

Course Objectives

To provide an introduction to knowledge discovery in databases and complex data repositories, and to present basic concepts relevant to real data mining applications, as well as reveal important research issues germane to the knowledge discovery domain and advanced mining applications.



Students will understand the fundamental concepts underlying knowledge discovery in databases and gain hands-on experience with implementation of some data mining algorithms applied to real world cases.

© Dr. Osmar R. Zaïane, 1999-2004

Principles of Knowledge Discovery in Data

University of Alberta 💽 5

Evaluation and Grading

There is no final exam for this course, but there are assignments, presentations, a midterm and a project.

I will be evaluating all these activities out of 100% and give a final grade based on the evaluation of the activities.

The midterm has two parts: a <u>take-home exam</u> + <u>oral exam</u>.

- Assignments (4) 20%
- Midterm 25%
- Project
 - Quality of presentation + quality of report + quality of demos

39%

- Preliminary project demo (week 12) and final project demo (week 16) have the same weight
- Class presentations 16%
 - Quality of presentation + quality of slides + peer evaluation

• A+ will be given only for outstanding achievement.

© Dr. Osmar R. Zaïane, 1999-2004

Principles of Knowledge Discovery in Data University of Alberta

More About Evaluation

Re-examination.

None, except as per regulation.

Collaboration.

Collaborate on assignments and projects, etc; do not merely copy.

Plagiarism.

Work submitted by a student that is the work of another student or any other person is considered plagiarism. Read **Sections 26.1.4** and **26.1.5** of the University of Alberta calendar. Cases of plagiarism are immediately referred to the Dean of Science, who determines what course of action is appropriate.

© Dr. Osmar R. Zaïane, 1999-2004 Principles of F

Principles of Knowledge Discovery in Data



About Plagiarism

Plagiarism, cheating, misrepresentation of facts and participation in such offences are viewed as serious academic offences by the University and by the Campus Law Review Committee (CLRC) of General Faculties Council.

Sanctions for such offences range from a reprimand to suspension or expulsion from the University.

© Dr. Osmar R. Zaïane, 1999-2004

Notes and Textbook

Course home page:

http://www.cs.ualberta.ca/~zaiane/courses/cmput695/

We will also have a mailing list for the course (probably also a newsgroup).

Textbook:

Data Mining: Concepts and Techniques Jiawei Han and Micheline Kamber Morgan Kaufmann Publisher, 2001 ISBN 1-55860-489-8





Other Books

Principles of Data Mining

© Dr. Osmar R. Zaïane, 1999-2004

• David Hand, Heikki Mannila, Padhraic Smyth, MIT Press, 2001, ISBN 0-262-08290-X 546 pages



- Data Mining: Introductory and Advanced Topics
 - Margaret H. Dunham, Prentice Hall, 2003, ISBN 0-13-088892-3 315 pages
- Dealing with the data flood: Mining data, text and multimedia

Principles of Knowledge Discovery in Data

• Edited by Jeroen Meij, SST Publications, 2002, ISBN 90-804496-6-0 896 pages



University of Alberta 🖉

Data Mining

© Dr. Osmar R. Zaïane, 1999-2004 Principles of Knowledge Discovery in Data





© Dr. Osmar R. Zaïane, 1999-2004

Presentation Schedule Presentation Review October November 17 17 22 22 24 24 29 29 31 31 5 5 7 7 19 19 21 21 26 26 28 28 Student 1 Student 2 Student 3 4 Student 4 Student 5 Student 6 Student 7 Student 8 Student 9 Student 10 Student 1 Student 12 Student 1 Student Student 1 Student 10 Student 17 Student 18 Student Student 20 Student 2 Student 22 © Dr. Osmar R. Zaïane, 1999-2004 Principles of Knowledge Discovery in Data University of Alberta 🖉

Projects

Principles of Knowledge Discovery in Data

On-line Resources

Course notes

Course slides

Student submitted resources

• Frequently asked questions

University of Alberta

• Web links

• Glossary

• U-Chat

• Newsgroup

	Choice	Deliverables
Ô.	Implement data mining project	Project proposal + 10' proposal presentation + project pre-demo + final demo + project report

Examples and details of data mining projects will be posted on the course web site.

Assignments

- 1- Competition in one algorithm implementation
- 2- evaluation of off the shelf data mining tools
- 3- Use of educational DM tool to evaluate algorithms
- 4- Review of a paper



More About Projects

Students should write a project proposal (1 or 2 pages).



All projects are demonstrated at the end of the semester. **December 2 and 7** to the whole class.

Preliminary project demos are private demos given to the instructor on **week November 22**.

Implementations: C/C++ or Java,

OS: Linux, Window XP/2000, or other systems.



Quick Overview of some Data Mining Operations

Association Rules Clustering Classification

Principles of Knowledge Discovery in Data

What Is Association Mining?

- Association rule mining searches for relationships between items in a dataset:
 - Finding association, correlation, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
 - Rule form: "Body → Head [support, confidence]".
- Examples:
 - buys(x, "bread") \rightarrow buys(x, "milk") [0.6%, 65%]
 - major(x, "CS") ^ takes(x, "DB") → grade(x, "A") [1%, 75%]

© Dr. Osmar R. Zaïane, 1999-2004 Principles of Knowledge Discovery in Data

University of Alberta 🚰

Basic Concepts

A transaction is a set of items: $T = \{i_a, i_b, \dots, i_t\}$

 $T \subset I$, where *I* is the set of all possible items $\{i_1, i_2, \dots i_n\}$

D, the task relevant data, is a set of transactions.

An association rule is of the form: $P \rightarrow Q$, where $P \subset I$, $Q \subset I$, and $P \cap Q = \emptyset$

 $P \rightarrow Q$ holds in *D* with <u>support</u> s and $P \rightarrow Q$ has a confidence c in the transaction set *D*.

Support($P \rightarrow Q$) = Probability($P \cup Q$) Confidence($P \rightarrow Q$)=Probability(Q/P)

© Dr. Osmar R. Zaïane, 1999-2004

© Dr. Osmar R. Zaïane, 1999-2004

Principles of Knowledge Discovery in Data



University of Alberta

Association Rule Mining





Bound by a support threshold

•Frequent itemset generation is still computationally expensive

© Dr. Osmar R. Zaïane, 1999-2004

Principles of Knowledge Discovery in Data

Frequent Itemset Generation



Frequent Itemset Generation

- Brute-force approach (Basic approach):
 - Each itemset in the lattice is a candidate frequent itemset
 - Count the support of each candidate by scanning the database





Grouping Clustering Partitioning

- We need a notion of similarity or closeness (what features?)

- Should we know apriori how many clusters exist?
- How do we characterize members of groups?
- How do we label groups?

© Dr. Osmar R. Zaïane, 1999-2004





Grouping Clustering Partitioning

What about objects that belong to different groups?

- We need a notion of similarity or closeness (what features?)
- Should we know apriori how many clusters exist?
- How do we characterize members of groups?
- How do we label groups?



Classification Methods

- Decision Tree Induction
- Neural Networks
- ✤ Bayesian Classification
- ✤ K-Nearest Neighbour
- Support Vector Machines
- ✤ Associative Classifiers
- Case-Based Reasoning
- ✤ Genetic Algorithms
- Rough Set Theory
- Fuzzy Sets
- ♦ Etc.

© Dr. Osmar R. Zaïane, 1999-2004 Principles of Knowledge Discovery in Data

very in Data Univ

