

Machine Learning, Data Mining, and Knowledge Discovery: An Introduction

Lesson Outline

- **Introduction: Data Flood**
- Data Mining Application Examples
- Data Mining & Knowledge Discovery
- Data Mining Tasks

2

Trends leading to Data Flood

- More data is generated:
 - Bank, telecom, other business transactions ...
 - Scientific data: astronomy, biology, etc
 - Web, text, and e-commerce



3

Big Data Examples

- Europe's Very Long Baseline Interferometry (VLBI) has 16 telescopes, each of which produces **1 Gigabit/second** of astronomical data over a 25-day observation session
 - storage and analysis a big problem
- AT&T handles billions of calls per day
 - so much data, it cannot be all stored -- analysis has to be done "on the fly", on streaming data

4

Largest databases in 2003

- Commercial databases:
 - Winter Corp. 2003 Survey: France Telecom has largest decision-support DB, ~30TB; AT&T ~ 26 TB
- Web
 - Alexa internet archive: 7 years of data, 500 TB
 - Google searches 4+ Billion pages, many hundreds TB
 - IBM WebFountain, 160 TB (2003)
 - Internet Archive (www.archive.org), ~ 300 TB

5

5 million terabytes created in 2002

- UC Berkeley 2003 estimate: 5 exabytes (5 million terabytes) of new data was created in 2002.
www.sims.berkeley.edu/research/projects/how-much-info-2003/
- US produces ~40% of new stored data worldwide

6

Data Growth Rate

- Twice as much information was created in 2002 as in 1999 (~30% growth rate)
- Other growth rate estimates even higher
- Very little data will ever be looked at by a human
- Knowledge Discovery is **NEEDED** to make sense and use of data.

7

Lesson Outline

- Introduction: Data Flood
- **Data Mining Application Examples**
- Data Mining & Knowledge Discovery
- Data Mining Tasks

8

Machine Learning / Data Mining Application areas

- Science
 - astronomy, bioinformatics, drug discovery, ...
- Business
 - advertising, CRM (Customer Relationship management), investments, manufacturing, sports/entertainment, telecom, e-Commerce, targeted marketing, health care, ...
- Web:
 - search engines, bots, ...
- Government
 - law enforcement, profiling tax cheaters, anti-terror(?)

9

Data Mining for Customer Modeling

- Customer Tasks:
 - attrition prediction
 - targeted marketing:
 - cross-sell, customer acquisition
 - credit-risk
 - fraud detection
- Industries
 - banking, telecom, retail sales, ...

10

Customer Attrition: Case Study

- Situation: Attrition rate at for mobile phone customers is around 25-30% a year!

Task:

- Given customer information for the past N months, predict who is likely to attrite next month.
- Also, estimate customer value and what is the cost-effective offer to be made to this customer.

11

Customer Attrition Results

- Verizon Wireless built a customer data warehouse
- Identified potential attriters
- Developed multiple, regional models
- Targeted customers with high propensity to accept the offer
- Reduced attrition rate from over 2%/month to under 1.5%/month (huge impact, with >30 M subscribers)

(Reported in 2003)

12

Assessing Credit Risk: Case Study

- Situation: Person applies for a loan
- Task: Should a bank approve the loan?
- Note: People who have the best credit don't need the loans, and people with worst credit are not likely to repay. Bank's best customers are in the middle

13

Credit Risk - Results

- Banks develop credit models using variety of machine learning methods.
- Mortgage and credit card proliferation are the results of being able to successfully predict if a person is likely to default on a loan
- Widely deployed in many countries

14

Successful e-commerce – Case Study

- A person buys a book (product) at Amazon.com.
- Task: Recommend other books (products) this person is likely to buy
- Amazon does clustering based on books bought:
 - customers who bought "**Advances in Knowledge Discovery and Data Mining**", also bought "**Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations**"
- Recommendation program is quite successful

15

Problems Suitable for Data-Mining

- require knowledge-based decisions
- have a changing environment
- have sub-optimal current methods
- have accessible, sufficient, and relevant data
- provides high payoff for the right decisions!

Privacy considerations important if personal data is involved

16

Lesson Outline

- Introduction: Data Flood
- Data Mining Application Examples
- **Data Mining & Knowledge Discovery**
- Data Mining Tasks

17

Knowledge Discovery Definition

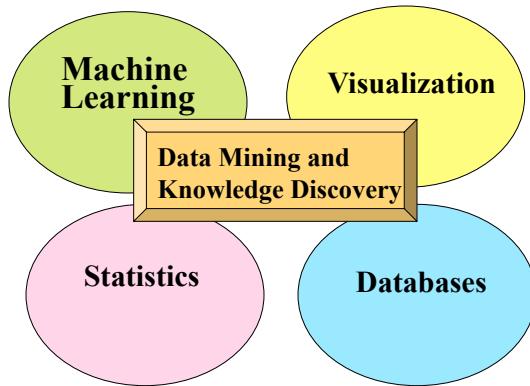
Knowledge Discovery in Data is the *non-trivial* process of identifying

- *valid*
- *novel*
- potentially *useful*
- and ultimately *understandable patterns* in data.

from *Advances in Knowledge Discovery and Data Mining*, Fayyad, Piatetsky-Shapiro, Smyth, and Uthurusamy, (Chapter 1), AAAI/MIT Press 1996

18

Related Fields



19

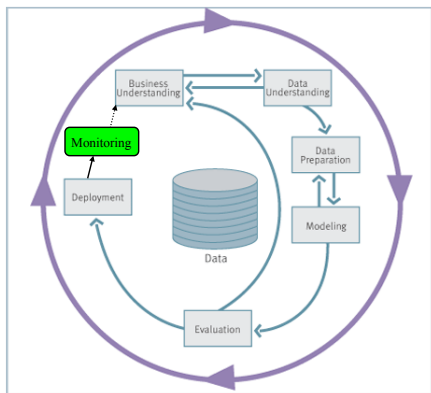
Statistics, Machine Learning and Data Mining

- Statistics:
 - more theory-based
 - more focused on testing hypotheses
- Machine learning
 - more heuristic
 - focused on improving performance of a learning agent
 - also looks at real-time learning and robotics – areas not part of data mining
- Data Mining and Knowledge Discovery
 - integrates theory and heuristics
 - focus on the entire process of knowledge discovery, including data cleaning, learning, and integration and visualization of results
- Distinctions are fuzzy

witten&keibe

20

Knowledge Discovery Process flow, according to CRISP-DM



see
www.crisp-dm.org
 for more
 information

21

Historical Note: Many Names of Data Mining

- Data Fishing, Data Dredging: 1960-
 - used by Statistician (as bad name)
- Data Mining :1990 --
 - used DB, business
 - in 2003 – bad image because of TIA
- Knowledge Discovery in Databases (1989-)
 - used by AI, Machine Learning Community
- also Data Archaeology, Information Harvesting, Information Discovery, Knowledge Extraction, ...



Currently: Data Mining and Knowledge Discovery are used interchangeably

22

Lesson Outline

- Introduction: Data Flood
- Data Mining Application Examples
- Data Mining & Knowledge Discovery
- **Data Mining Tasks**

23

Major Data Mining Tasks

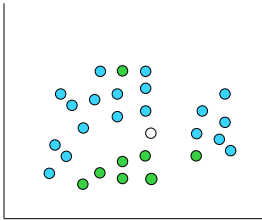
- **Classification:** predicting an item class
- **Clustering:** finding clusters in data
- **Associations:** e.g. A & B & C occur frequently
- **Visualization:** to facilitate human discovery
- **Summarization:** describing a group
- **Deviation Detection:** finding changes
- Estimation: predicting a continuous value
- Link Analysis: finding relationships

▪ ...

24

Data Mining Tasks: Classification

Learn a method for predicting the instance class from pre-labeled (classified) instances

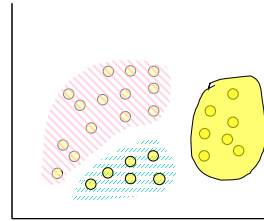


Many approaches:
Statistics,
Decision Trees,
Neural Networks,
...

25

Data Mining Tasks: Clustering

Find “natural” grouping of instances given un-labeled data



26

Summary:

- Technology trends lead to data flood
 - data mining is needed to make sense of data
- Data Mining has many applications, successful and not
- Knowledge Discovery Process
- Data Mining Tasks
 - classification, clustering, ...

27

More on Data Mining and Knowledge Discovery

KDnuggets.com

- News, Publications
- Software, Solutions
- Courses, Meetings, Education
- Publications, Websites, Datasets
- Companies, Jobs
- ...

28