

# CSC 411

## MACHINE LEARNING and DATA MINING

Lectures:	Monday, Wednesday 12-1 (section 1), 3-4 (section 2)
Lecture Room:	MP 134 (section 1); Bahen 1200 (section 2)
Instructor (section 1):	Richard Zemel
Instructor (section 2):	Raquel Urtasun
Instructor email:	<csc411prof@cs.toronto.edu>
Office hours:	Tuesday 3-4 Pratt 290E (Urtasun); Thursday 3-4 Pratt 290D (Zemel)
TA email:	<csc411ta@cs.toronto.edu>
Tutorials:	Fridays 12-1 (section 1); 3-4 (section 2)
Tutorial Room:	Same as lecture
Class URL:	<a href="http://www.cs.toronto.edu/~zemel/Courses/CS411.html">www.cs.toronto.edu/~zemel/Courses/CS411.html</a>

## Overview

Machine learning research aims to build computer systems that learn from experience. Learning systems are not directly programmed by a person to solve a problem, but instead they develop their own program based on examples of how they should behave, or from trial-and-error experience trying to solve the problem. These systems require learning algorithms that specify how the system should change its behavior as a result of experience. Researchers in machine learning develop new algorithms, and try to understand which algorithms should be applied in which circumstances.

Machine learning is an exciting interdisciplinary field, with historical roots in computer science, statistics, pattern recognition, and even neuroscience and physics. In the past 10 years, many of these approaches have converged and led to rapid theoretical advances and real-world applications.

This course will focus on the machine learning methods that have proven valuable and successful in practical applications. This course will contrast the various methods, with the aim of explaining the circumstances under which each is most appropriate. We will also discuss basic issues that confront any machine learning method.

## Pre-requisites

You should understand basic probability and statistics, (STA 107, 250), and college-level algebra and calculus. For example it is expected that you know about standard probability distributions (Gaussians, Poisson), and also how to calculate derivatives. Knowledge

of linear algebra is also expected, and knowledge of mathematics underlying probability models (STA 255, 261) will be useful. For the programming assignments, you should have some background in programming (CSC 270), and it would be helpful if you know Matlab or Python. Some introductory material for Matlab will be available on the course website as well as in the first tutorial.

## Readings

There is no required textbook for this course. There are several recommended books. On the course webpage we will readings from *Introduction to Machine Learning* by Ethem Alpaydin, and from *Pattern Recognition and Machine Learning* by Chris Bishop. We also provide pointers to other online resources.

## Course requirements and grading

The format of the class will be lecture, with some discussion. I strongly encourage interaction and questions. There are assigned readings for each lecture that are intended to prepare you to participate in the class discussion for that day.

The grading in the class will be divided up as follows:

Assignments	50%
Mid-Term Exam	20%
Final Exam	30%

There will be four assignments; each is worth 12.5% of your grade.

## Homework assignments

The best way to learn about a machine learning method is to program it yourself and experiment with it. So the assignments will generally involve implementing machine learning algorithms, and experimentation to test your algorithms on some data. You will be asked to summarize your work, and analyze the results, in brief (3-4 page) write ups. The implementations may be done in any language, but Matlab or Python is recommended. A brief tutorial on Matlab is available from the course web-site.

Collaboration on the assignments is not allowed. Each student is responsible for his or her own work. Discussion of assignments and programs should be limited to clarification of the handout itself, and should not involve any sharing of pseudocode or code or simulation results. Violation of this policy is grounds for a semester grade of F, in accordance with university regulations.

The schedule of assignments is included in the syllabus. Assignments are due at the beginning of class/tutorial on the due date. Because they may be discussed in class that day, it is important that you have completed them by that day. Assignments handed in late but before 5 pm of that day will be penalized by 5% (i.e., total points multiplied by 0.95); a late penalty of 10% per day will be assessed thereafter. Extensions will be granted only in special situations, and you will need a Student Medical Certificate or a written request approved by the instructor at least one week before the due date.

For the final assignment, we will have a *bake-off*: a competition between machine learning algorithms. We will give everyone some data for training a machine learning system, and you will try to develop the best method. We will then determine which system performs best on some unseen test data.

## **Exams**

There will be a mid-term in tutorial on October 24<sup>th</sup>, which will be a closed book exam on all material covered up to that point in the lectures, tutorials, required readings, and assignments.

The final will not be cumulative, except insofar as concepts from the first half of the semester are essential for understanding the later material.

The exams will cover material presented in lectures, tutorials, and assignments. You will not be responsible for topics in the reading not covered in any of these.

## **Attendance**

We expect students to attend all classes, and all tutorials. This is especially important because we will cover material in class that is not included in the textbook. Also, the tutorials will not only be for review and answering questions, but new material will also be covered.

## **Electronic Communication**

If you have questions about the assignments, you should send email to the TA account, and cc me on it. You should include your full name in the email, and it will also be useful to include your CDF account name and/or student number. Feel free to email me with questions or comments about the material covered in the course, or other related questions.

For questions about marks on the assignments, please first contact the TA. Questions about the exams should be addressed to me.

## CLASS SCHEDULE, Part 1

Shown below are the topics for lectures and tutorials (in italics), as are the dates that each assignment will be handed out and is due. The notes from each lecture and tutorial will be available on the class web-site the day of the class meeting. The assigned readings are specific sections from the book. All of these are subject to change.

<b>Date</b>	<b>Topic</b>	<b>Assignments</b>
Sep 8	Introduction	
Sep 10	Linear Regression	
<i>Sep 12</i>	<i>Probability for ML &amp; Linear regression</i>	
Sep 15	Linear Discrimination	
Sep 17	Logistic Regression	
<i>Sep 19</i>	<i>Optimization for ML</i>	
Sep 22	Decision Trees	Asst 1 Out
Sep 24	Nonparametric Methods	
<i>Sep 26</i>	<i>Parametric vs. Nonparametric</i>	
Sep 29	Multi-class Classification	
Oct 1	Probabilistic Classifiers	Asst 1 In
<i>Oct 3</i>	<i>(No tutorial)</i>	
Oct 6	Probabilistic Classifiers II	Asst 2 Out
Oct 8	Naive Bayes	
<i>Oct 10</i>	<i>Naive Bayes &amp; Gaussian Bayes classifiers</i>	
[Oct 13]	Thanksgiving: No class	
Oct 15	Neural Networks I	
<i>Oct 17</i>	<i>Mid-term review</i>	

## CLASS SCHEDULE, Part 2

Date	Topic	Assignments
Oct 20	Neural Networks II	
Oct 22	Clustering	Asst 2 In
<i>Oct 24</i>	<i>Mid-term</i>	
Oct 27	EM & Mixtures of Gaussians	
Oct 29	PCA & Autoencoders	Asst 3 Out
<i>Oct 31</i>	<i>Mixture of Gaussians</i>	
Nov 3	Support Vector Machines	
Nov 5	Kernels and Margins	
<i>Nov 7</i>	<i>SVMs</i>	
Nov 10	Ensemble Methods	
Nov 12	Ensemble Methods II	Asst 3 In
<i>Nov 14</i>	<i>Bagging &amp; Boosting</i>	Asst 4 Out
[Nov 17]	Mid-term break: No class	
Nov 19	Bayesian Learning	
<i>Nov 21</i>		
Nov 24	Computational Learning Theory	
Nov 26	Computational Learning Theory II	
<i>Nov 28</i>		
Dec 1	Reinforcement Learning I	Asst 4 In