

1

650

Introduction and Motivation

651 Machine learning is about designing algorithms that automatically extract
652 valuable from data. The emphasis here is on “automatic”, i.e., machine
653 learning is concerned about general-purpose methodologies that can be
654 applied to many datasets, while producing something that is meaningful.
655 There are three concepts that are at the core of machine learning: data, a
656 model and learning

657 Since machine learning is inherently data driven, *data* is at the core data
658 of machine learning. The goal of machine learning is to design general-
659 purpose methodologies to extract valuable patterns from data, ideally
660 without much domain-specific expertise. For example, given a large corpus
661 of documents (e.g., books in many libraries), machine learning methods
662 can be used to automatically find relevant topics that are shared across
663 documents (Hoffman et al., 2010). To achieve this goal, we design *mod- model*
664 *els* that are typically related to the process that generates data, similar to
665 the dataset we are given. For example, in a regression setting, the model
666 would describe a function that maps inputs to real-valued outputs. To
667 paraphrase Mitchell (1997): A model is said to learn from data if its per-
668 formance on a given task improves after the data is taken into account.
669 The goal is to find good models that generalize well to yet unseen data,
670 which we may care about in the future. *Learning* can be understood as a learning
671 way to automatically find patterns and structure in data by optimizing the
672 parameters of the model.

673 While machine learning has seen many success stories, and software is
674 readily available to design and train rich and flexible machine learning
675 systems, we believe that the mathematical foundations of machine learn-
676 ing are important in order to understand fundamental principles upon
677 which more complicated machine learning systems are built. Understand-
678 ing these principles can facilitate creating new machine learning solutions,
679 understanding and debugging existing approaches and learning about the
680 inherent assumptions and limitations of the methodologies we are work-
681 ing with.

1.1 Finding Words for Intuitions

682

683 A challenge we face regularly in machine learning is that concepts and
 684 words are slippery, and a particular component of the machine learning
 685 system can be abstracted to different mathematical concepts. For example,
 686 the word “algorithm” is used in at least two different senses in the con-
 687 text of machine learning. In the first sense, we use the phrase “machine
 688 learning algorithm” to mean a system that makes predictions based on in-
 689 put data. We refer to these algorithms as *predictors*. In the second sense,
 690 we use the exact same phrase “machine learning algorithm” to mean a
 691 system that adapts some internal parameters of the predictor so that it
 692 performs well on future unseen input data. Here we refer to this adapta-
 693 tion as *training* a system.

predictors

training

694 This book will not resolve the issue of ambiguity, but we want to high-
 695 light upfront that, depending on the context, the same expressions can
 696 mean different things. However, we attempt to make the context suffi-
 697 ciently clear to reduce the level of ambiguity.

698 The first part of this book introduces the mathematical concepts and
 699 foundations needed to talk about the three main components of a machine
 700 learning system: data, models and learning. We will briefly outline these
 701 components here, and we will revisit them again in Chapter 8 once we
 702 have discussed the necessary mathematical concepts.

703 While not all data is numerical it is often useful to consider data in
 704 a number format. In this book, we assume that *data* has already been
 705 appropriately converted into a numerical representation suitable for read-
 706 ing into a computer program. Therefore, we think of data as vectors. As
 707 another illustration of how subtle words are, there are (at least) three
 708 different ways to think about vectors: a vector as an array of numbers (a
 709 computer science view), a vector as an arrow with a direction and magni-
 710 tude (a physics view), and a vector as an object that obeys addition and
 711 scaling (a mathematical view).

data as vectors

model

712 A *model* is typically used to describe a process for generating data, sim-
 713 ilar to the dataset at hand. Therefore, good models can also be thought of
 714 simplified versions of the real (unknown) data-generating process, captur-
 715 ing aspects that are relevant for modeling the data and extracting hidden
 716 patterns from it. A good model can then be used to predict what would
 717 happen in the real world without performing real-world experiments.

learning

718 We now come to the crux of the matter, the *learning* component of
 719 machine learning. Assume we are given a dataset and a suitable model.
 720 *Training* the model means to use the data available to optimize some pa-
 721 rameters of the model with respect to a utility function that evaluates how
 722 well the model predicts the training data. Most training methods can be
 723 thought of as an approach analogous to climbing a hill to reach its peak.
 724 In this analogy, the peak of the hill corresponds to a maximum of some
 725 desired performance measure. However, in practice, we are interested in

726 the model to perform well on unseen data. Performing well on data that
727 we have already seen (training data) may only mean that we found a
728 good way to memorize the data. However, this may not generalize well to
729 unseen data, and, in practical applications, we often need to expose our
730 machine learning system to situations that it has not encountered before.

731 Let us summarize the main concepts of machine learning that we cover
732 in this book:

- 733 • We represent data as vectors.
- 734 • We choose an appropriate model, either using the probabilistic or opti-
735 mization view.
- 736 • We learn from available data by using numerical optimization methods
737 with the aim that the model performs well on data not used for training.

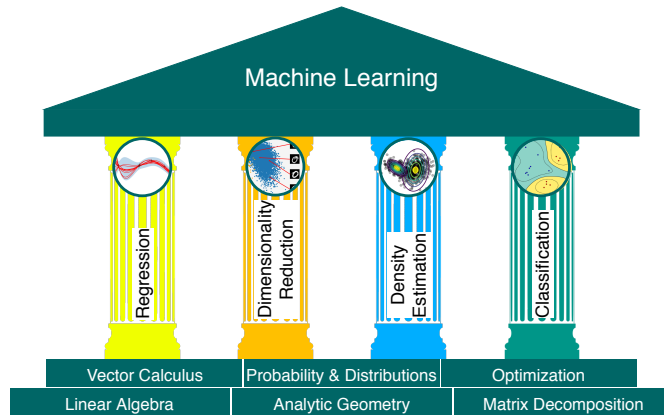
738 1.2 Two Ways to Read this Book

739 We can consider two strategies for understanding the mathematics for
740 machine learning:

- 741 • **Bottom-up:** Building up the concepts from foundational to more ad-
742 vanced. This is often the preferred approach in more technical fields,
743 such as mathematics. This strategy has the advantage that the reader
744 at all times is able to rely on their previously learned concepts. Unfor-
745 tunately, for a practitioner many of the foundational concepts are not
746 particularly interesting by themselves, and the lack of motivation means
747 that most foundational definitions are quickly forgotten.
- 748 • **Top-down:** Drilling down from practical needs to more basic require-
749 ments. This goal-driven approach has the advantage that the reader
750 knows at all times why they need to work on a particular concept, and
751 there is a clear path of required knowledge. The downside of this strat-
752 egy is that the knowledge is built on potentially shaky foundations, and
753 the reader has to remember a set of words for which they do not have
754 any way of understanding.

755 We decided to write this book in a modular way to separate foundational
756 (mathematical) concepts from applications so that this book can be read
757 in both ways. The book is split into two parts, where Part I lays the math-
758 ematical foundations and Part II applies the concepts from Part I to a set
759 of fundamental machine learning problems, which form four pillars of
760 machine learning as illustrated in Figure 1.1: regression, dimensionality
761 reduction, density estimation, and classification. Chapters in Part I mostly
762 build upon the previous ones, but it is possible to skip a chapter and work
763 backward if necessary. Chapters in Part II are only loosely coupled and
764 can be read in any order. There are many pointers forward and backward
765 between the two parts of the book to link mathematical concepts with
766 machine learning algorithms.

Figure 1.1 The foundations and four pillars of machine learning.



767 *Of course there are more than two ways to read this book.* Most readers
 768 learn using a combination of top-down and bottom-up approaches, some-
 769 times building up basic mathematical skills before attempting more com-
 770 plex concepts, but also choosing topics based on applications of machine
 771 learning.

772 *Part I is about Mathematics*

773 The four pillars of machine learning we cover in this book (see Figure 1.1)
 774 require a solid mathematical foundation, which is laid out in Part I.

linear algebra

775 We represent numerical data as vectors and represent a table of such
 776 data as a matrix. The study of vectors and matrices is called *linear algebra*,
 777 which we introduce in Chapter 2. The collection of vectors as a matrix is
 778 also described there.

analytic geometry

779 Given two vectors representing two objects in the real world we want
 780 to make statements about their similarity. The idea is that vectors that
 781 are similar should be predicted to have similar outputs by our machine
 782 learning algorithm (our predictor). To formalize the idea of similarity be-
 783 tween vectors, we need to introduce operations that take two vectors as
 784 input and return a numerical value representing their similarity. The con-
 785 struction of similarity and distances is central to *analytic geometry* and is
 786 discussed in Chapter 3.

matrix
decomposition

787 In Chapter 4, we introduce some fundamental concepts about matri-
 788 ces and *matrix decomposition*. Some operations on matrices are extremely
 789 useful in machine learning, and they allow for an intuitive interpretation
 790 of the data and more efficient learning.

791 We often consider data to be noisy observations of some true underly-
 792 ing signal. We hope that by applying machine learning we can identify the
 793 signal from the noise. This requires us to have a language for quantify-
 794 ing what “noise” means. We often would also like to have predictors that
 795 allow us to express some sort of uncertainty, e.g., to quantify the confi-
 796 dence we have about the value of the prediction at a particular test data

797 point. Quantification of uncertainty is the realm of *probability theory* and
798 is covered in Chapter 6. probability theory

799 To train machine learning models, we typically find parameters that
800 maximize some performance measure. Many optimization techniques re-
801 quire the concept of a gradient, which tells us the direction in which to
802 search for a solution. Chapter 5 is about *vector calculus* and details the con- vector calculus
803 cept of gradients, which we subsequently use in Chapter 7, where we talk
804 about *optimization* to find maxima/minima of functions. optimization

805 *Part II is about Machine Learning*

806 The second part of the book introduces *four pillars of machine learning* as four pillars of
807 shown in Figure 1.1. machine learning

808 We illustrate how the mathematical concepts introduced in the first part
809 of the book are the foundation for each pillar. Broadly speaking, the chap-
810 ters are ordered by difficulty (in ascending order).

811 In Chapter 8, we restate the three components of machine learning
812 (data, models and parameter estimation) in a mathematical fashion. In
813 addition, we provide some guidelines for building experimental set-ups
814 that guard against overly optimistic evaluations of machine learning sys-
815 tems. Recall that the goal is to build a predictor that performs well on
816 unseen data.

817 In Chapter 9, we will have a close look at *linear regression*, where our linear regression
818 objective is to find functions that map inputs $\mathbf{x} \in \mathbb{R}^D$ to corresponding
819 observed function values $y \in \mathbb{R}$, which we can interpret as the labels of
820 their respective inputs. We will discuss classical model fitting (parameter
821 estimation) via maximum likelihood and maximum a posteriori estimation
822 as well as Bayesian linear regression where we integrate the parameters
823 out instead of optimizing them.

824 Chapter 10 focuses on *dimensionality reduction*, the second pillar in Fig- dimensionality
825 ure 1.1, using principal component analysis. The key objective of dimen- reduction
826 sionality reduction is to find a compact, lower-dimensional representation
827 of high-dimensional data $\mathbf{x} \in \mathbb{R}^D$, which is often easier to analyze than
828 the original data. Unlike regression, dimensionality reduction is only con-
829 cerned about modeling the data – there are no labels associated with a
830 data point \mathbf{x} .

831 In Chapter 11, we will move to our third pillar: *density estimation*. The density estimation
832 objective of density estimation is to find a probability distribution that de-
833 scribes a given dataset. We will focus on Gaussian mixture models for this
834 purpose, and we will discuss an iterative scheme to find the parameters of
835 this model. As in dimensionality reduction, there are no labels associated
836 with the data points $\mathbf{x} \in \mathbb{R}^D$. However, we do not seek a low-dimensional
837 representation of the data. Instead, we are interested in a density model
838 that describes the data.

839 Chapter 12 concludes the book with an in-depth discussion of the fourth classification
840 pillar: *classification*. We will discuss classification in the context of support classification

841 vector machines. Similar to regression (Chapter 9) we have inputs x and
842 corresponding labels y . However, unlike regression where the labels were
843 real-valued, the labels in classification are integers, which requires special
844 care.

845 **1.3 Exercises and Feedback**

846 We provide some exercises in Part I, which can be done mostly by pen and
847 paper. For Part II we provide programming tutorials (jupyter notebooks)
848 to explore some properties of the machine learning algorithms we discuss
849 in this book.

850 We appreciate that Cambridge University Press strongly supports our
851 aim to democratize education and learning by making this book freely
852 available for download at

853 <https://mml-book.com>

854 where tutorials, errata and additional materials can be found. Mistakes
855 can be reported and feedback provided using the URL above.