

# • ISE 365/465: Applied Data Mining

## Syllabus: Spring Semester 2020

### Course Description

This course is an introduction to Data Mining. Due to the explosion of data being collected and the advancement of computer technology, the field of Data Mining has evolved to find useful patterns of information in these large collections of data. The course focuses primarily on the data, modeling, and mathematical techniques used for the application of Data Mining. Students will gain a conceptual understanding of algorithms used for Data Mining and hands-on experience with these algorithms using IBM SPSS Modeler.

### Course Objectives

Upon completion of this course, students will:

- Understand the CRISP-DM and SEMMA process for applying data mining
- Understand the importance of *data* in the data mining process and how to collect and set up data for a successful data mining project
- Have a conceptual understanding of several common data mining algorithms drawn from the fields of statistics and machine learning. This understanding will allow the students to apply appropriate algorithms correctly to solve data mining problems
- Be able to use IBM SPSS Modeler to collect, aggregate, explore and manipulate data from multiple sources to facilitate modeling
- Be able to use IBM SPSS Modeler to model and solve data mining problems using data mining algorithms
- Be able to use IBM SPSS Modeler to evaluate and present the results from a data mining project
- Understand the wide array of fields where data mining is applied today
- Have basic skills needed to be hired for a job in the field of data mining

### Prerequisites

**A basic statistics course** covering topics such as statistical significance and linear regression - While small samples are not a problem in data mining, basic statistical techniques will be important when using many data mining techniques. Students are expected to have basic statistical knowledge before taking this course. For advanced data mining techniques, theoretical background will be taught to allow their correct application to data mining problems.

## Contact Information

### Professor

Brent L. Peterson

E-mail: blp219@lehigh.edu

Office: Mohler Lab #350

Office Hours: 9:15 – 11:00, Tuesdays and Thursdays

## Textbook (required)

**Data Mining: Concepts and Techniques, 3rd ed.** By Han, Jiawei, Kamber, Michelle, and Pei, Jian. Elsevier, 2012. ISBN 978-0-12-381479-1.

Coverage will be in chapters 1-4, chapter 6, chapters 8-10, and chapter 13 for this course.

## Course Philosophy

The goal of this course is to impart knowledge through traditional lectures, in-class computer labs, readings, homework, exams, and a project. In data mining, the best way to learn is through hands-on learning. Therefore, this course will strive to give sufficient background through readings and lectures to allow the students to spend as much time as possible using IBM SPSS Modeler to develop solutions to data mining problems from a number of different problem domains. **While the course outline below is organized by solution technique, we will cover many topics throughout the course including people issues, data collection and set-up, modeling tips and tricks, and deployment that are critical to the data mining process. Students will have the opportunity to practice what they learn in labs, homework, and the project. These assignments will be heavy in hands-on use of IBM SPSS Modeler. In addition, knowledge of topics covered will be tested in written exams.**

## Academic Honesty

Integrity and Honesty are vital in life, especially for data miners, since we may have access to sensitive data and the systems we design or modify can improve people's quality of life, or can do irreparable harm. Using data mining ethically requires that we state all of the facts and assumptions in as clear a manner as possible, to avoid "lying with statistics". We are also bound by honor to give credit where it is due. In this class, you might ask others for help with a homework assignment. Once you write up your answer in your own words to turn in, it is a good idea to include a mention of their help on any particular problem. It is dishonest to copy homework solutions from past years that you might obtain or have. On quizzes and exams, of course, your work should be entirely your own. **Violations of academic honesty will result in a grade of 0 on the assignment where dishonest behavior occurred and possible disciplinary proceedings.**

Here is a statement of the Lehigh Student Senate on academic integrity: We, the Lehigh University Student Senate, as the standing representative body of all undergraduates, reaffirm the duty and obligation of students to meet and uphold the highest principles and values of personal, moral and ethical conduct. As partners in our educational community, both students and faculty share the responsibility for promoting and helping ensure an environment of academic integrity. As such, each student is expected to complete all academic course work in accordance to the standards set forth by the faculty and in compliance with the University's Code of Conduct.

## **Grading Policy**

Your final numeric score will be determined as follows:

40% : Homework  
5% : Class participation  
15% : Midterm Exam  
20% : Project  
20% : Final Exam

**Late Homework Policy: Late homework will incur an immediate 20% penalty if not handed in on time and an additional 1% for each hour late after the due time.**

**Missed Exams: Any exam missed without a legitimate excuse will result in a score of 0 on that exam.**

Plus and minus grading will be used for final grades. Final grades will be "curved".

**Accommodations for Students with Disabilities:** Lehigh University is committed to maintaining an equitable and inclusive community and welcomes students with disabilities into all of the University's educational programs. In order to receive consideration for reasonable accommodations, a student with a disability must contact Disability Support Services (DSS), provide documentation, and participate in an interactive review process. If the documentation supports a request for reasonable accommodations, DSS will provide students with a Letter of Accommodations. Students who are approved for accommodations at Lehigh should share this letter and discuss their accommodations and learning needs with instructors as early in the semester as possible. For more information or to request services, please contact Disability Support Services in person in Williams Hall, Suite 301, via phone at 610-758-4152, via email at [indss@lehigh.edu](mailto:indss@lehigh.edu), or online at <https://studentaffairs.lehigh.edu/disabilities>.

## **Principles of Equitable Community:**

Lehigh University endorses The Principles of Our Equitable Community [[http://www.lehigh.edu/~inprv/initiatives/PrinciplesEquity\\_Sheet\\_v2\\_032212.pdf](http://www.lehigh.edu/~inprv/initiatives/PrinciplesEquity_Sheet_v2_032212.pdf)]. We expect each member of this class to acknowledge and practice these Principles. Respect for each other and for differing viewpoints is a vital component of the learning environment inside and outside the classroom.

**Tentative Course Outline – This will almost certainly change. Official due dates and topics will be announced in class by the instructor throughout the semester. This outline is for your reference only to give an idea of topics that may be covered.**

<b>Tuesday Date</b>	<b>Tuesday</b>	<b>Thursday</b>
Jan 21	<ul style="list-style-type: none"> <li>• Course Introduction</li> <li>• Pre-course survey</li> <li>• Chapter 1 – Intro to Data Mining</li> <li>• <b>Assign Chapters 1, 4.1-2 Reading</b></li> </ul>	<ul style="list-style-type: none"> <li>• Chapter 1, Ch. 4.1 and 4.2 – Intro to Data Mining and Data Warehousing</li> <li>• US Presidential Election Lecture</li> <li>• <b>Assign Chapter 2 &amp;3 Reading</b></li> </ul>
Jan 28	<ul style="list-style-type: none"> <li>• Chapter 2 – Getting to Know Your Data</li> </ul>	<ul style="list-style-type: none"> <li>• IBM SPSS Modeler Graphs and Exploratory Statistics</li> <li>• <b>Assign Data Mining Paper and “10 Data Mining Mistakes” Reading</b></li> </ul>
Feb 4	<ul style="list-style-type: none"> <li>• Chapter 3– Data Preprocessing</li> <li>• Feature Selection</li> <li>• IBM SPSS Modeler Data Manipulation</li> <li>• <b>Assign “Getting Started with SPSS Modeler” Reading: <a href="#">IBM Introduction.html</a></b></li> </ul>	<ul style="list-style-type: none"> <li>• Merge Overview Discussion</li> <li>• <b>Data Understanding and Preprocessing Lab</b></li> </ul>
Feb 11	<ul style="list-style-type: none"> <li>• <b>Data Understanding and Preprocessing Lab</b></li> <li>• <u>Assign Homework 1</u></li> </ul>	<ul style="list-style-type: none"> <li>• Applied Linear Regression Lecture</li> </ul>
Feb 18	<ul style="list-style-type: none"> <li>• Pop Quiz Results</li> <li>• Top 10 Data Mining Mistakes</li> <li>• Linear Regression Lab</li> </ul>	<ul style="list-style-type: none"> <li>• Linear Regression Lab Continued</li> <li>• <b>Assign Chapter 8.1, 8.2, 8.5, and 8.6 Reading</b></li> <li>• <u>Assign Homework 2</u></li> <li>• <b>Homework 1 Due</b></li> </ul>
Feb 25	<ul style="list-style-type: none"> <li>• Principal Components Analysis</li> <li>• Model Evaluation – <b>Chapter 8.5</b></li> <li>• <b>Chapter 8.1 and 8.2</b> – Decision Trees</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Chapter 8.1 and 8.2</b> – Decision Trees</li> <li>• <b>Homework 2 due</b></li> </ul>
Mar 3	<ul style="list-style-type: none"> <li>• <b>Decision Tree Lab</b></li> <li>• Mid-class Survey and <u>Exam Review</u></li> </ul>	<ul style="list-style-type: none"> <li>• <b>MID-TERM EXAM</b></li> </ul>
Mar 10	<ul style="list-style-type: none"> <li>• <b>SPRING BREAK</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>SPRING BREAK</b></li> </ul>
Mar 17	<ul style="list-style-type: none"> <li>• Mid-term Exam results</li> <li>• <b>Chapter 8.6</b> - Bagging, Boosting, Ensemble Models</li> <li>• <b>Project Description</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Chapter 9.3</b> - Neural Networks</li> <li>• <b>Chapter 9.4</b> - Support Vector Machines</li> <li>• <u>Assign Homework 3</u></li> </ul>
Mar 24	<ul style="list-style-type: none"> <li>• <b>Neural Network/Support Vector Machines/Bagging and Boosting Lab</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Chapter 6</b> - Market Basket Analysis &amp; Association Rules - Apriori</li> </ul>
Mar 31	<ul style="list-style-type: none"> <li>• <b>Chapter 6</b> - Market Basket Analysis &amp; Association Rules - Apriori</li> <li>• <b>MBA/Apriori Lab</b></li> <li>• <b>Homework 3 Due</b></li> <li>• <u>Assign Homework 4</u></li> </ul>	<ul style="list-style-type: none"> <li>• <b>MBA/Apriori Lab</b></li> <li>• <b>Chapter 10</b> - Cluster Analysis</li> </ul>
April 7	<ul style="list-style-type: none"> <li>• <b>Chapter 10</b> - Cluster Analysis</li> <li>• <b>Cluster Analysis Lab</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Cluster Analysis Lab</b></li> <li>• <b>Assign Ch. 9.5.1 Reading – k-NN</b></li> <li>• <b>Homework 4 Due</b></li> </ul>
April 14	<ul style="list-style-type: none"> <li>• Project Update</li> <li>• k-Nearest Neighbors</li> </ul>	<ul style="list-style-type: none"> <li>• <b>k-Nearest Neighbor Lab</b></li> </ul>
April 21	<ul style="list-style-type: none"> <li>• <b>Chapter 13</b> - Advanced Data Mining Topics</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Chapter 13</b> - Advanced Data Mining Topics</li> <li>• Final Exam Review</li> </ul>
April 28	<ul style="list-style-type: none"> <li>• <b>PROJECT PRESENTATIONS</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>PROJECT PRESENTATIONS</b></li> </ul>

	<ul style="list-style-type: none"><li>• <b>Project Write-up Due</b></li></ul>	
--	---	--