







Thanks To

Without my students my research work wouldn't have been possible.

Current: (No particular order) Maria-Luiza Antonie 0069 • Jiyang Chen • Andrew Foss Seyed-Vahid Jazayeri Past: (No particular order) Stanley Oliveira Yue Zhang Yuan Ji Yang Wang Chi-Hoon Lee • Yi Li Lisheng Sun Weinan Wang Yaling Pei Jia Li Ayman Ammoura Yan Jin ٠ Hang Cui Alex Strilets Mohammad El-Hajj ٠ William Cheung • Jun Luo Maria-Luiza Antonie Jivang Chen Andrew Foss



Course Requirements

- Understand the basic concepts of database systems
- Understand the basic concepts of artificial intelligence and machine learning
- Be able to develop applications in C/C^{++} and/or Java



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.

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 is either a take-home exam or an oral exam.

- Assignments
- Midterm
 - 25% 39%
- Project
 - Quality of presentation + quality of report and proposal + quality of demos

20% (2 assignments)

- Preliminary project demo (week 11) and final project demo (week 15) have the same weight (could be week 16)
- Class presentations 16%
 - Quality of presentation + quality of slides + peer evaluation
- A+ will be given only for outstanding achievement.

Principles of Knowledge Discovery in Data



Projects

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

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

Assignments

Competition in one algorithm implementation (in C/C⁺⁺)
Devising Exercises with solutions

More About Projects

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



All projects are demonstrated at the end of the semester. **December 11-12** to the whole class.

Preliminary project demos are private demos given to the instructor on **week November 19**.

Implementations: C/C⁺⁺ or Java,

© Dr. Osmar R. Zaïane, 1999-2007

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

© Dr. Osmar R. Zaïane, 1999-2007

Principles of Knowledge Discovery in Data

University of Alberta 💽 13

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, 1999-2007



About Plagiarism

Principles of Knowledge Discovery in Data

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.

University of Alberta 🤇

Notes and Textbook

Course home page:

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

We will also have a mailing list and newsgroup for the course.

No Textbook but recommended books.

Data Mining: Concepts and Techniques Jiawei Han and Micheline Kamber Morgan Kaufmann Publisher





http://www-faculty.cs.uiuc.edu/~hanj/bk2/

ISBN 1-55860-489 550 pages

University of Alberta () 17

© Dr. Osmar R. Zaïane, 1999-2007

Principles of Knowledge Discovery in Data

800 pages



Course Schedule (Tentative, subject to changes)

There are 13 weeks from Sept 6th to December 4th.

Week 1:	Sept 6	: Introduction to Data Mining	
Week 2:	Sept 11-13	: Association Rules	
Week 3:	Sept 18-20	: Association Rules (advanced topics)	
Week 4:	Sept 25-27	: Sequential Pattern Analysis	
Week 5:	Oct 2-4	: Classification (Neural Networks)	
Week 6:	Oct 9-11	: Classification (Decision Trees and +)	
Week 7:	Oct 16-18	: Data Clustering	
Week 8:	Oct 23-25	: Outlier Detection	
Week 9:	Oct 30-Nov 1	: Data Clustering in subspaces	Due dates
Week 10:	Nov 6-8	: Contrast sets + Web Mining	-Midterm
Week 11:	Nov 13-15	: Web Mining + Class Presentations	week 8
Week 12:	Nov 20-22	: Class Presentations	-Assignment 1
Week 12:	Nov 27-29	: Class Presentations	week 6
Week 13:	Dec 4	: Class Presentations	-Assignment 2
Week 15:	Dec 11	: Project Demos	variable dates
			-
© Dr. Osm	nar R. Zaïane, 1999-2007	Principles of Knowledge Discovery in Data	University of Alberta

Course Content

- Introduction to Data Mining
- Association analysis
- Sequential Pattern Analysis
- Classification and prediction
- Contrast Sets
- Data Clustering
- Outlier Detection
- Web Mining
- Other topics if time permits (spatial data, biomedical data, etc.)









• Dealing with the data flood: Mining data, text

• Pang-Ning Tan, Michael Steinbach, Vipin Kumar Addison Wesley, ISBN: 0-321-32136-7, 769 pages

Other Books

• David Hand, Heikki Mannila, Padhraic Smyth,

• Data Mining: Introductory and Advanced Topics

Prentice Hall, 2003, ISBN 0-13-088892-3, 315 pages

SST Publications, 2002, ISBN 90-804496-6-0, 896 pages

© Dr. Osmar R. Zaïane, 1999-2007

• Principles of Data Mining

• Margaret H. Dunham,

· Edited by Jeroen Meij,

and multimedia

Principles of Knowledge Discovery in Data

University of Alberta



For those of you who watch what you eat... Here's the final word on nutrition and health. It's a relief to know the truth after all those conflicting medical studies.

- The Japanese eat very little fat and suffer fewer heart attacks than the British or Americans.
- The Mexicans eat a lot of fat and suffer fewer heart attacks than the British or Americans.
- The Japanese drink very little red wine and suffer fewer heart attacks than the British or Americans
- The Italians drink excessive amounts of red wine and suffer fewer heart attacks than the British or Americans.
- The Germans drink a lot of beer and eat lots of sausages and fats and suffer fewer heart attacks than the British or Americans.

CONCLUSION:

Eat and drink what you like. Speaking English is apparently what kills you.

© Dr. Osmar R. Zaïane, 1999-2007

Principles of Knowledge Discovery in Data

University of Alberta

Quick Overview of some Data Mining Operations

Association Rules Clustering Classification Outlier Detection

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%]

University of Alberta

Basic Concepts



Grouping

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?

ta University of Alberta 🔛

Grouping

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?

© Dr. Osmar R. Zaïane, 1999-2007

Principles of Knowledge Discovery in Data University of Alberta

What is Classification?

The goal of data classification is to organize and categorize data in distinct classes.

- A model is first created based on the data distribution.
- ▶ The model is then used to classify new data.
- Given the model, a class can be predicted for new data.

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.

Outlier Detection

- To find exceptional data in various datasets and uncover the implicit patterns of rare cases
- Inherent variability reflects the natural variation
- Measurement error (inaccuracy and mistakes)
- Long been studied in statistics
- An active area in data mining in the last decade
- Many applications
 - Detecting credit card fraud
 - Discovering criminal activities in E-commerce
 - Identifying network intrusion
 - Monitoring video surveillance

© Dr. Osmar R. Zaïane, 1999-2007

- ...

Outliers are Everywhere

- Data values that appear inconsistent with the rest of the data.
- Some types of outliers

ARE SIGNED WHEN YOUR EMPLOYEES FOUND A YOU'RE OUT SICK. TELECOMMUTE . CORRELATION WE HAVE THEY DO? DOGBERT CONSULTS IT SAYS YOU'VE THEN IT SAYS. BEEN STEALING "HA HA, THAT WASN'T PUDDING!" MY DATA-MINING LUNCHES FROM THE SOFTWARE HAS REFRIGERATOR IN FOUND ANOTHER THE BREAK MESSAGE ROOM. FROM GOD DOGBERT CONSULTS IF YOU MINE THE DATA HARD ENOUGH, YOU CAN ALSO FIND ... SALES TO LEFT-HANDED SQUIRRELS ARE UP... AND GOD YOU NEED TO DO DATA MINING SAYS YOUR TIE TO UNCOVER MESSAGES FROM DOESN'T GO WITH HIDDEN SALES THAT SHIRT. GOD. TRENDS. **Contrasting Sequence Sets** Collaborative e-learning **E**-commerce sites What makes a buyer buy & a non-buyer leave empty handed? What contrasts genes or proteins? Finding emerging sequences and contrasting the sets of sequences would give insight about what behaviour leads to success and what doesn't pata University of Alberta

AND 100% OF ALL

EXPENSE VOUCHERS

WHEN YOU'RE ON VACATION, ALL

THE DATA MINER

EUREKAL I

Breast Cancer

Recommender Systems

WebViz

Clustering with Constraints

