

# Principles of Data Mining

Instructor: Sargur N. Srihari

University at Buffalo  
The State University of New York  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Introduction: Topics

1. Introduction to Data Mining
2. Nature of Data Sets
3. Types of Structure  
Models and Patterns
4. Data Mining Tasks (What?)
5. Components of Data Mining Algorithms(How?)
6. Statistics vs Data Mining

# Flood of Data

New York Times, January 11, 2010

## Military Is Awash in Data From Drones

By CHRISTOPHER DREW

So much video intelligence is coming out of Afghanistan that military analysts are finding it hard to keep up.



Video and Image Data  
“Unstructured”

Government agencies are still having trouble making sense of the flood of data they collect for intelligence purposes, a point underscored by the 9/11 Commission and, more recently, by [President Obama](#) after the attempted bombing of a Detroit-bound passenger flight on Christmas Day.

“Structured and Unstructured”  
(Text) Data

# Large Data Sets are Ubiquitous

1. Due to advances in digital data acquisition and storage technology

## Business

- Supermarket transactions
- Credit card usage records
- Telephone call details
- Government statistics

## Scientific

- Images of astronomical bodies
- Molecular databases
- Medical records

International organizations produce more information in a week than many people could read in a lifetime

2. *Automatic data production* leads to need for *automatic data consumption*

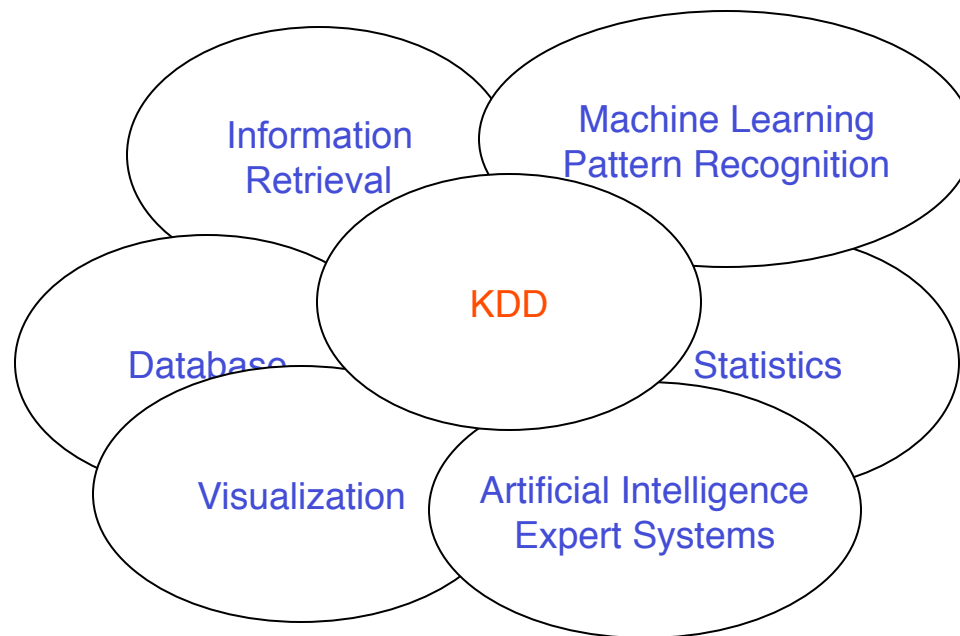
3. Large databases mean vast amounts of information

4. Difficulty lies in accessing it

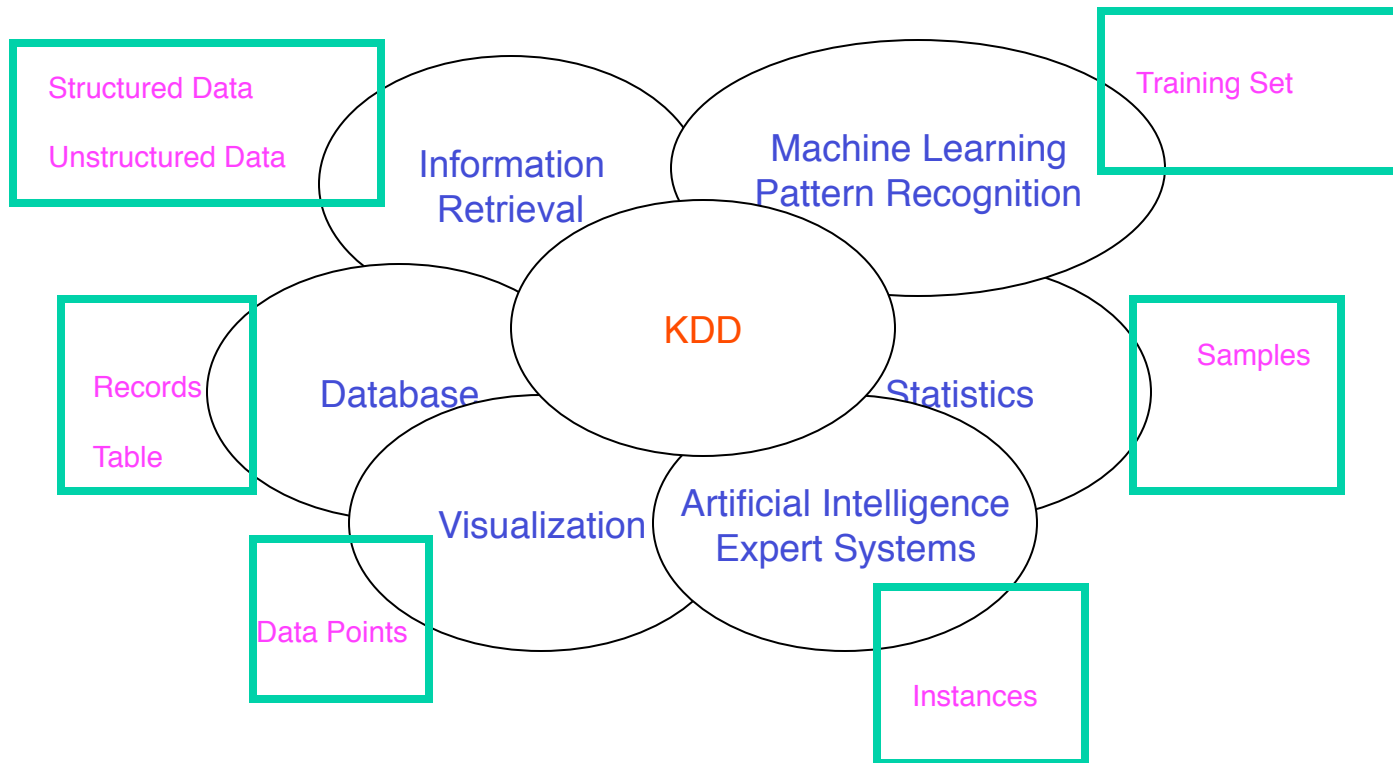
# Data Mining as Discovery

- Data Mining is
  - Science of extracting useful information from large data sets or databases
- Also known as KDD
  - Knowledge Discovery and Data Mining
  - Knowledge Discovery in Databases

# KDD is a multidisciplinary field



# Terminology for Data



# Data Mining Definition

Analysis of (often large) *Observational Data* to find unsuspected relationships and Summarize data in novel ways that are understandable and useful to data owner

## Unsuspected Relationships

non-trivial, implicit, previously unknown

Ex of Trivial: Those who are pregnant are female

## Relationships and Summary

are in the form of Patterns and Models

Linear Equations, Rules, Clusters, Graphs, Tree Structures, Recurrent Patterns in Time Series

## Usefulness:

meaningful: lead to some advantage, usually economic

## Analysis:

Process of discovery (Extraction of knowledge)

Automatic or Semi-automatic



# Observational Data

- **Observational Data**
  - Objective of data mining exercise plays no role in data collection strategy
  - E.g., Data collected for Transactions in a Bank
- **Experimental Data**
  - Collected in Response to Questionnaire
  - Efficient strategies to Answer Specific Questions
- In this way it differs from much of statistics
- For this reason, data mining is referred to as secondary data analysis

# KDD Process

- Stages:
  - Selecting Target Data
  - Preprocessing
  - Transforming them
  - Data Mining to Extract Patterns and Relationships
  - Interpreting Assesses Structures
- KDD more complicated than initially thought
  - 80% preparing data
  - 20% mining data

# Seeking Relationships

- Finding accurate, convenient and useful representations of data involves these steps:
  - Determining nature and structure of representation
    - E.g., linear regression
  - Deciding how to quantify and compare two different representation
    - E.g., sum of squared errors
  - Choosing an algorithmic process to optimize score function
    - E.g., gradient descent optimization
  - Efficient Implementation using data management

# Example of Regression Analysis

1. Representation
2. Score function
3. Process to optimize score
4. Implementation:  
data management,  
efficiency

## EXAMPLE of Model

1. Regression:

$$y = a + bx$$

Predictor variable =  $x$   
(income)

Response variable =  $y$   
(credit card spending)

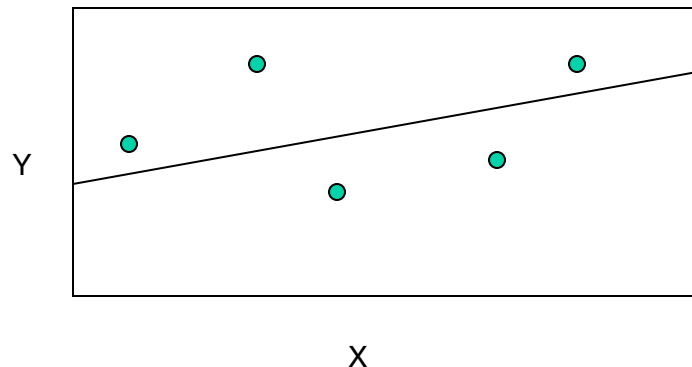
2. Score: sum of squared errors

# Linear Regression Process: Extracting a Linear Model

## Linear regression with one variable

Data of the form  $(x_i, y_i)$ ,  $i = 1, \dots, n$  samples

Need to find  $a$  and  $b$  such that  $y = a + bx$



What is involved in calculating  $a$  and  $b$   
So that the line fits the points the best?

Data Representation

$y$	$x$
1	3
8	9
11	11
4	5
3	2

# Score: Sum of Squared Errors

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - y'_i)^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

Where  $y_i$  is the response value obtained from the model

We wish to minimize SSE

# Minimizing SSE for Regression

Differentiating SSE with respect to  $a$  and  $b$  we have

$$\frac{\partial(SSE)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i)$$

$$\frac{\partial(SSE)}{\partial b} = -2 \sum_{i=1}^n (y_i - a - bx_i)x_i$$

Setting partial derivatives equal to zero and rearranging terms

$$na + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

Which we solve for  $a$  and  $b$ ,  
the *regression coefficients*

# Regression Coefficients

$$b = \frac{\sum_{i=1}^n (x_i - \text{mean}_x) \cdot (y_i - \text{mean}_y)}{\sum_{i=1}^n (x_i - \text{mean}_x)^2}$$

$$a = \text{mean}_y - b \cdot \text{mean}_x$$

To calculate  $a$  and  $b$  we need to find the *means* of the  $x$  and  $y$  values. Then we calculate  $b$  as a function of the  $x$  and  $y$  values and the means

$a$  as a function of the means and  $b$



# Application to Data

$y$	$x$
1	3
8	9
11	11
4	5
3	2

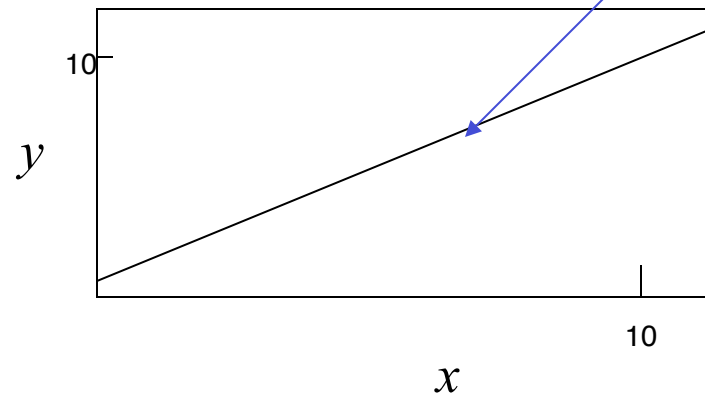
17

$$\text{mean}_y = 5 \qquad \text{mean}_x = 6$$

$$a = 0.8, b = 1.04$$

Optimal regression line is

$$y = 0.8 + 1.04x$$



Linear Regression  
For the data set

# Multiple Regression

$p$  predictor variables

$y$	$x_1$	$x_2$	.....	$x_p$
$y(1)$	$x_1(1)$			
$y(n)$	$x_1(n)$			

$$y = a_0 + \sum_{j=1}^p a_j x_j$$

$n$  objects

$X = n \times d+1$  matrix

Where a column of 1's are added to incorporate  $a_0$  in model

$$y = Xa + e$$

$y$  is a column vector,  $a = (a_0, \dots, a_p)$   
 $e$  is a  $n$  by 1 vector containing residuals

Solution:

18

$$a = (X^T X)^{-1} \cdot X^T y$$

# Implementation of Regression

$$y = a_0 + \sum_{j=1}^p a_j x_j$$

Solution:

$$a = (X^t X)^{-1} \cdot X^t y$$

Simple summaries of the data; sums, sums of squares and sums of products of  $X$  and  $Y$  are sufficient to compute estimates of  $a$  and  $b$

Implies single pass through the data will yield estimates

## 2. Nature of Data Sets

- Structured Data
  - set of measurements from an environment or process
- Simple case
  - $n$  objects with  $d$  measurements each:  $n \times d$  matrix
  - $d$  columns are called *variables, features, attributes or fields*

# Structured Data and Data Types

## US Census Bureau Data

### Public Use Microdata Sample data sets (PUMS)

ID	Age	Sex	Marital Status	Education	Income
	Quantitative Continuous	Categorical Nominal		Categorical Ordinal	
248	54	Male	Married	High School grad	100000
	??	Missing data			Noisy data A guess?
249	??	Female	Married	HS grad	12000
250	29	Male	Married	Some College	23000
251	9	Male	Not Married	Child	0

PUMS Data has identifying information removed.

Available in 5% and 1% sample sizes. 1% sample has 2.7 million records

# Unstructured Data

## 1. Structured Data

- Well-defined tables, attributes (columns), tuples (rows)
- UC Irvine data set

## 2. Unstructured Data

- World wide web
  - Documents and hyperlinks
    - HTML docs represent tree structure with text and attributes embedded at nodes
    - XML pages use metadata descriptions
- Text Documents
  - Document viewed as sequence of words and punctuations
    - Mining Tasks
      - » Text categorization
      - » Clustering Similar Documents
      - » Finding documents that match a query
      - » Automatic Essay Scoring (AES)
    - Reuters collection is at <http://www.research.att.com/~lewis>

# Representations of Text Documents

- Boolean Vector

- Document is a vector where each element is a bit representing presence/absence of word
- A set of documents
  - can be represented as matrix  $(d,w)$ 
    - where document  $d$  and word  $w$  has value 1 or 0  
(sparse matrix)

- Vector Space Representation

- Each element has a value such as no. of occurrences or frequency
- A set of documents represented as a document-term matrix

# Vector Space Example

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

t1	database
t2	SQL
t3	index
t4	regression
t5	likelihood
t6	linear

$d_{ij}$  represents number of times that term appears in that document



# Mixed Data: Structured & Unstructured

## Medical Patient Data

- Blood Pressure at different times of day
- Image data (x-ray or MRI)
- Specialist's comments (text)
- Hierarchy of relationships between patients, doctors, hospitals

$N \times d$  data matrix is oversimplification of what occurs in practice

# Transaction Data

List of store purchases: date, customer ID, list of items and prices

Web transaction log -sequence of triples: (user id, web page, time)

Can be transformed into binary-valued matrix

Individuals

1	1		1	1	1				1
1					1				
1	1	1			1	1	1	1	
1							1		
1		1		1		1			1
1			1						1

Web Page Visited

# 3.Types of Structures: Models and Patterns

- Representations sought in data mining
  - Global Model
  - Local Pattern

# Models and Patterns

- *Global Model*

- *Make a statement about any point in  $d$ -space*
  - *E.g., assign a point to a cluster*
  - *Even when some values are missing*
- *Simple model:  $Y = aX + c$* 
  - *Functional model is linear*
  - *Linear in variables rather than parameters*

- *Local Patterns*

- *Make a statement about restricted regions of space spanned by variables*
  - *E.g.1: if  $X > \text{thresh1}$  then  $\text{Prob}(Y > \text{thresh2}) = p$*
  - *E.g.2: certain classes of transactions do not show peaks and troughs (bank discovers dead people's open accounts)*

# 4. Data Mining Tasks (What?)

- Not so much a single technique
- Idea that there is more knowledge hidden in the data than shows itself on the surface
- Any technique that helps to extract more out of data is useful
- Five major task types:
  1. Exploratory Data Analysis (Visualization)
  2. Descriptive Modeling (Density estimation, Clustering)
  3. Predictive Modeling (Classification and Regression)
  4. Discovering Patterns and Rules (Association rules)
  5. Retrieval by Content (Retrieve items similar to pattern of interest)

Model  
building

# Exploratory Data Analysis

- Interactive and Visual
- Pie Charts (angles represent size)
- Cox Comb Charts (radii represent size)
- Intricate spatial displays of users of Google around the world

# Descriptive Modeling

- Describe all the data or a process for generating the data
- Probability Distribution using Density Estimation
- Clustering and Segmentation
  - Partitioning  $p$ -dimensional space into groups
  - Similar people are put in same group

# Predictive Modeling

- Classification and Regression
- Market value of a stock, disease, brittleness of a weld
- Machine Learning Approaches
- A unique variable is the objective in prediction unlike in description.



# Discovering Patterns and Rules

- Detecting fraudulent behavior by determining data that differs significantly from rest
- Finding combinations of transactions that occur frequently in transactional data bases
  - Grocery items purchased together

# Retrieval by Content

- User has pattern of interest and wishes to find that pattern in database, Ex:
  - Text Search
    - Estimate the relative importance of web pages using a feature vector whose elements are derived from the Query-URL pair
  - Image Search
    - Search a large database of images by using content descriptors such as color, texture, relative position

# Components of Data Mining Algorithms (How?)

Four basic components in each algorithm\*

## 1. Model or Pattern Structure

Determining underlying structure or functional form we seek from data

## 2. Score Function

Judging the quality of the fitted model

## 3. Optimization and Search Method

Searching over different model and pattern structures

## 4. Data Management Strategy

Handling data access efficiently

\*Illustrated in Regression example

# Statistics vs Data Mining

- Size of data set (large in data mining)
  - Eyeballing not an option (terabytes of data)
  - Entire dataset rather than a sample
- Many variables
  - Curse of dimensionality
- Make predictions
  - Small sample sizes can lead to spurious discovery:
    - Superbowl winner conference correlates to stock market (up/down)

# Searching Data Base vs Data Mining

Data Base: When you know **exactly** what you are looking for

- Query Tool: SQL (Structured Query Language) example

Table called Persons

<b>LastName</b>	<b>FirstName</b>	<b>Address</b>	<b>City</b>
Hansen	Ola	Timoteivn 10	Sandnes
Svendson	Tove	Borgvn 23	Sandnes
Pettersen	Kari	Storgt 20	Stavanger

- Query: SELECT **LastName** FROM **Persons** results in

**LastName**

Hansen
Svendson
Pettersen

Data Mining: When you only **vaguely** know what you are looking for

# Reference Textbooks

1. Hand, David, Heikki Mannila, and Padhraic Smyth, *Principles of Data Mining*, MIT Press 2001.
2. Bishop, Christopher, *Pattern Recognition and Machine Learning*, Springer 2006

## Approach:

Fundamental principles  
Emphasis on Theory and Algorithms

## Many other textbooks:

Emphasize business applications, case studies

# Many Other Textbooks

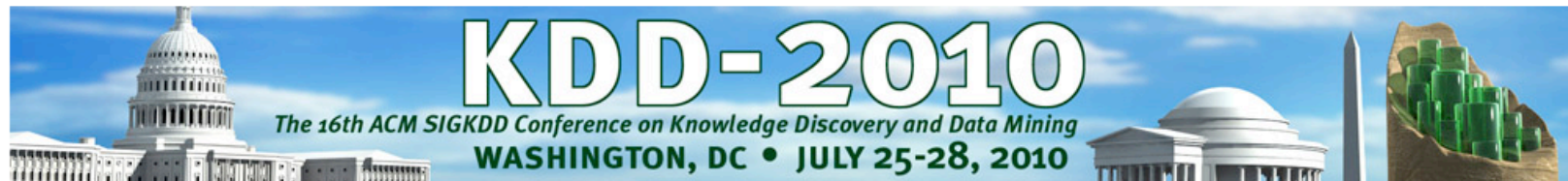
1. Han and Kamber, *Data Mining Concepts and Techniques*, Morgan Kaufmann, 2000 (Data Base Perspective)
2. Witten, I. H., and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, 2000. (Machine Learning Perspective)
3. Adriaans, P., and D. Zantinge, *Data Mining*, Addison- Wesley, 1998. (Layman Perspective)
4. Groth, R., *Data Mining: A Hands-on Approach for Business Professionals*, Prentice-Hall PTR, 1997. (Business Perspective)
5. Kennedy, R., Y. Lee, et al., *Solving Data Mining Problems through Pattern Recognition*, Prentice-Hall PTR, 1998. (Pattern Recognition Perspective)
6. Weiss, S., and N. Indurkha, *Predictive Data Mining: A Practical Guide*, Morgan Kaufmann, 1998. (Statistical Perspective)

# More Data Mining Textbooks

7. S.Chakrabarti, *Mining the web*, Morgan Kaufman, 2003 (Emphasis on webpages and hyperlinks)
8. T. Dasu and T. Johnson, *Exploratory Data Mining and Data Cleaning*, Wiley, 2003 (Focus on data quality)
9. K. Cios, W. Pedrycz and R. Swiniarski, *Data Mining Methods for Knowledge Discovery*, Kluwer, 1998, (Focus on Mathematical issues, e.g., rough sets)
10. M. Kantardzic, *Data Mining: Concepts, Models and Algorithms*, IEEE-Wiley, 2003 (Focus on Machine Learning)
11. A. K. Pujari, *Data Mining Techniques*, Universities Press, 2001, (Data Base Perspective)
12. R. Groth, *Data Mining: A hands-on approach for business professionals*, Prentice Hall, 1998 (Business user perspective including software CD)



# Premier Data Mining Conference



## Submit

- [Home](#)
- [Dates](#)
- [Call for Papers](#)
- [Call for Workshops](#)
- [Call for Tutorials](#)
- [Call for Panels](#)
- [Call for Demos](#)
- [Proposals for KDD Cup](#)

## Attend KDD

- [Program](#)
- [Registration](#)
- [Travel](#)
- [Organizers](#)
- [Awards](#)

## Important Dates

Jul 25, 2010	KDD-2010 Conference
Jun 25, 2010	Early Registration
Jun 4, 2010	Camera-ready workshop proceedings due
May 28, 2010	Camera-ready papers due
May 14, 2010	Demo, Tutorial, Panel acceptance
Apr 30, 2010	Paper acceptance
Apr 16, 2010	Panel acceptance
Feb 26, 2010	Tutorial and Panel Proposal Deadline
Mar 1, 2010	Workshop acceptance
Feb 5, 2010	Full Paper deadline
<b>Feb 2, 2010</b>	<b>Paper abstract deadline</b>
Feb 1, 2010	Workshop proposal deadline

Sign up today to become a Corporate Sponsor and/or be featured with exhibits!  
[Sponsors](#)  
[Exhibits](#)

