
Introduction and Motivation

1.1 Finding Words for Intuitions

573

574 Machine learning is about designing algorithms that learn from data. The
 575 goal is to find good models that generalize well to future data. The chal-
 576 lenge is that the concepts and words are slippery, and a particular compo-
 577 nent of the machine learning system can be abstracted to different math-
 578 ematical concepts. For example, the word “algorithm” is used in at least
 579 two different senses in the context of machine learning. In the first sense,
 580 we use the phrase “machine learning algorithm” to mean a system that
 581 makes predictions based on input data. We refer to these algorithms as
 582 *predictors*. In the second sense, we use the exact same phrase “machine
 583 learning algorithm” to mean a system that adapts some internal parame-
 584 ters of the predictor so that it performs well on future unseen input data.
 585 Here we refer to this adaptation as *training* a predictor.

predictors

training

586 The first part of this book describes the mathematical concepts and
 587 foundations needed to talk about the three main components of a machine
 588 learning system: data, models, and learning. We will briefly outline these
 589 components here, and we will revisit them again in Chapter 8 once we
 590 have the mathematical language under our belt. Adding to the challenge
 591 is the fact that the same English word could mean different mathematical
 592 concepts, and we can only work out the precise meaning via the context.
 593 We already remarked about the overloaded use of the word “algorithm”,
 594 and the reader will be faced with other such phrases. We advise the reader
 595 to use the idea of “type checking” from computer science and apply it
 596 to machine learning concepts. Type checking allows the reader to sanity
 597 check whether the equation that they are considering contains inputs and
 598 outputs of the correct type, and whether they are mixing different types
 599 of objects.

600 While not all data is numerical it is often useful to consider data in a
 601 number format. In this book, we assume that the *data* has already been
 602 appropriately converted into a numerical representation suitable for read-
 603 ing into a computer program. In this book, we think of data as vectors.
 604 As another illustration of how subtle words are, there are three different
 605 ways to think about vectors: a vector as an array of numbers (a computer
 606 science view), a vector as an arrow with a direction and magnitude (a

data

data as vectors

607 physics view), and a vector as an object that obeys addition and scaling (a
608 mathematical view).

model

609 What is a *model*? Models are simplified versions of reality, which capture
610 aspects of the real world that are relevant to the task. Users of the model
611 need to understand what the model does not capture, and hence obtain
612 an appreciation of the limitations of it. Applying models without knowing
613 their limitations is like driving a vehicle without knowing whether it can
614 turn left or not. Machine learning algorithms adapt to data, and therefore
615 their behavior will change as it learns. Applying machine learning models
616 without knowing their limitations is like sitting in a self-driving vehicle
617 without knowing whether it has encountered enough left turns during its
618 training phase. In this book, we use the word “model” to distinguish be-
619 tween two schools of thought about the construction of machine learning
620 predictors: the probabilistic view and the optimization view. The reader
621 is referred to Domingos (2012) for a more general introduction to the five
622 schools of machine learning.

learning

623 We now come to the crux of the matter, the *learning* component of
624 machine learning. Assume we have a way to represent data as vectors
625 and that we have an appropriate model. We are interested in training
626 our model based on data so that it performs well on unseen data. Pre-
627 dicting well on data that we have already seen (training data) may only
628 mean that we found a good way to memorize the data. However, this may
629 not generalize well to unseen data, and in practical applications we often
630 need to expose our machine learning system to situations that it has not
631 encountered before. We use numerical methods to find good parameters
632 that “fit” the model to data, and most training methods can be thought of
633 as an approach analogous to climbing a hill to reach its peak. The peak
634 of the hill corresponds to a maximization of some desired performance
635 measure. The challenge is to design algorithms that learn from past data
636 but generalizes well.

637 Let us summarize the main concepts of machine learning:

- 638 • We use domain knowledge to represent data as vectors.
- 639 • We choose an appropriate model, either using the probabilistic or opti-
640 mization view.
- 641 • We learn from past data by using numerical optimization methods with
642 the aim that it performs well on unseen data.

643 1.2 Two Ways to Read this Book

644 We can consider two strategies for understanding the mathematics for
645 machine learning:

- 646 • Building up the concepts from foundational to more advanced. This is
647 often the preferred approach in more technical fields, such as mathe-
648 matics. This strategy has the advantage that the reader at all times is

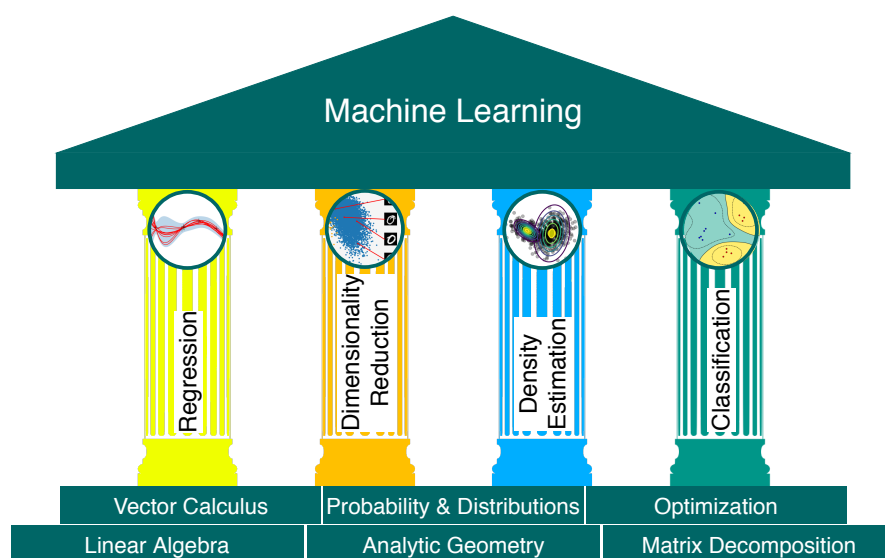


Figure 1.1 The foundations and four pillars of machine learning.

649 able to rely on their previously learned definitions, and there are no
 650 murky hand-wavy arguments that the reader needs to take on faith.
 651 Unfortunately, for a practitioner many of the foundational concepts are
 652 not particularly interesting by themselves, and the lack of motivation
 653 means that most foundational definitions are quickly forgotten.

- 654 • Drilling down from practical needs to more basic requirements. This
 655 goal-driven approach has the advantage that the reader knows at all
 656 times why they need to work on a particular concept, and there is a
 657 clear path of required knowledge. The downside of this strategy is that
 658 the knowledge is built on shaky foundations, and the reader has to
 659 remember a set of words for which they do not have any way of under-
 660 standing.

661 This book is split into two parts, where Part I lays the mathematical
 662 foundations and Part II applies the concepts from Part I to a set of basic
 663 machine learning problems, which form four pillars of machine learning
 664 as illustrated in Figure 1.1.

665 *Part I is about Mathematics*

666 We represent numerical data as vectors and represent a table of such data
 667 as a matrix. The study of vectors and matrices is called *linear algebra*,
 668 which we introduce in Chapter 2. The collection of vectors as a matrix is
 669 also described there. Given two vectors, representing two objects in the
 670 real world, we want to be able to make statements about their similarity.
 671 The idea is that vectors that are similar should be predicted to have similar
 672 outputs by our machine learning algorithm (our predictor). To formalize
 673 the idea of similarity between vectors, we need to introduce operations

linear algebra

674 that take two vectors as input and return a numerical value represent-
 675 ing their similarity. This construction of similarity and distances is called
 analytic geometry 676 *analytic geometry* and is discussed in Chapter 3. In Chapter 4, we introduce
 matrix 677 some fundamental concepts about matrices and *matrix decomposition*. It
 decomposition 678 turns out that operations on matrices are extremely useful in machine
 679 learning, and we use them for representing data as well as for modeling.

680 We often consider data to be noisy observations of some true underlying
 681 signal, and hope that by applying machine learning we can identify
 682 the signal from the noise. This requires us to have a language for quanti-
 683 fying what noise means. We often would also like to have predictors that
 684 allow us to express some sort of uncertainty, e.g., to quantify the confi-
 685 dence we have about the value of the prediction for a particular test data
 probability theory 686 point. Quantification of uncertainty is the realm of *probability theory* and
 687 is covered in Chapter 6. Instead of considering a predictor as a single func-
 688 tion, we could consider predictors to be probabilistic models, i.e., models
 689 describing the distribution of possible functions.

690 To apply hill-climbing approaches for training machine learning models,
 691 we need to formalize the concept of a gradient, which tells us the direc-
 692 tion which to search for a solution. This idea of the direction to search
 calculus 693 is formalized by *calculus*, which we present in Chapter 5. How to use a
 694 sequence of these search directions to find the top of the hill is called
 optimization 695 *optimization*, which we introduce in Chapter 7.

696 It turns out that the mathematics for discrete categorical data is differ-
 697 ent from the mathematics for continuous real numbers. Most of machine
 698 learning assumes continuous variables, and except for Chapter 6 the other
 699 chapters in Part I of the book only discuss continuous variables. However,
 700 for many application domains, data is categorical in nature, and naturally
 701 there are machine learning problems that consider categorical variables.
 702 For example, we may wish to model sex (male/female). Since we assume
 703 that our data is numerical, we encode sex as the numbers -1 and $+1$
 704 for male and female, respectively. However, it is worth keeping in mind
 705 when modeling that sex is a categorical variable, and the actual differ-
 706 ence in value between the two numbers should not have any meaning in
 707 the model. This distinction between continuous and categorical variables
 708 gives rise to different machine learning approaches.

709 *Part II is about Machine Learning*

four pillars of 710 The second part of the book introduces *four pillars of machine learning* as
 machine learning 711 listed in Table 1.1. The rows in the table distinguish between problems
 712 where the variable of interest is continuous or categorical. We illustrate
 713 how the mathematical concepts introduced in the first part of the book
 714 can be used to design machine learning algorithms. In Chapter 8, we re-
 715 state the three components of machine learning (data, models and param-
 716 eter estimation) in a mathematical fashion. In addition, we provide some
 717 guidelines for building experimental setups that guard against overly op-

	Supervised	Unsupervised
Continuous latent variables	Regression (Chapter 9)	Dimensionality reduction (Chapter 10)
Categorical latent variables	Classification (Chapter 12)	Density estimation (Chapter 11)

Table 1.1 The four pillars of machine learning

718 timistic evaluations of machine learning systems. Recall that the goal is to
719 build a predictor that performs well on future data.

720 The terms “supervised” and “unsupervised” (the columns in Table 1.1)
721 learning refer to the question of whether or not we provide the learning
722 algorithm with labels during training. An example use case of *supervised*
723 *learning* is when we build a classifier to decide whether a tissue biopsy is
724 cancerous. For training, we provide the machine learning algorithm with
725 a set of images and a corresponding set of annotations by pathologists.
726 This expert annotation is called a *label* in machine learning, and for many
727 supervised learning tasks it is obtained at great cost or effort. After the
728 classifier is trained, we show it an image from a new biopsy and hope that
729 it can accurately predict whether the tissue is cancerous. An example use
730 case of unsupervised learning (using the same cancer biopsy problem) is
731 if we want to visualize the properties of the tissue around which we have
732 found cancerous cells. We could choose two particular features of these
733 images and plot them in a scatter plot. Alternatively we could use all the
734 features and find a two dimensional representation that approximates all
735 the features, and plot this instead. Since this type of machine learning task
736 does not provide a label during training, it is called *unsupervised learning*.
737 The second part of the book provides a brief overview of two fundamental
738 supervised (*regression* and *classification*) and unsupervised (*dimensionality*
739 *reduction* and *density estimation*) machine learning problems.

740 *Of course there are more than two ways to read this book.* Most read-
741 ers learn using a combination of top-down and bottom-up approaches,
742 sometimes building up basic mathematical skills before attempting more
743 complex concepts, but also choosing topics based on applications of ma-
744 chine learning. Chapters in Part I mostly build upon the previous ones, but
745 the reader is encouraged to skip to a chapter that covers a particular gap
746 the reader’s knowledge and work backwards if necessary. Chapters in Part
747 II are loosely coupled and are intended to be read in any order. There are
748 many pointers forward and backward between the two parts of the book
749 to assist the reader in finding their way.

750 1.3 Exercises and Feedback

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

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

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

759 where you can also find the tutorials, errata and additional materials. You
760 can also report mistakes and provide feedback using the URL above.