	1.95	Medium
Kim	F	1.9	Medium
Amy	F	1.8	Medium
Wynette	F	1.75	Medium

51

## KDD Issues (cont'd)

- Multimedia Data
- Missing Data
- Irrelevant Data
- Noisy Data
- Changing Data
- Integration
- Application

52

# WARNING

- With data mining you don't always know what you are looking for.
- There is not one right answer.
- The data you are using is noisy
- Data Mining is a very applied discipline.
- A data mining course provides you tools to use to analyze data.
- Experience provides you knowledge of how to use these tools.



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

53

53



<http://ieeexplore.ieee.org/iel5/6/32236/01502526.pdf?tp=&number=1502526&isnumber=32236>

54



## Social Implications of DM

- Privacy
- Profiling
- Unauthorized use
- Invalid results and claims

---

## Data Mining Metrics

- Usefulness
- Return on Investment (ROI)
- Accuracy
- ...
- Space/Time



---

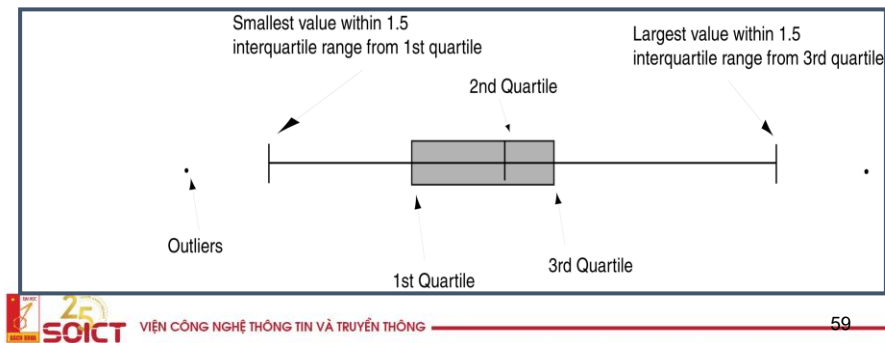
## Visualization Techniques

- Graphical
- Geometric
- Icon-based
- Pixel-based
- Hierarchical
- Hybrid



## Models Based on Summarization

- Visualization: Frequency distribution, mean, variance, median, mode, etc.
- Box Plot:



59

## DM Tools

- XLMiner – Easy addin to Excel  
<http://www.solver.com/xlminer/index.html>
- Weka – Open Source; Visualization, Functionality, Interface  
<http://www.cs.waikato.ac.nz/ml/weka/>
- SAS (JMP) – Commercial Product
- SPSS – Commercial Product
- MATLAB – Statistical/Math Applications
- R – Programming

61



62



63