



ORIGINAL COURSE IMPLEMENTATION DATE: January 2012  
 REVISED COURSE IMPLEMENTATION DATE: September 2021  
 COURSE TO BE REVIEWED (six years after UEC approval): March 2025  
 Course outline form version: 05/18/2018

## OFFICIAL UNDERGRADUATE COURSE OUTLINE FORM

Note: The University reserves the right to amend course outlines as needed without notice.

<b>Course Code and Number:</b> STAT 431	<b>Number of Credits:</b> 3 <a href="#">Course credit policy (105)</a>														
<b>Course Full Title:</b> Data Mining <b>Course Short Title:</b> <i>(Transcripts only display 30 characters. Departments may recommend a short title if one is needed. If left blank, one will be assigned.)</i>															
<b>Faculty:</b> Faculty of Science	<b>Department (or program if no department):</b> Mathematics & Statistics														
<b>Calendar Description:</b> Data mining provides the techniques of extracting useful information and hidden patterns from a massive amount of data. Main topics include data exploration, classification, decision trees, Bayesian classifiers, frequent item sets, association rules, clustering, K-means, EM algorithm, and anomaly detection.  Note: This course is offered as STAT 431 and COMP 431. Students may take only one of these for credit.															
<b>Prerequisites (or NONE):</b>	COMP 230 (formerly CIS 230), STAT 271, and STAT 331/COMP 331.														
<b>Corequisites (if applicable, or NONE):</b>															
<b>Pre/corequisites (if applicable, or NONE):</b>															
<b>Antirequisite Courses</b> <i>(Cannot be taken for additional credit.)</i> Former course code/number: <b>MATH 431</b> Cross-listed with: <b>COMP 431</b> Dual-listed with: Equivalent course(s): <i>(If offered in the previous five years, antirequisite course(s) will be included in the calendar description as a note that students with credit for the antirequisite course(s) cannot take this course for further credit.)</i>	<b>Special Topics</b> <i>(Double-click on boxes to select.)</i> This course is offered with different topics: <input checked="" type="checkbox"/> No <input type="checkbox"/> Yes <i>(If yes, topic will be recorded when offered.)</i>														
<b>Typical Structure of Instructional Hours</b> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>Lecture/seminar hours</td><td style="text-align: center;">25</td></tr> <tr><td>Tutorials/workshops</td><td></td></tr> <tr><td>Supervised laboratory hours</td><td style="text-align: center;">25</td></tr> <tr><td>Experiential (field experience, practicum, internship, etc.)</td><td></td></tr> <tr><td>Supervised online activities</td><td></td></tr> <tr><td>Other contact hours:</td><td></td></tr> <tr><td style="text-align: right;"><b>Total hours</b></td><td style="text-align: center;"><b>50</b></td></tr> </table> Labs to be scheduled independent of lecture hours: <input checked="" type="checkbox"/> No <input type="checkbox"/> Yes	Lecture/seminar hours	25	Tutorials/workshops		Supervised laboratory hours	25	Experiential (field experience, practicum, internship, etc.)		Supervised online activities		Other contact hours:		<b>Total hours</b>	<b>50</b>	<b>Independent Study</b> If offered as an Independent Study course, this course may be repeated for further credit: <i>(If yes, topic will be recorded.)</i> <input checked="" type="checkbox"/> No <input type="checkbox"/> Yes, repeat(s) <input type="checkbox"/> Yes, no limit
	Lecture/seminar hours	25													
	Tutorials/workshops														
	Supervised laboratory hours	25													
Experiential (field experience, practicum, internship, etc.)															
Supervised online activities															
Other contact hours:															
<b>Total hours</b>	<b>50</b>														
<b>Transfer Credit</b> Transfer credit already exists: <i>(See <a href="http://bctransferguide.ca">bctransferguide.ca</a>.)</i> <input checked="" type="checkbox"/> No <input type="checkbox"/> Yes Submit outline for (re)articulation: <input checked="" type="checkbox"/> No <input type="checkbox"/> Yes <i>(If yes, fill in transfer credit form.)</i>	<b>Grading System</b> <input checked="" type="checkbox"/> Letter Grades <input type="checkbox"/> Credit/No Credit														
	<b>Maximum enrolment (for information only):</b> 36 <b>Expected Frequency of Course Offerings:</b> Annually <i>(Every semester, Fall only, annually, etc.)</i>														
<b>Department / Program Head or Director:</b> Ian Affleck	<b>Date approved:</b> June 15, 2020														
<b>Faculty Council approval</b>	<b>Date approved:</b> September 11, 2020														
<b>Dean/Associate VP:</b> Lucy Lee	<b>Date approved:</b> September 11, 2020														
<b>Campus-Wide Consultation (CWC)</b>	<b>Date of posting:</b> n/a														
<b>Undergraduate Education Committee (UEC) approval</b>	<b>Date of meeting:</b> January 29, 2021														

**Learning Outcomes:**

Upon successful completion of this course, students will be able to:

1. Use summary statistics to describe univariate and multivariate data sets.
2. Apply different definitions of distances to measure the similarity and dissimilarity between data objects.
3. Classify discrete and continuous data using decision trees.
4. Employ different classifiers (rule-based, nearest-neighbor, naïve Bayes) to solve various classification problems.
5. Generate various frequent itemsets and set up association rules in market basket analysis.
6. Evaluate association patterns by interest factor.
7. Use the K-means clustering technique to classify data in different clusters.
8. Determine the correct number of clusters.
9. Apply the EM algorithm to cluster data by estimating the parameters of a mixture model.
10. Detect outliers in univariate and multivariate normal distributions.
11. Collaborate with peers to complete and present a project which involves a large real-life data set and requires the skills and abilities above.
12. Integrate feedback and suggestions from peers, faculty, and supervisors in completion and presentation of final project findings.

**Prior Learning Assessment and Recognition (PLAR)**

Yes       No, PLAR cannot be awarded for this course because

**Typical Instructional Methods** (*Guest lecturers, presentations, online instruction, field trips, etc.; may vary at department's discretion.*)  
Lectures, and use of computer software. All classes take place in a computer lab.

**NOTE: The following sections may vary by instructor. Please see course syllabus available from the instructor.**

**Typical Text(s) and Resource Materials** (*If more space is required, download Supplemental Texts and Resource Materials form.*)

Author (surname, initials)	Title (article, book, journal, etc.)	Current ed.	Publisher	Year
1. P. Tan et al.	Introduction to Data Mining	<input checked="" type="checkbox"/>	Pearson	2018
2. M. Kantardzic	Data Mining---Concepts, Models, Methods and Algorithms (2 <sup>nd</sup> edition)	<input checked="" type="checkbox"/>	Wiley-IEEE Press	2011
3.		<input type="checkbox"/>		
4.		<input type="checkbox"/>		
5.		<input type="checkbox"/>		

**Required Additional Supplies and Materials** (*Software, hardware, tools, specialized clothing, etc.*)

Statistical software such as SAS Enterprise Miner.

**Typical Evaluation Methods and Weighting**

Final exam:	40%	Assignments:	15%	Field experience:	%	Portfolio:	%
Midterm exam:	20%	Project:	25%	Practicum:	%	Other:	%
Quizzes/tests:	%	Lab work:	%	Shop work:	%	Total:	100%

**Details (if necessary):**

The above percentages may vary among instructors and years, but the project component will constitute at least 10% of the overall grade. The final exam is comprehensive. Students must obtain at least 40% on the final exam to pass the course.

**Typical Course Content and Topics**

**Introduction:** The origins of data mining, data mining tasks.

**Exploring data:** Types of data, data quality, data preprocessing, measures of similarity and dissimilarity, Minkowski distance, summary statistics, visualization, OLAP and multi-dimensional data analysis.

**Classification:** Decision tree induction, model over-fitting, evaluating the performance of a classifier, methods for comparing classifiers, rule-based classifiers, nearest-neighbor classifiers, Bayesian classifiers, artificial neural network, support vector machine, ensemble methods.

**Association analysis:** Frequent itemset generation, rule generation, maximal frequent itemsets, closed frequent itemsets, FP-Growth algorithm, evaluation of association patterns, interest factor, handling categorical and continuous attributes, sequential patterns.

**Cluster analysis:** Different types of clusters, K-means, bisecting K-means, agglomerative hierarchical clustering, density-based clustering algorithm, determining the correct number of clusters, supervised measures of cluster validity, EM algorithm, minimum spanning tree clustering.

**Anomaly detection:** Causes of anomalies, determining outliers in univariate and multivariate normal distributions, proximity-based outlier detection, density-based outlier detection, clustering-based techniques.