# A Course in Machine Learning



Hal Daumé III

For my students and teachers.
  Often the same.

# Table of Contents

Machine learning is a broad and fascinating field. It has been called one of the sexiest fields to work in[1]. It has applications in an incredibly wide variety of application areas, from medicine to advertising, from military to pedestrian. Its importance is likely to grow, as more and more areas turn to it as a way of dealing with the massive amounts of data available.

[1]

## 0.1 How to Use this Book

## 0.2 Why Another Textbook?

The purpose of this book is to provide a *gentle* and *pedagogically organized* introduction to the field. This is in contrast to most existing machine learning texts, which tend to organize things topically, rather than pedagogically (an exception is Mitchell's book[2], but unfortunately that is getting more and more outdated). This makes sense for researchers in the field, but less sense for learners. A second goal of this book is to provide a view of machine learning that focuses on ideas and models, not on math. It is not possible (or even advisable) to avoid math. But math should be there to *aid* understanding, not hinder it. Finally, this book attempts to have minimal dependencies, so that one can fairly easily pick and choose chapters to read. When dependencies exist, they are listed at the start of the chapter, as well as the list of dependencies at the end of this chapter.

[2] ?

The *audience* of this book is anyone who knows differential calculus and discrete math, and can program reasonably well. (A little bit of linear algebra and probability will not hurt.) An undergraduate in their fourth or fifth semester should be fully capable of understanding this material. However, it should also be suitable for first year graduate students, perhaps at a slightly faster pace.

## 0.3   Organization and Auxilary Material

There is an associated web page, `http://ciml.info/`, which contains an online copy of this book, as well as associated code and data. It also contains errate. For instructors, there is the ability to get a solutions manual.

This book is suitable for a single-semester undergraduate course, graduate course or two semester course (perhaps the latter supplemented with readings decided upon by the instructor). Here are suggested course plans for the first two courses; a year-long course could be obtained simply by covering the entire book.

## 0.4   Acknowledgements