



Data Mining Introduction

Dr. N. Abdolvand

Data Mining

Data Mining: Concepts and Techniques

(3rd ed.)

— Chapter 1 —

Jiawei Han, Micheline Kamber, and Jian Pei
University of Illinois at Urbana-Champaign &
Simon Fraser University

©2013 Han, Kamber & Pei. All rights reserved.

Chapter 1. Introduction

- ◆ Why Data Mining?
- ◆ What Is Data Mining?
- ◆ A Multi-Dimensional View of Data Mining
- ◆ What Kinds of Data Can Be Mined?
- ◆ What Kinds of Patterns Can Be Mined?
- ◆ What Kinds of Technologies Are Used?
- ◆ What Kinds of Applications Are Targeted?
- ◆ Major Issues in Data Mining
- ◆ A Brief History of Data Mining and Data Mining Society
- ◆ Summary

Information is crucial

- ◆ Example 1: *in vitro* fertilization
 - Given: embryos described by 60 features
 - Problem: selection of embryos that will survive
 - Data: historical records of embryos and outcome
- ◆ Example 2: cow culling
 - Given: cows described by 700 features
 - Problem: selection of cows that should be culled
 - Data: historical records and farmers' decisions

Why Data Mining?

- ◆ The Explosive Growth of Data: from terabytes to petabytes
 - Data collection and data availability
 - Automated data collection tools, database systems, Web, computerized society
 - Major sources of abundant data
 - Business: Web, e-commerce, transactions, stocks, ...
 - Science: Remote sensing, bioinformatics, scientific simulation, ...
 - Society and everyone: news, digital cameras, YouTube
- ◆ We are drowning in data, but starving for knowledge!
- ◆ “Necessity is the mother of invention”—Data mining—Automated analysis of massive data sets
- ◆ Data Opportunities:
 - Volume of Data
 - Variety of Data
 - Powerful Computers
 - Better Algorithms

5

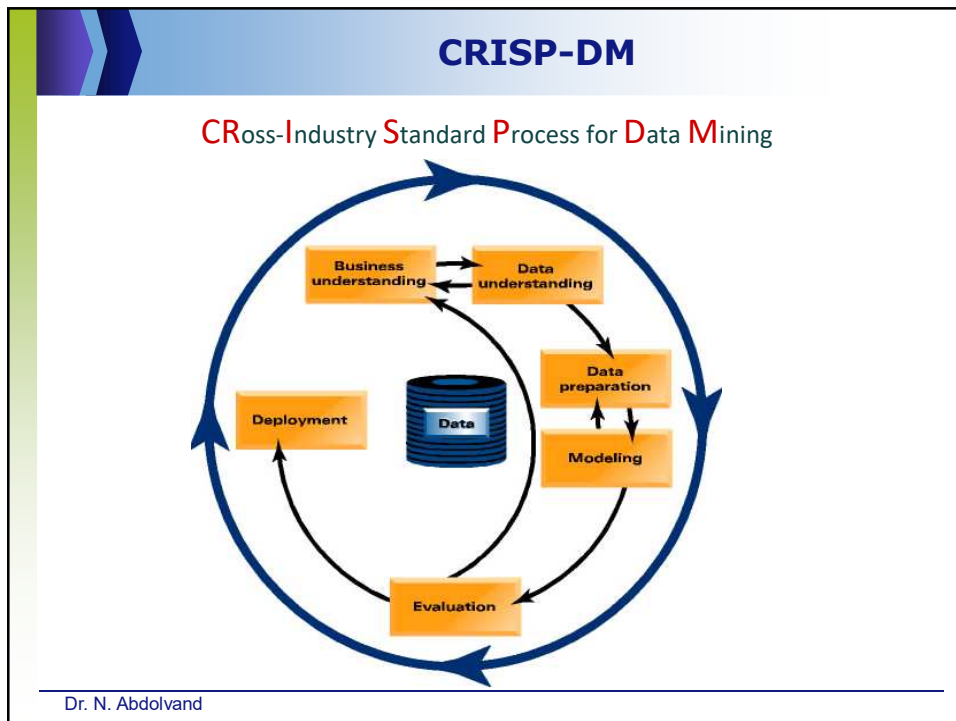
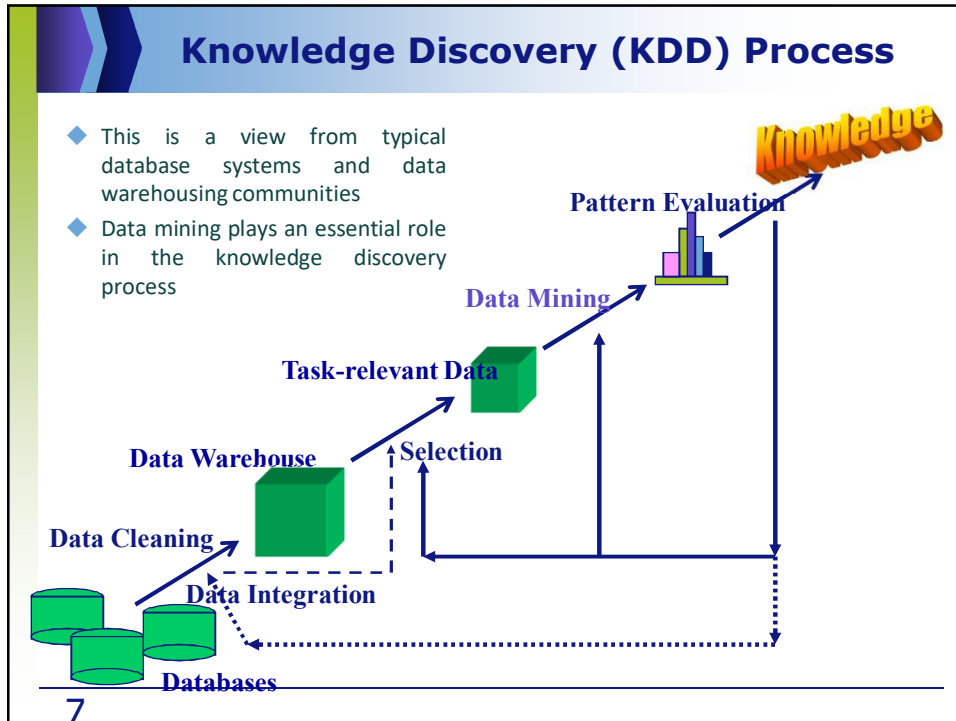
What Is Data Mining?

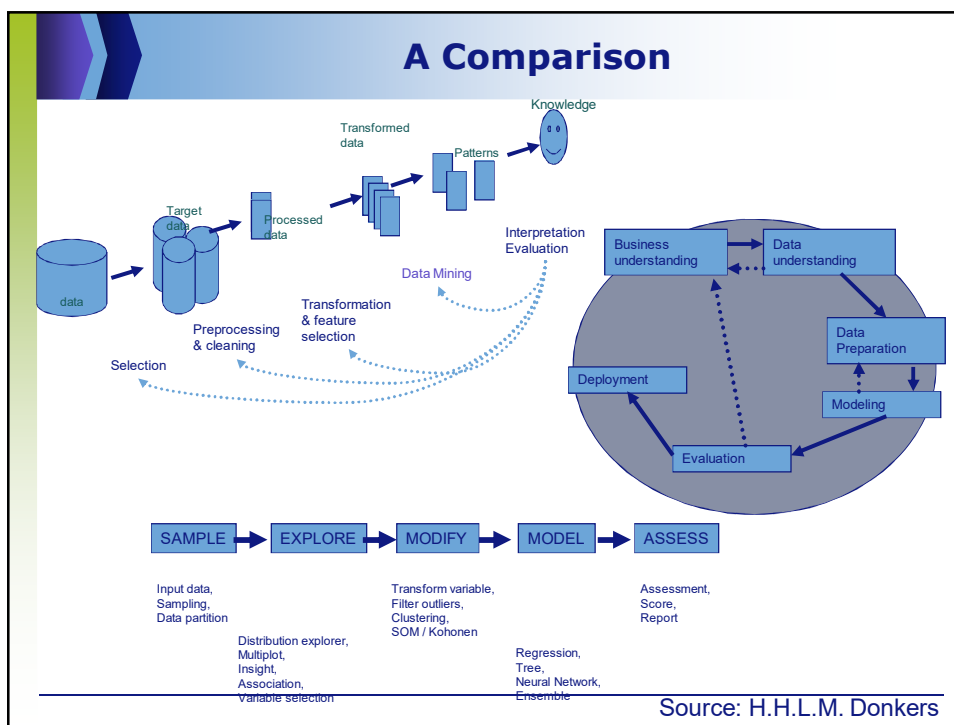
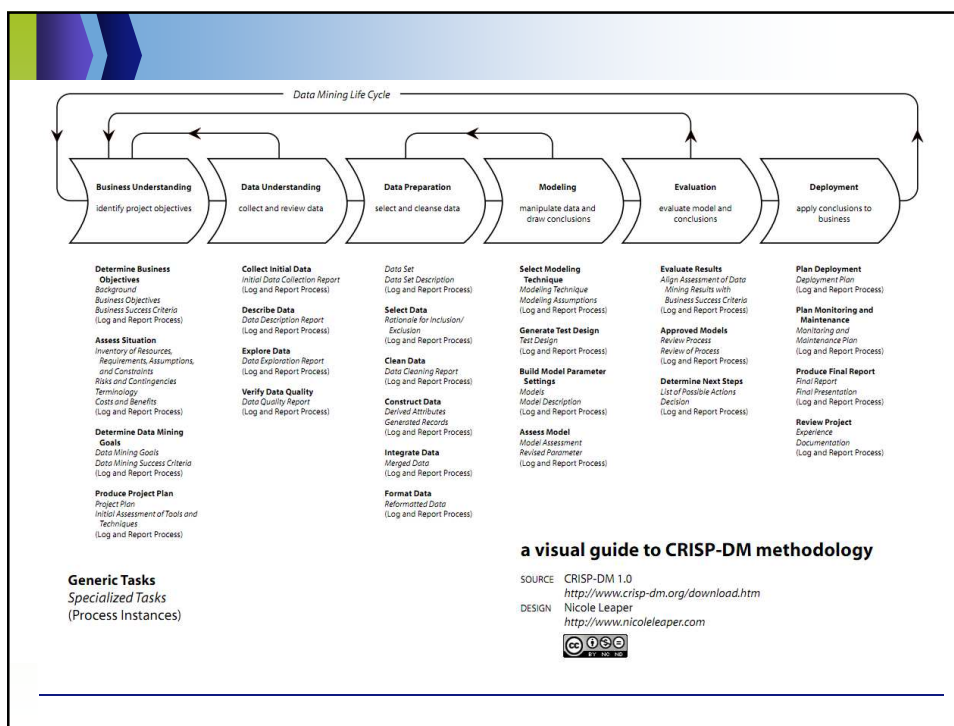


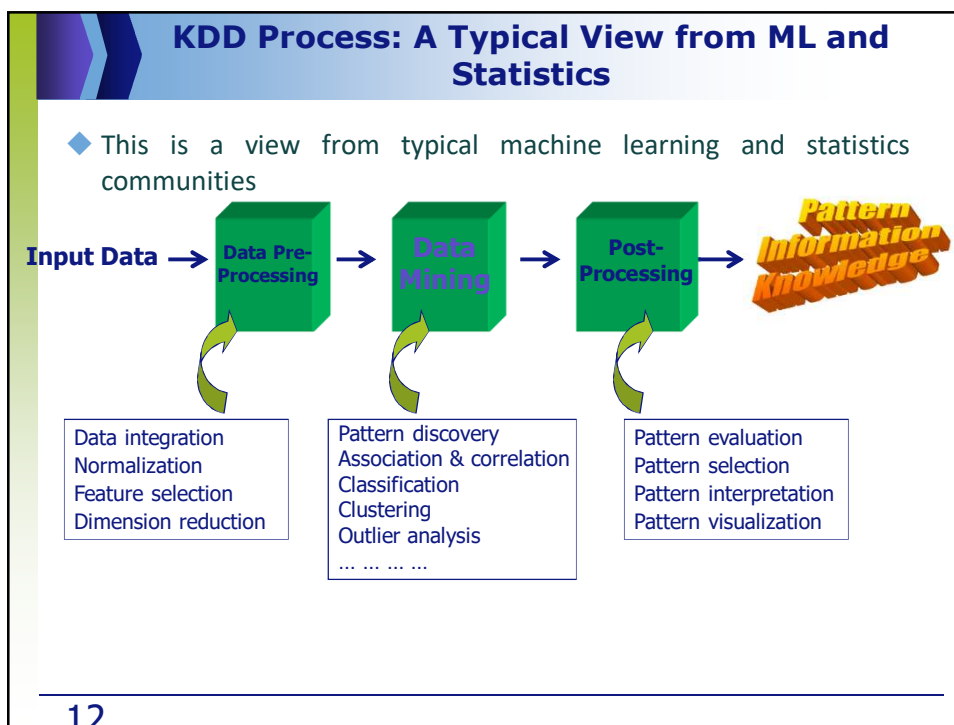
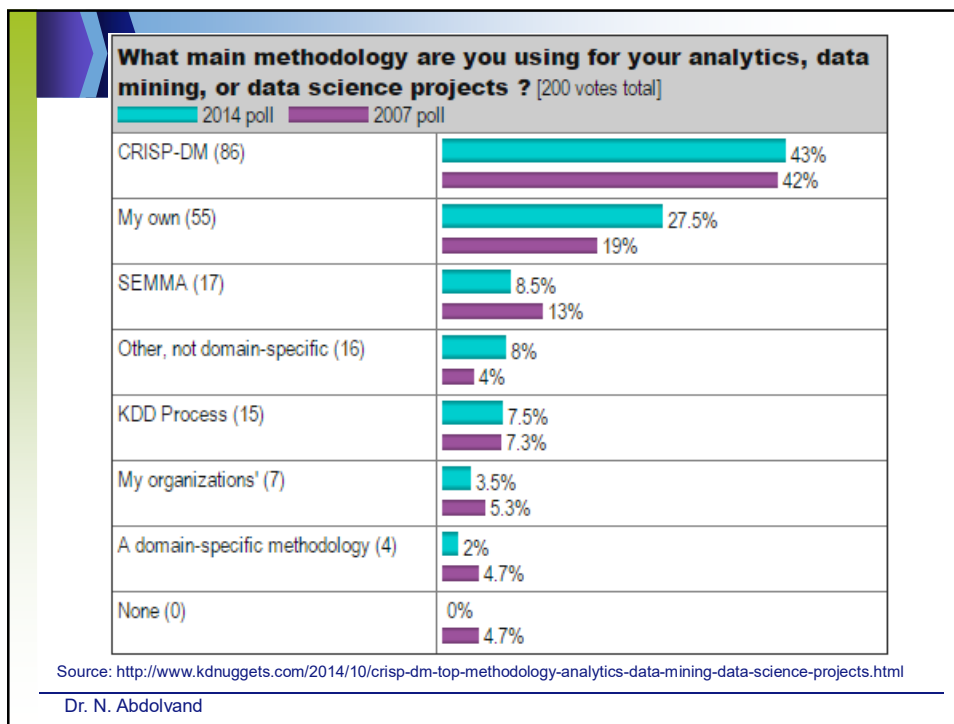
- ◆ Data mining (knowledge discovery from data)
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
 - Data mining: a misnomer?
- ◆ Alternative names
 - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- ◆ Watch out: Is everything “data mining”?
 - Simple search and query processing
 - (Deductive) expert systems



6







Multi-Dimensional View of Data Mining

- ◆ Data to be mined
 - Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks
- ◆ Knowledge to be mined (or: Data mining functions)
 - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
 - Descriptive vs. predictive data mining
 - Multiple/integrated functions and mining at multiple levels
- ◆ Techniques utilized
 - Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.
- ◆ Applications adapted
 - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

15

Data Mining: On What Kinds of Data?

- ◆ Database-oriented data sets and applications
 - Relational database, data warehouse, transactional database
 - Object-relational databases, Heterogeneous databases and legacy databases
- ◆ Advanced data sets and advanced applications
 - Data streams and sensor data
 - Time-series data, temporal data, sequence data (incl. bio-sequences)
 - Structure data, graphs, social networks and information networks
 - Spatial data and spatiotemporal data
 - Multimedia database
 - Text databases
 - The World-Wide Web

16

Data Mining Function: (1) Association and Correlation Analysis

- ◆ Frequent patterns (or frequent itemsets)
 - What items are frequently purchased together in your Walmart?
- ◆ Association, correlation vs. causality
 - A typical association rule
 - Diaper → Beer [0.5%, 75%] (support, confidence)
 - Are strongly associated items also strongly correlated?
- ◆ How to mine such patterns and rules efficiently in large datasets?
- ◆ How to use such patterns for classification, clustering, and other applications?

18

Data Mining Function: (2) Classification

- ◆ Classification and label prediction
 - Construct models (functions) based on some training examples
 - Describe and distinguish classes or concepts for future prediction
 - E.g., classify countries based on (climate), or classify cars based on (gas mileage)
 - Predict some unknown class labels
- ◆ Typical methods
 - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...
- ◆ Typical applications:
 - Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, ...

19

Data Mining Function: (3) Cluster Analysis

- ◆ Unsupervised learning (i.e., Class label is unknown)
- ◆ Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
- ◆ Principle: Maximizing intra-class similarity & minimizing interclass similarity
- ◆ Many methods and applications

20

Data Mining Function: (4) Outlier Analysis

- ◆ Outlier analysis
 - Outlier: A data object that does not comply with the general behavior of the data
 - Noise or exception? — One person's garbage could be another person's treasure
 - Methods: by product of clustering or regression analysis, ...
 - Useful in fraud detection, rare events analysis

21

Time and Ordering: Sequential Pattern, Trend and Evolution Analysis

- ◆ Sequence, trend and evolution analysis
 - Trend, time-series, and deviation analysis: e.g., regression and value prediction
 - Sequential pattern mining
 - e.g., first buy digital camera, then buy large SD memory cards
 - Periodicity analysis
 - Motifs and biological sequence analysis
 - Approximate and consecutive motifs
 - Similarity-based analysis
- ◆ Mining data streams
 - Ordered, time-varying, potentially infinite, data streams

22

Structure and Network Analysis

- ◆ Graph mining
 - Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)
- ◆ Information network analysis
 - Social networks: actors (objects, nodes) and relationships (edges)
 - e.g., author networks in CS, terrorist networks
 - Multiple heterogeneous networks
 - A person could be multiple information networks: friends, family, classmates, ...
 - Links carry a lot of semantic information: Link mining
- ◆ Web mining
 - Web is a big information network: from PageRank to Google
 - Analysis of Web information networks
 - Web community discovery, opinion mining, usage mining, ...

23

Evaluation of Knowledge

- ◆ Are all mined knowledge interesting?
 - One can mine tremendous amount of “patterns”
 - Some may fit only certain dimension space (time, location, ...)
 - Some may not be representative, may be transient, ...
- ◆ Evaluation of mined knowledge → directly mine only interesting knowledge?
 - Descriptive vs. predictive
 - Coverage
 - Typicality vs. novelty
 - Accuracy
 - Timeliness
 - ...

24

Major Issues in Data Mining (1)

- ◆ Mining Methodology
 - Mining various and new kinds of knowledge
 - Mining knowledge in multi-dimensional space
 - Data mining: An interdisciplinary effort
 - Boosting the power of discovery in a networked environment
 - Handling noise, uncertainty, and incompleteness of data
 - Pattern evaluation and pattern- or constraint-guided mining
- ◆ User Interaction
 - Interactive mining
 - Incorporation of background knowledge
 - Presentation and visualization of data mining results

28

Major Issues in Data Mining (2)

- ◆ Efficiency and Scalability
 - Efficiency and scalability of data mining algorithms
 - Parallel, distributed, stream, and incremental mining methods
- ◆ Diversity of data types
 - Handling complex types of data
 - Mining dynamic, networked, and global data repositories
- ◆ Data mining and society
 - Social impacts of data mining
 - Privacy-preserving data mining
 - Invisible data mining

29

Where to Find References? DBLP, CiteSeer, Google

- ◆ Data mining and KDD (SIGKDD: CDROM)
 - Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
 - Journal: Data Mining and Knowledge Discovery, KDD Explorations, ACM TKDD
- ◆ Database systems (SIGMOD: ACM SIGMOD Anthology—CD ROM)
 - Conferences: ACM-SIGMOD, ACM-PODS, VLDB, IEEE-ICDE, EDBT, ICDT, DASFAA
 - Journals: IEEE-TKDE, ACM-TODS/TOIS, JIIS, J. ACM, VLDB J., Info. Sys., etc.
- ◆ AI & Machine Learning
 - Conferences: Machine learning (ML), AAAI, IJCAI, COLT (Learning Theory), CVPR, NIPS, etc.
 - Journals: Machine Learning, Artificial Intelligence, Knowledge and Information Systems, IEEE-PAMI, etc.
- ◆ Web and IR
 - Conferences: SIGIR, WWW, CIKM, etc.
 - Journals: WWW: Internet and Web Information Systems,
- ◆ Statistics
 - Conferences: Joint Stat. Meeting, etc.
 - Journals: Annals of statistics, etc.
- ◆ Visualization
 - Conference proceedings: CHI, ACM-SIGGraph, etc.
 - Journals: IEEE Trans. visualization and computer graphics, etc.

30

Summary

- ◆ Data mining: Discovering interesting patterns and knowledge from massive amount of data
- ◆ A natural evolution of science and information technology, in great demand, with wide applications
- ◆ A KDD process includes data cleaning, data integration, data selection, transformation, data mining, pattern evaluation, and knowledge presentation
- ◆ Mining can be performed in a variety of data
- ◆ Data mining functionalities: characterization, discrimination, association, classification, clustering, trend and outlier analysis, etc.
- ◆ Data mining technologies and applications
- ◆ Major issues in data mining

31

Recommended Reference Books

- ◆ E. Alpaydin. *Introduction to Machine Learning*, 2nd ed., MIT Press, 2011
- ◆ S. Chakrabarti. *Mining the Web: Statistical Analysis of Hypertext and Semi-Structured Data*. Morgan Kaufmann, 2002
- ◆ R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, 2ed., Wiley-Interscience, 2000
- ◆ T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley & Sons, 2003
- ◆ U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy. *Advances in Knowledge Discovery and Data Mining*. AAAI/MIT Press, 1996
- ◆ U. Fayyad, G. Grinstein, and A. Wierse, *Information Visualization in Data Mining and Knowledge Discovery*, Morgan Kaufmann, 2001
- ◆ J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*. Morgan Kaufmann, 3rd ed., 2011
- ◆ T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., Springer, 2009
- ◆ B. Liu, *Web Data Mining*, Springer 2006
- ◆ T. M. Mitchell, *Machine Learning*, McGraw Hill, 1997
- ◆ Y. Sun and J. Han, *Mining Heterogeneous Information Networks*, Morgan & Claypool, 2012
- ◆ P.-N. Tan, M. Steinbach and V. Kumar, *Introduction to Data Mining*, Wiley, 2005
- ◆ S. M. Weiss and N. Indurkha, *Predictive Data Mining*, Morgan Kaufmann, 1998
- ◆ I. H. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, 2nd ed. 2005

32



Questions?