

# Chapter 1

## Introduction

Inventors have long dreamed of creating machines that think. This desire dates back to at least the time of ancient Greece. The mythical figures Pygmalion, Daedalus, and Hephaestus may all be interpreted as legendary inventors, and Galatea, Talos, and Pandora may all be regarded as artificial life (Ovid and Martin , 2004 Sparkes; , 1996 Tandy; , 1997).

When programmable computers were first conceived, people wondered whether such machines might become intelligent, over a hundred years before one was built (Lovelace, 1842). Today, **artificial intelligence** (AI) is a thriving field with many practical applications and active research topics. We look to intelligent software to automate routine labor, understand speech or images, make diagnoses in medicine and support basic scientific research.

In the early days of artificial intelligence, the field rapidly tackled and solved problems that are intellectually difficult for human beings but relatively straightforward for computers—problems that can be described by a list of formal, mathematical rules. The true challenge to artificial intelligence proved to be solving the tasks that are easy for people to perform but hard for people to describe formally—problems that we solve intuitively, that feel automatic, like recognizing spoken words or faces in images.

This book is about a solution to these more intuitive problems. This solution is to allow computers to learn from experience and understand the world in terms of a hierarchy of concepts, with each concept defined through its relation to simpler concepts. By gathering knowledge from experience, this approach avoids the need for human operators to formally specify all the knowledge that the computer needs. The hierarchy of concepts enables the computer to learn complicated concepts by building them out of simpler ones. If we draw a graph showing how these concepts

# 第一章

## 介绍

发明家们一直梦想着创造出能够思考的机器，这一愿望至少可以追溯到古希腊时期，神话人物皮格马利翁、代达罗斯和赫菲斯托斯都可能被视为传奇发明家，而伽拉忒亚、塔洛斯和潘多拉都可能被视为人工生命（Ovid and Martin 2004 Sparkes 1996 Tandy 1997；，；，）。

当可编程计算机首次被设想出来时，人们想知道这种机器是否会变得智能，这比实际制造出一台机器早了一百多年（Lovelace 1842）。今天，人工智能（AI）是一个蓬勃发展的领域，拥有许多实际应用和活跃的研究课题。我们期待智能软件能够自动化日常劳动、理解语音或图像、进行医学诊断并支持基础科学的研究。

在人工智能发展的早期，该领域迅速处理和解决了那些对人类来说在智力上很困难但对计算机来说相对简单的问题——这些问题可以用一系列正式的数学规则来描述。事实证明，人工智能面临的真正挑战是解决那些对人类来说很容易执行但难以正式描述的任务——我们可以直观地、自动地解决这些问题，比如识别口头单词或图像中的人脸。

本书旨在解决这些更直观的问题。该解决方案旨在让计算机从经验中学习，并根据概念的层级结构理解世界，每个概念都通过其与更简单概念的关系来定义。通过从经验中收集知识，这种方法避免了人类操作员正式指定计算机所需的所有知识。概念的层级结构使计算机能够通过从简单概念构建复杂概念来学习它们。如果我们画一个图表来展示这些概念如何

are built on top of each other, the graph is deep, with many layers. For this reason, we call this approach to AI **deeplearning**.

Many of the early successes of AI took place in relatively sterile and formal environments and did not require computers to have much knowledge about the world. For example, IBM's Deep Blue chess-playing system defeated world champion Garry Kasparov in 1997 (Hsu, 2002). Chess is of course a very simple world, containing only sixty-four locations and thirty-two pieces that can move in only rigidly circumscribed ways. Devising a successful chess strategy is a tremendous accomplishment, but the challenge is not due to the difficulty of describing the set of chess pieces and allowable moves to the computer. Chess can be completely described by a very brief list of completely formal rules, easily provided ahead of time by the programmer.

Ironically, abstract and formal tasks that are among the most difficult mental undertakings for a human being are among the easiest for a computer. Computers have long been able to defeat even the best human chess player but only recently have begun matching some of the abilities of an average human being to recognize objects or speech. A person's everyday life requires an immense amount of knowledge about the world. Much of this knowledge is subjective and intuitive, and therefore difficult to articulate in a formal way. Computers need to capture this same knowledge in order to behave in an intelligent way. One of the key challenges in artificial intelligence is how to get this informal knowledge into a computer.

Several artificial intelligence projects have sought to hard-code knowledge about the world in formal languages. A computer can reason automatically about statements in these formal languages using logical inference rules. This is known as the **knowledgebase** approach to artificial intelligence. None of these projects has led to major success. One of the most famous such projects is Cyc (Lenat and Guha, 1989). Cyc is an inference engine and a database of statements in a language called CycL. These statements are entered by a staff of human supervisors. It is an unwieldy process. People struggle to devise formal rules with enough complexity to accurately describe the world. For example, Cyc failed to understand a story about a person named Fred shaving in the morning (Linde, 1992). Its inference engine detected an inconsistency in the story: it knew that people do not have electrical parts, but because Fred was holding an electric razor, it believed the entity "Fred While Shaving" contained electrical parts. It therefore asked whether Fred was still a person while he was shaving.

The difficulties faced by systems relying on hard-coded knowledge suggest that AI systems need the ability to acquire their own knowledge, by extracting

彼此叠加，图形很深，有很多层。因此，我们将这种人工智能方法称为深度学习。人工智能的许多早期成功都发生在相对贫瘠和正式的环境中，并且不需要计算机对世界有太多了解。例如，IBM 的 DeepBlue 国际象棋系统在 1997 年击败了世界冠军加里卡斯帕罗夫 (,)。国际象棋当然是一个非常简单的世界，只包含 64 个位置和 32 个棋子，只能以严格限定的方式移动。设计一个成功的国际象棋策略是一项巨大的成就，但挑战并不在于难以向计算机描述棋子集和允许的移动。国际象棋可以通过非常简短的完全正式规则列表轻松描述由程序员提前提供。

讽刺的是，对于人类来说最困难的脑力劳动中的抽象和正式任务对于计算机来说是最容易的。计算机早已能够击败最好的人类棋手，但直到最近才开始匹配普通人类的一些能力，例如识别物体或语音。一个人的日常生活需要大量关于世界的知识。这些知识大部分是主观和直觉的，因此很难以正式的方式表达。计算机需要捕捉这些相同的知识才能以智能的方式运行。人工智能的关键挑战之一是如何将这些非正式知识输入计算机。

一些人工智能项目试图将关于世界的知识硬编码到非正式语言中。计算机可以使用逻辑推理规则自动推理这些正式语言中的语句。这被称为人工智能的知识库方法。这些项目都没有取得重大成功。最著名的此类项目之一是 Cyc (Lenat and Guha 1989, )。Cyc 是一个推理引擎和名为 CycL 的语言语句数据库。这些语句由人类主管人员输入。这是一个笨拙的过程。人们努力设计足够复杂的正式规则来准确描述世界。例如，Cyc 未能理解一个名叫 Fred 的人早上刮胡子的故事 (,)。它的推理 Linde 1992 引擎检测到故事中的不一致之处：它知道人没有电气部件，但是因为弗雷德拿着一把电动剃须刀，它认为实体“弗雷德在剃须时”包含电气部件。因此它询问弗雷德在剃须时是否仍然是一个人。

依赖硬编码知识的系统所面临的困难表明，人工智能系统需要能够通过提取

patterns from raw data. This capability is known as **machine learning**. The introduction of machine learning enabled computers to tackle problems involving knowledge of the real world and made decisions that appear subjective. A simple machine learning algorithm called **logistic regression** can determine whether to recommend cesarean delivery ([Mor-Yosef et al. , 1990](#)). A simple machine learning algorithm called **naive Bayes** can separate legitimate e-mail from spam e-mail.

The performance of these simple machine learning algorithms depends heavily on the **representation** of the data they are given. For example, when logistic regression is used to recommend cesarean delivery, the AI system does not examine the patient directly. Instead, the doctor tells the system several pieces of relevant information, such as the presence or absence of a uterine scar. Each piece of information included in the representation of the patient is known as a **feature**. Logistic regression learns how each of these features of the patient correlates with various outcomes. However, it cannot influence how features are defined in any way. If logistic regression were given an MRI scan of the patient, rather than the doctor's formalized report, it would not be able to make useful predictions. Individual pixels in an MRI scan have negligible correlation with any complications that might occur during delivery.

This dependence on representations is a general phenomenon that appears throughout computer science and even daily life. In computer science, operations such as searching a collection of data can proceed exponentially faster if the collection is structured and indexed intelligently. People can easily perform arithmetic on Arabic numerals but find arithmetic on Roman numerals much more time consuming. It is not surprising that the choice of representation has an enormous effect on the performance of machine learning algorithms. For a simple visual example, see figure [1.1](#).

Many artificial intelligence tasks can be solved by designing the right set of features to extract for that task, then providing those features to a simple machine learning algorithm. For example, a useful feature for speaker identification from sound is an estimate of the size of the speaker's vocal tract. This feature gives a strong clue as to whether the speaker is a man, woman, or child.

For many tasks, however, it is difficult to know what features should be extracted. For example, suppose that we would like to write a program to detect cars in photographs. We know that cars have wheels, so we might like to use the presence of a wheel as a feature. Unfortunately, it is difficult to describe exactly what a wheel looks like in terms of pixel values. A wheel has a simple geometric shape, but its image may be complicated by shadows falling on the wheel, the sun glaring off the metal parts of the wheel, the fender of the car or an object in the

模式来自原始数据。这种能力被称为机器学习。机器学习的引入使计算机能够解决涉及现实世界知识的问题并做出看似主观的决策。一种称为逻辑回归的简单机器学习算法可以确定是否建议剖腹产和分娩 (Mor-Yosef 1990etal.,)。一种称为朴素贝叶斯的简单机器学习算法可以将合法电子邮件与垃圾邮件区分开来。

这些简单的机器学习算法的性能在很大程度上取决于它们所给出的数据的表示。例如，当使用逻辑回归来推荐剖腹产时，人工智能系统不会直接检查病人。相反，医生会告诉系统几条相关信息，例如是否存在子宫疤痕。病人表示中包含的每条信息都称为一个特征。

逻辑回归可以显示患者的各项特征与各种结果之间的关联。然而，它无法以任何方式影响这些特征的定义。如果使用患者的 MRI 扫描结果而非医生的正式报告进行逻辑回归分析，则无法做出有用的预测。MRI 图像中的单个像素与分娩过程中可能发生的任何并发症之间的关联微乎其微。

这种对表示的依赖是整个计算机科学甚至日常生活中出现的普遍现象。在计算机科学中，如果数据集合是结构化的并且索引是智能的，那么搜索数据集合等操作可以以指数级的速度进行。人们可以轻松地对阿拉伯数字进行算术运算，但发现对罗马数字进行算术运算要耗费得多。表示的选择对机器学习算法的性能有很大的影响，这并不奇怪。对于一个简单的视觉示例，请参见图 1.1  
许多人工智能任务可以通过设计正确的一组特征来提取，然后将这些特征提供给简单的机器学习算法来解决。例如，从声音中识别说话人的有用特征是估计说话人声道的大小。这个特征提供了一个强有力的线索来表明说话人是男人、女人还是孩子。

然而，对于许多任务来说，很难知道应该提取哪些特征。例如，假设我们想编写一个程序来检测照片中的汽车。我们知道汽车有轮子，所以我们可能想使用轮子的存在作为安全特征。不幸的是，很难用像素值来准确描述轮子的样子。轮子有简单的几何形状，但它的图像可能会因为落在轮子上的阴影、轮子金属部分上的阳光、汽车的挡泥板或车内的物体而变得复杂。

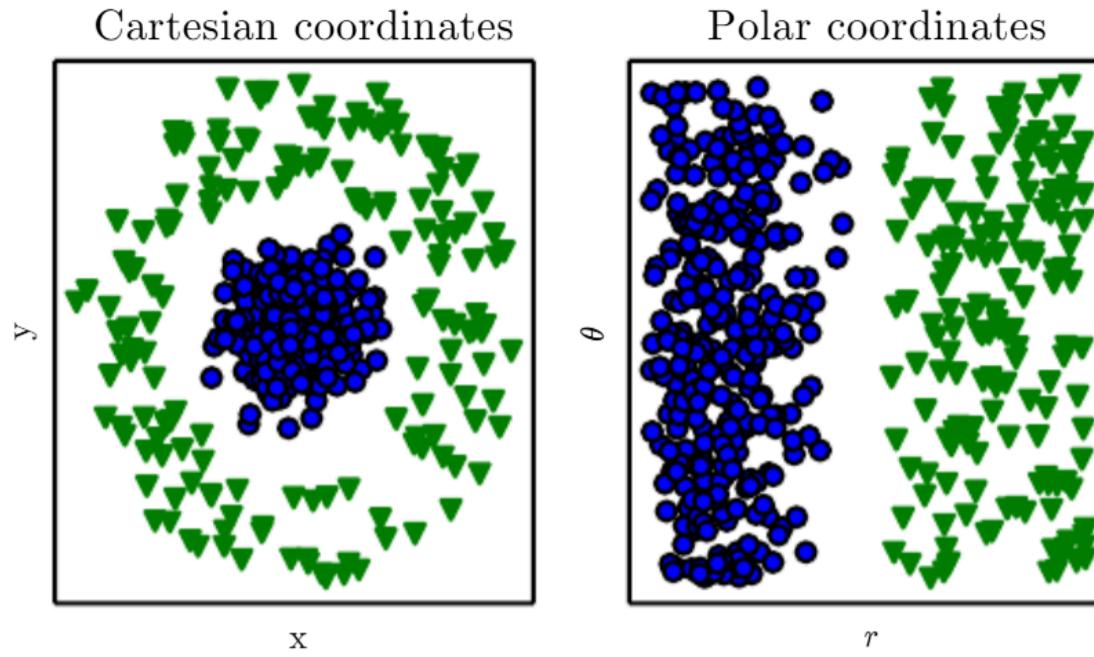


Figure 1.1: Example of different representations: suppose we want to separate two categories of data by drawing a line between them in a scatter plot. In the plot on the left, we represent some data using Cartesian coordinates, and the task is impossible. In the plot on the right, we represent the data with polar coordinates and the task becomes simple to solve with a vertical line. (Figure reproduced in collaboration with David Warde-Farley.)

foreground obscuring part of the wheel, and soon.

One solution to this problem is to use machine learning to discover not only the mapping from representation to output but also the representation itself. This approach is known as **representation learning**. Learned representations often result in much better performance than can be obtained with hand-designed representations. They also enable AI systems to rapidly adapt to new tasks, with minimal human intervention. A representation learning algorithm can discover a good set of features for a simple task in minutes, or for a complex task in hours to months. Manually designing features for a complex task requires a great deal of human time and effort; it can take decades for an entire community of researchers.

The quintessential example of a representation learning algorithm is the **autoencoder**. An autoencoder is the combination of an **encoder** function, which converts the input data into a different representation, and a **decoder** function, which converts the new representation back into the original format. Autoencoders are retrained to preserve as much information as possible when an input is run through the encoder and then the decoder, but they are also trained to make the new representation have various nice properties. Different kinds of autoencoders aim to achieve different kinds of properties.

When designing features or algorithms for learning features, our goal is usually to separate the **factors of variation** that explain the observed data. In this

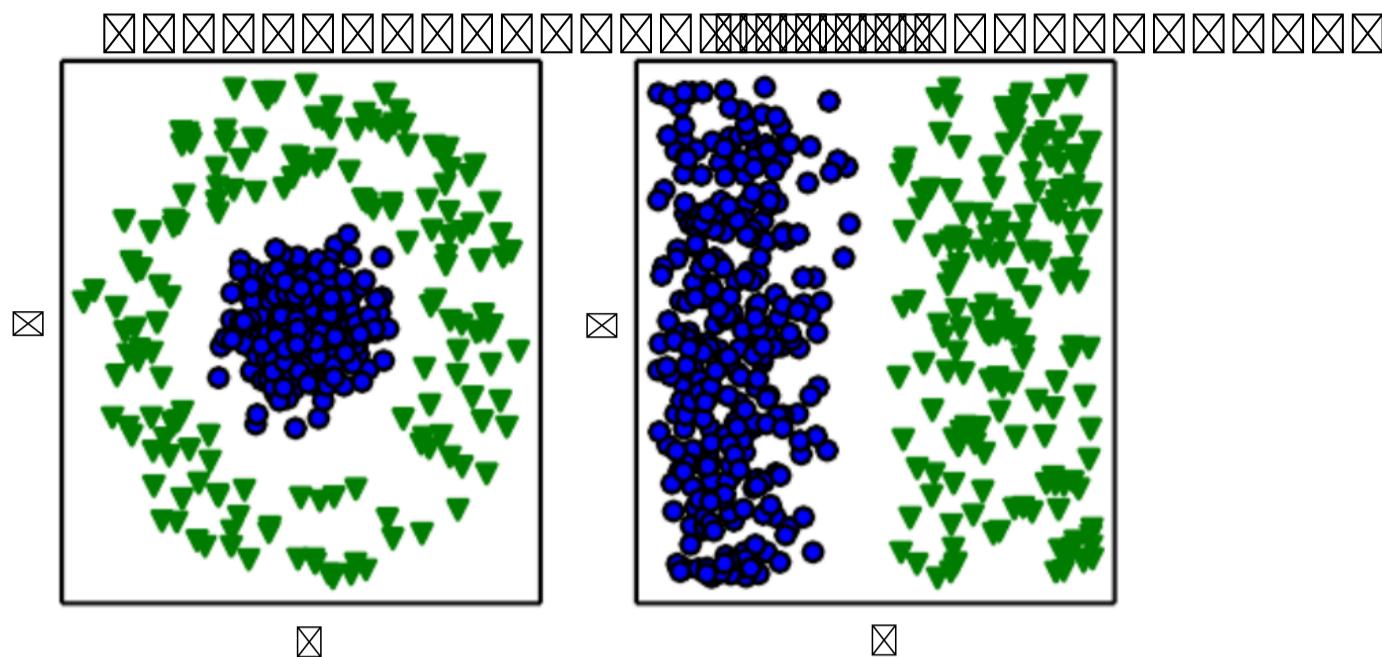


图 1.1：不同表示的示例：假设我们要通过在散点图之间画一条线来分离两类数据。在左侧的图中，我们使用笛卡尔坐标表示一些数据，这个任务是不可能完成的。在右侧的图中，我们用极坐标表示数据，用垂直线就可以很容易地解决这个任务。（该图是与 DavidWarde-Farley 合作制作的。）

前景遮挡了部分车轮，很快。

解决这个问题的一个方法是利用机器学习，不仅发现从表征到输出的映射，还发现表征本身。这种方法被称为表征学习。学习到的表征通常比手工设计的表征性能更好。它们还能使人工智能系统快速适应新任务，并最大限度地减少人工干预。表征学习算法可以在几分钟内为简单任务发现一组良好的特征，或者在几小时到几个月内为复杂任务发现一组良好的特征。手动设计复杂任务的特征需要大量的人力和精力；这可能需要整个研究群体数十年的时间。

演示学习算法的典型例子是自动编码器。自动编码器是编码器功能和解码器功能的组合，编码器功能将输入数据转换为不同的表示，解码器功能将新的表示转换回原始格式。自动编码器经过训练，在输入通过编码器和解码器时保留尽可能多的信息，但它们也经过训练，使新的表示具有各种良好的属性。不同类型的自动编码器旨在实现不同类型的属性。

在设计特征或学习特征的算法时，我们的目标通常是分离变异因素  
解释观察到的数据。在此

context, we use the word “factors” simply to refer to separate sources of influence; the factors are usually not combined by multiplication. Such factors are often not quantities that are directly observed. Instead, they may exist as either unobserved objects or unobserved forces in the physical world that affect observable quantities. They may also exist as constructs in the human mind that provide useful simplifying explanations or inferred causes of the observed data. They can be thought of as concepts or abstractions that help us make sense of the rich variability in the data. When analyzing a speech recording, the factors of variation include the speaker’s age, their sex, their accent and the words they are speaking. When analyzing an image of a car, the factors of variation include the position of the car, its color, and the angle and brightness of the sun.

A major source of difficulty in many real-world artificial intelligence applications is that many of the factors of variation influence every single piece of data we are able to observe. The individual pixels in an image of a red car might be very close to black at night. The shape of the car’s silhouette depends on the viewing angle. Most applications require us to *disentangle* the factors of variation and discard the ones that we do not care about.

Of course, it can be very difficult to extract such high-level, abstract features from raw data. Many of these factors of variation, such as a speaker’s accent, can be identified only using sophisticated, nearly human-level understanding of the data. When it is nearly as difficult to obtain a representation as to solve the original problem, representation learning does not, at first glance, seem to help us.

**Deep learning** solves this central problem in representation learning by introducing representations that are expressed in terms of other, simpler representations. Deep learning enables the computer to build complex concepts out of simpler concepts. Figure 1.2 shows how a deep learning system can represent the concept of an image of a person by combining simpler concepts, such as corners and contours, which are in turn defined in terms of edges.

The quintessential example of a deep learning model is the feedforward deep network, or **multilayer perceptron** (MLP). A multilayer perceptron is just a mathematical function mapping some set of input values to output values. The function is formed by composing many simpler functions. We can think of each application of a different mathematical function as providing a new representation of the input.

The idea of learning the right representation for the data provides one perspective on deep learning. Another perspective on deep learning is that depth enables the computer to learn a multi-step computer program. Each layer of the representation can be thought of as the state of the computer’s memory after

在本文中，我们使用“因素”一词仅指独立的影响源；这些因素通常不会通过乘法组合而成。这些因素通常不是直接观察到的量，而是存在于物理世界中影响可观察量的未观察到的物体或未观察到的力。它们也可能存在于人类思维中，为观察到的数据提供有用的简化解释或推断原因。它们可以被认为是概念或抽象概念，帮助我们理解数据中丰富的可变性。分析语音录音时，变化因素包括说话者的年龄、性别、口音和他们所说的词语。分析汽车图像时，变化因素包括汽车的位置、颜色以及太阳的角度和亮度。

许多现实世界人工智能应用中的一个主要困难来源是，许多变化因素会影响我们能够观察到的每一条数据。汽车图像中的各个像素在夜间可能非常接近黑色。汽车轮廓的形状取决于视角。

大多数应用程序都要求我们考虑变化的因素，  
并丢弃那些我们不关心的因素。

当然，从原始数据中提取这种高级抽象特征是非常困难的。许多这样的变化因素，比如说话者的口音，只有使用复杂的、接近人类水平的数据理解才能识别。当获得一个表征几乎和解决原始问题一样困难时，表征学习乍一看似乎对我们没有帮助。

深度学习通过引入用其他更简单的表示形式来表达的表示形式，解决了表示学习中的这一核心问题。深度学习使计算机能够从简单的概念构建复杂的概念。如图所示，深度学习系统可以通过组合更简单的概念（例如角和轮廓，这些概念又用边来定义）来表示人物图像的概念。

深度学习模型的典型例子是前馈深度网络，或多层感知器（MLP）。多层感知器只是一个将一些输入值映射到输出值的数学函数。该函数由许多更简单的函数组合而成。我们可以将不同数学函数的每个应用视为提供输入的新表示。

学习数据的正确表示的想法为深度学习提供了一个视角。深度学习的另一个视角是深度使计算机能够学习多步骤的计算机程序。表示的每一层都可以被认为 是计算机内存状态的快速变化。

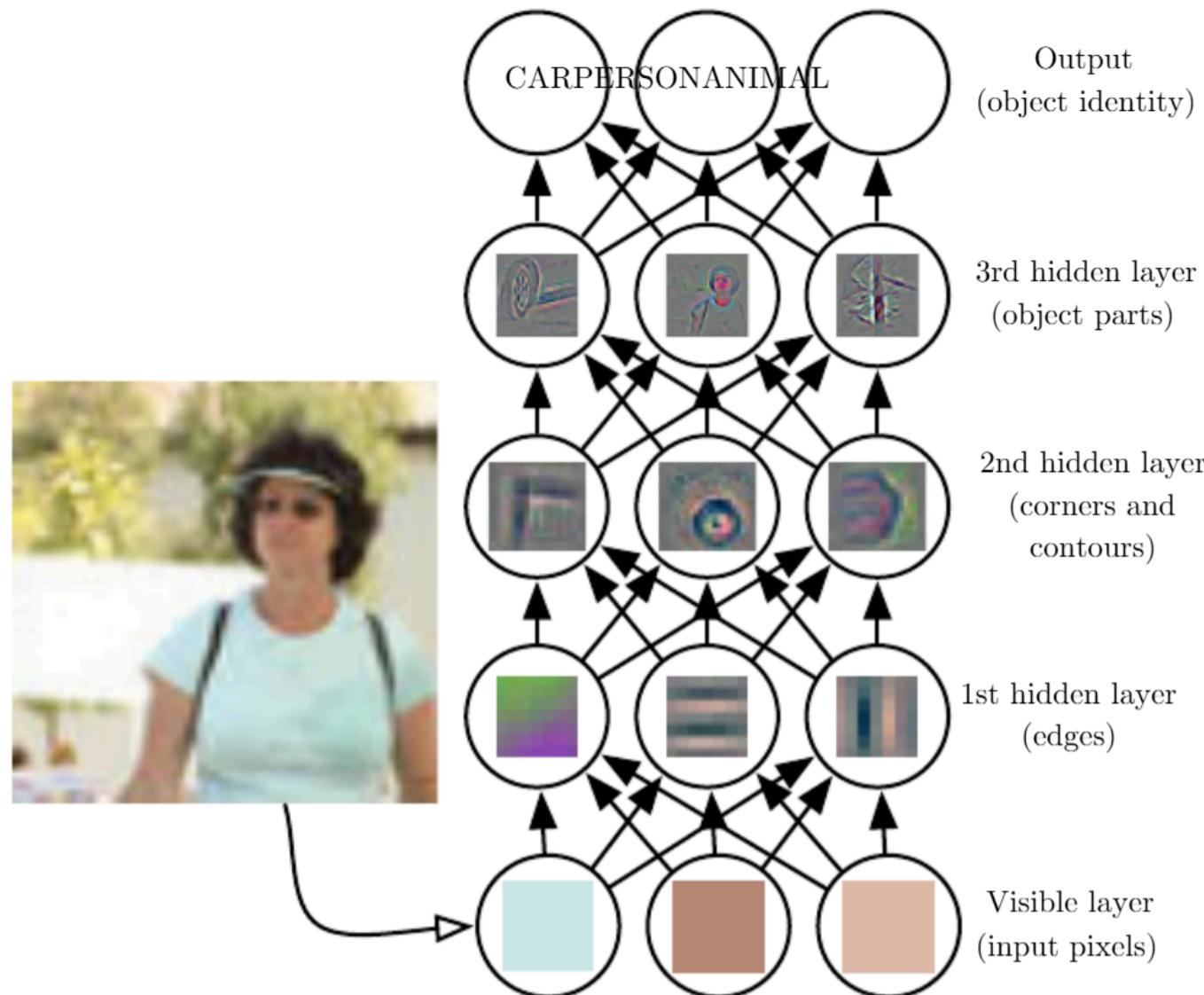


Figure 1.2: Illustration of a deep learning model. It is difficult for a computer to understand the meaning of raw sensory input data, such as this image represented as a collection of pixel values. The function mapping from a set of pixels to an object identity is very complicated. Learning or evaluating this mapping seems insurmountable if tackled directly. Deep learning resolves this difficulty by breaking the desired complicated mapping into a series of nested simple mappings, each described by a different layer of the model. The input is presented at the **visible layer**, so named because it contains the variables that we are able to observe. Then a series of **hidden layers** extracts increasingly abstract features from the image. These layers are called “hidden” because their values are not given in the data; instead the model must determine which concepts are useful for explaining the relationships in the observed data. The images here are visualizations of the kind of feature represented by each hidden unit. Given the pixels, the first layer can easily identify edges, by comparing the brightness of neighboring pixels. Given the first hidden layer’s description of the edges, the second hidden layer can easily search for corners and extended contours, which are recognizable as collections of edges. Given the second hidden layer’s description of the image in terms of corners and contours, the third hidden layer can detect entire parts of specific objects, by finding specific collections of contours and corners. Finally, this description of the image in terms of the object parts it contains can be used to recognize the objects present in the image. Images reproduced with permission from [Zeiler and Fergus 2014](#) () .

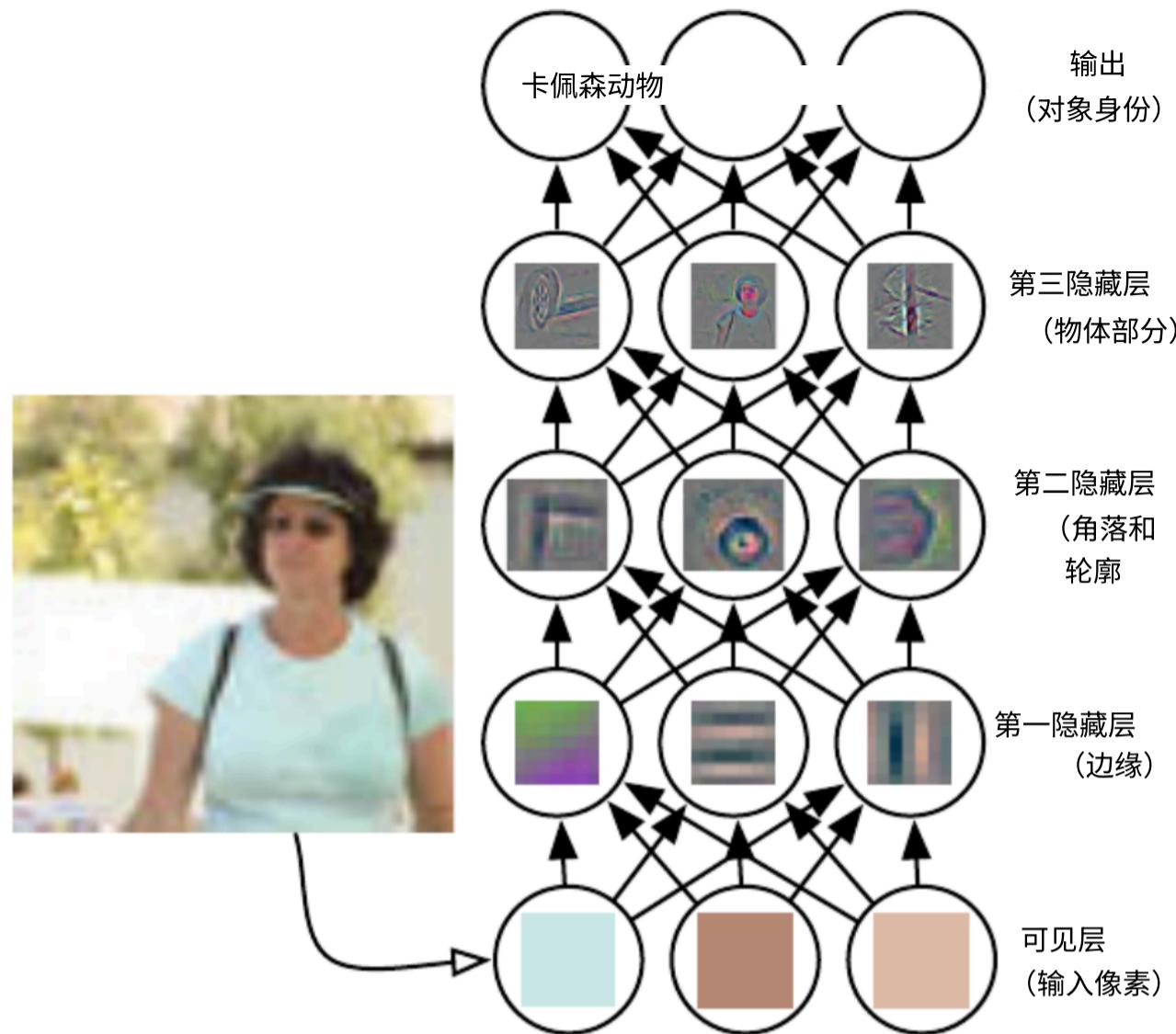


图 1.2：深度学习模型图解。计算机很难理解原始感官输入数据的含义，例如这幅以像素值集合表示的图像。从一组像素到对象身份的功能映射非常复杂。如果直接处理，学习或评估这种映射似乎是无法克服的。深度学习通过将所需的复杂映射分解为一系列嵌套的简单映射来解决这一难题，每个映射由模型的不同层描述。输入呈现在可见层，之所以这样命名，是因为它包含我们能够观察到的变量。一系列隐藏层从图像中提取越来越抽象的特征。这些层被称为“隐藏”，因为它们的值未在数据中给出；相反，模型必须确定哪些概念对于解释观察到的数据中的关系是有用的。这里的图像是每个隐藏单元所代表的特征类型的可视化。给定像素，第一层可以通过比较相邻像素的亮度轻松识别边缘。给定第一个隐藏层的边缘描述，第二个隐藏层可以轻松搜索可识别为边缘集合的角和扩展轮廓。给定第二个隐藏层对图像的角和轮廓的描述，第三个隐藏层可以通过查找特定的轮廓和角的集合来检测特定对象的整个部分。最后，此图像对其所包含的对象部分的描述可用于识别图像中存在的对象。图像经 Zeiler and Fergus 2014 () 许可复制。

executing another set of instructions in parallel. Networks with greater depth can execute more instructions in sequence. Sequential instructions offer great power because later instructions can refer back to the results of earlier instructions. According to this view of deep learning, not all the information in a layer's activations necessarily encodes factors of variation that explain the input. The representation also stores state information that helps to execute a program that can make sense of the input. This state information could be analogous to a counter or pointer in a traditional computer program. It has nothing to do with the content of the inputs specifically, but it helps the model to organize its processing.

There are two main ways of measuring the depth of a model. The first view is based on the number of sequential instructions that must be executed to evaluate the architecture. We can think of this as the length of the longest path through a flowchart that describes how to compute each of the model's outputs given its inputs. Just as two equivalent computer programs will have different lengths depending on which language the program is written in, the same function may be drawn as a flowchart with different depths depending on which functions we allow to be used as individual steps in the flowchart. Figure 1.3 illustrates how this choice of language can give two different measurements for the same architecture.

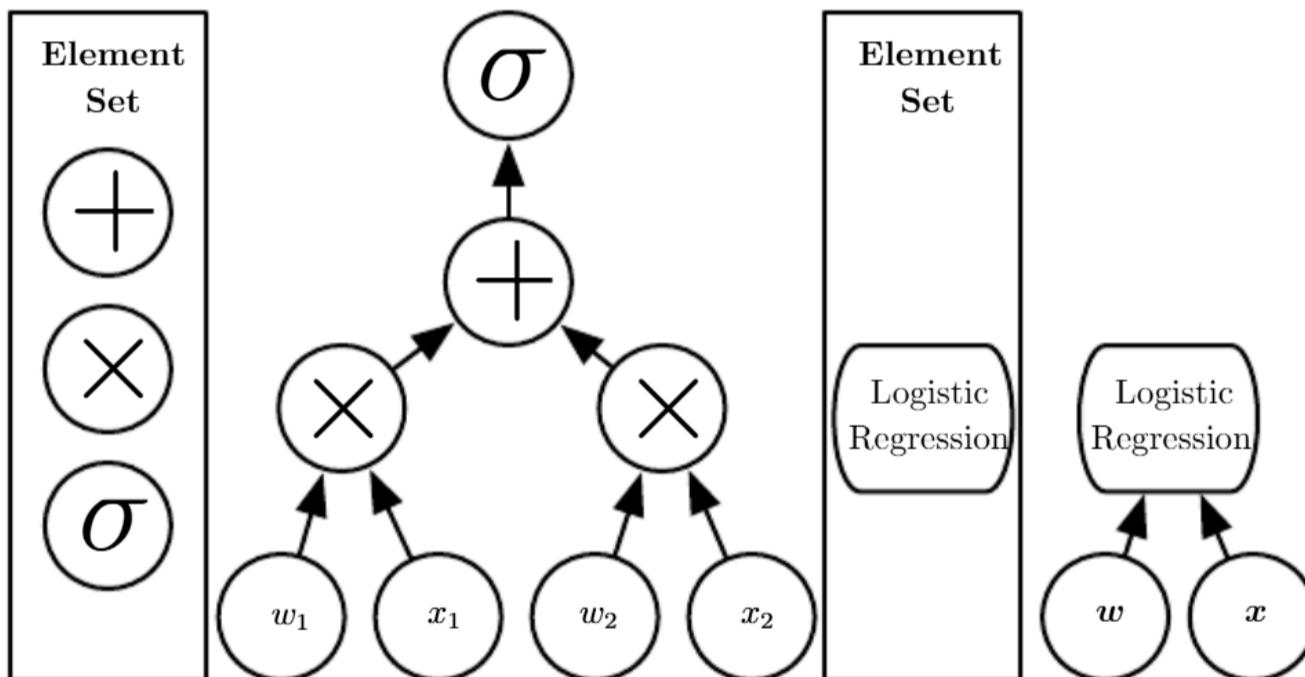


Figure 1.3: Illustration of computational graphs mapping an input to an output where each node performs an operation. Depth is the length of the longest path from input to output but depends on the definition of what constitutes a possible computational step. The computation depicted in these graphs is the output of a logistic regression model,  $\sigma(\mathbf{w}^T \mathbf{x})$ , where  $\sigma$  is the logistic sigmoid function. If we use addition, multiplication and logistic sigmoid as the elements of our computer language, then this model has depth three. If we view logistic regression as an element itself, then this model has depth one.

并行执行另一组指令。深度更大的网络可以按顺序执行更多指令。顺序指令具有强大的功能，因为后续指令可以参考先前指令的其他结果。根据这种深度学习观点，层中激活的信息并非一定都编码了解释输入的变化因素。这种表示还存储了有助于执行能够理解输入的程序的状态信息。这种状态信息可以类似于传统计算机程序中的计数器或指针。它与输入的内容无关，但它有助于模型组织其处理。

有两种主要方法可以测量模型的深度。第一种观点基于评估架构所必须执行的顺序指令的数量。我们可以将其视为流程图中最长路径的长度，该流程图描述了如何根据输入计算每个模型的输出。正如两个等效的计算机程序将具有不同的长度（取决于程序编写的语言）一样，相同的功能可以绘制为具有不同深度的流程图，具体取决于我们允许将哪些功能用作流程图中的单独步骤。该图说明了这种 1.3 语言选择如何为相同的架构提供两种不同的测量结果。

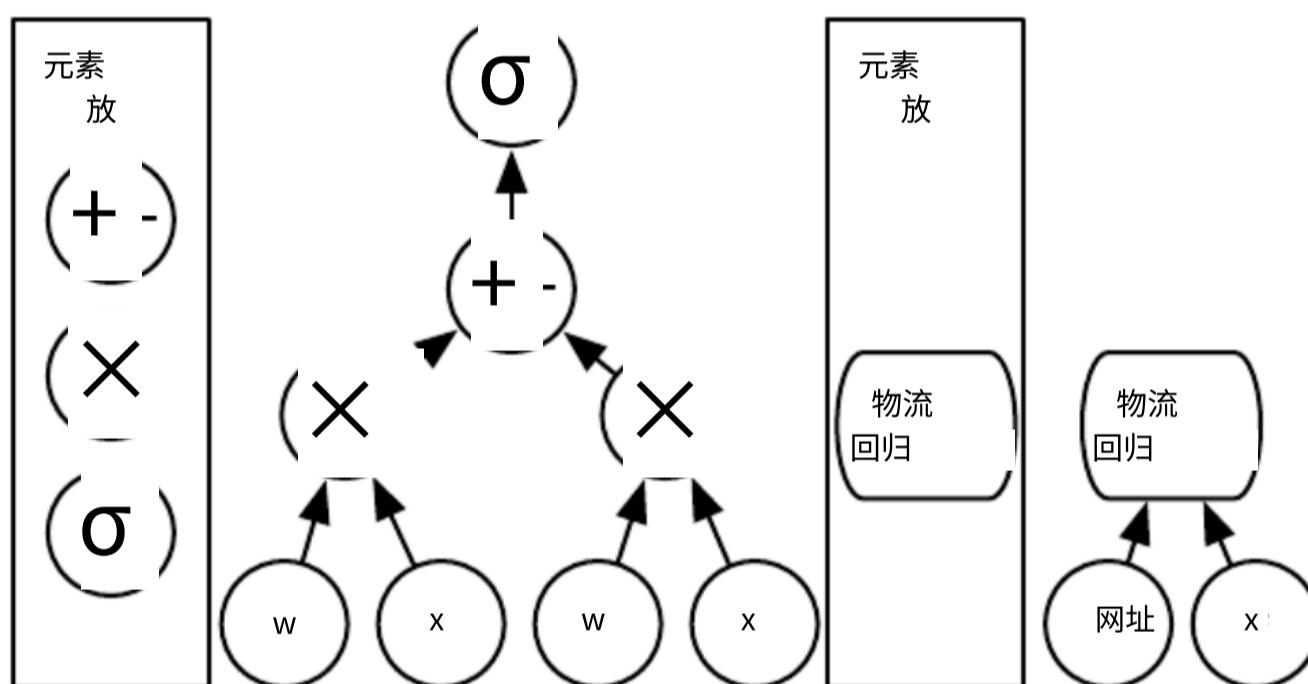


图 1.3：计算图的示意图，将输入映射到输出，其中每个节点执行一项操作。深度是从输入到输出的最长路径的长度，但取决于构成可能计算步骤的定义。这些图中描述的计算是逻辑回归模型  $\sigma(wx)$  的输出，其中  $\sigma$  是逻辑 S 形函数。如果我们使用加法、乘法和逻辑 S 形函数作为计算机语言的元素，则该模型的深度为三。如果我们将逻辑回归本身视为一个元素，则该模型的深度为一。

Another approach, used by deep probabilistic models, regards the depth of a model as being not the depth of the computational graph but the depth of the graph describing how concepts are related to each other. In this case, the depth of the flowchart of the computations needed to compute the representation of each concept may be much deeper than the graph of the concepts themselves. This is because the system's understanding of the simpler concepts can be refined given information about the more complex concepts. For example, an AI system observing an image of a face with one eye in shadow may initially see only one eye. After detecting that a face is present, the system can then infer that a second eye is probably present as well. In this case, the graph of concepts includes only two layers—a layer for eyes and a layer for faces—but the graph of computations includes  $2n$  layers if we refine our estimate of each concept given the other  $n$  times.

Because it is not always clear which of these two views—the depth of the computational graph, or the depth of the probabilistic modeling graph—is most relevant, and because different people choose different sets of small elements from which to construct their graphs, there is no single correct value for the depth of an architecture, just as there is no single correct value for the length of a computer program. Nor is there a consensus about how much depth a model requires to qualify as “deep.” However, deep learning can be safely regarded as the study of models that involve a greater amount of composition of either learned functions or learned concepts than traditional machine learning does.

To summarize, deep learning, the subject of this book, is an approach to AI. Specifically, it is a type of machine learning, a technique that enables computer systems to improve with experience and data. We contend that machine learning is the only viable approach to building AI systems that can operate in complicated real-world environments. Deep learning is a particular kind of machine learning that achieves great power and flexibility by representing the world as a nested hierarchy of concepts, with each concept defined in relation to simpler concepts, and more abstract representations computed in terms of less abstract ones. Figure 1.4 illustrates the relationship between these different AI disciplines. Figure 1.5 gives a high-level schematic of how each works.

## 1.1 Who Should Read This Book?

This book can be useful for a variety of readers, but we wrote it with two target audiences in mind. One of the target audiences is university students (undergraduate or graduate) learning about machine learning, including those who are beginning a career in deep learning and artificial intelligence research. The other

深度概率模型使用的另一种方法认为模型的深度不是计算图的深度，而是描述概念如何相互关联的图的深度。在这种情况下，计算每个概念的表示所需的计算流程图的深度可能比概念本身的图深得多。

这是因为系统对简单概念的理解可以通过提供有关更复杂概念的信息来改进。例如，一个人工智能系统在阴影中观察一张只有一只眼睛的脸的图像时，最初可能只看到一只眼睛。在检测到一张脸后，系统可以推断第二只眼睛也可能存在。在这种情况下，概念图只包括两层——眼睛层和脸层——但如果我们将每个概念的估计都考虑到另一个概念，那么计算图包括  $2n$  层

$n$  次。

由于计算图的深度和概率建模图的深度这两种观点中哪一种最相关并不总是很清楚，而且不同的人会选择不同的最小元素集来构建他们的图，因此对于架构的深度没有单一的正确值，就像对于计算机程序的长度没有单一的正确值一样。对于模型需要多深才能被称为“深度”，目前也没有达成共识。然而，深度学习可以安全地视为研究比传统机器学习涉及更多学习函数或学习概念组合的模型。

总而言之，本书的主题是深度学习，是一种人工智能方法。

具体来说，它是一种机器学习，一种使计算机系统能够通过经验和数据进行改进的技术。我们认为，机器学习是构建可以在复杂的现实世界环境中运行的人工智能系统的唯一可行方法。深度学习是一种特殊的机器学习，它通过将世界表示为概念的嵌套层次结构来实现强大的功能和灵活性，每个概念都相对于更简单的概念进行定义，并且更抽象的表示是根据不太抽象的概念计算出来的。图 1.4 说明了这些不同人工智能学科之间的关系。图 1.5 给出了每个工作原理的高级示意图。

## 1.1 谁应该读这本书？

本书适合各类读者，但我们在撰写时主要针对两类目标读者。一类目标读者是正在学习机器学习的大学生（本科生或研究生），包括那些刚开始从事深度学习和人工智能研究的人士。另一类目标读者是……

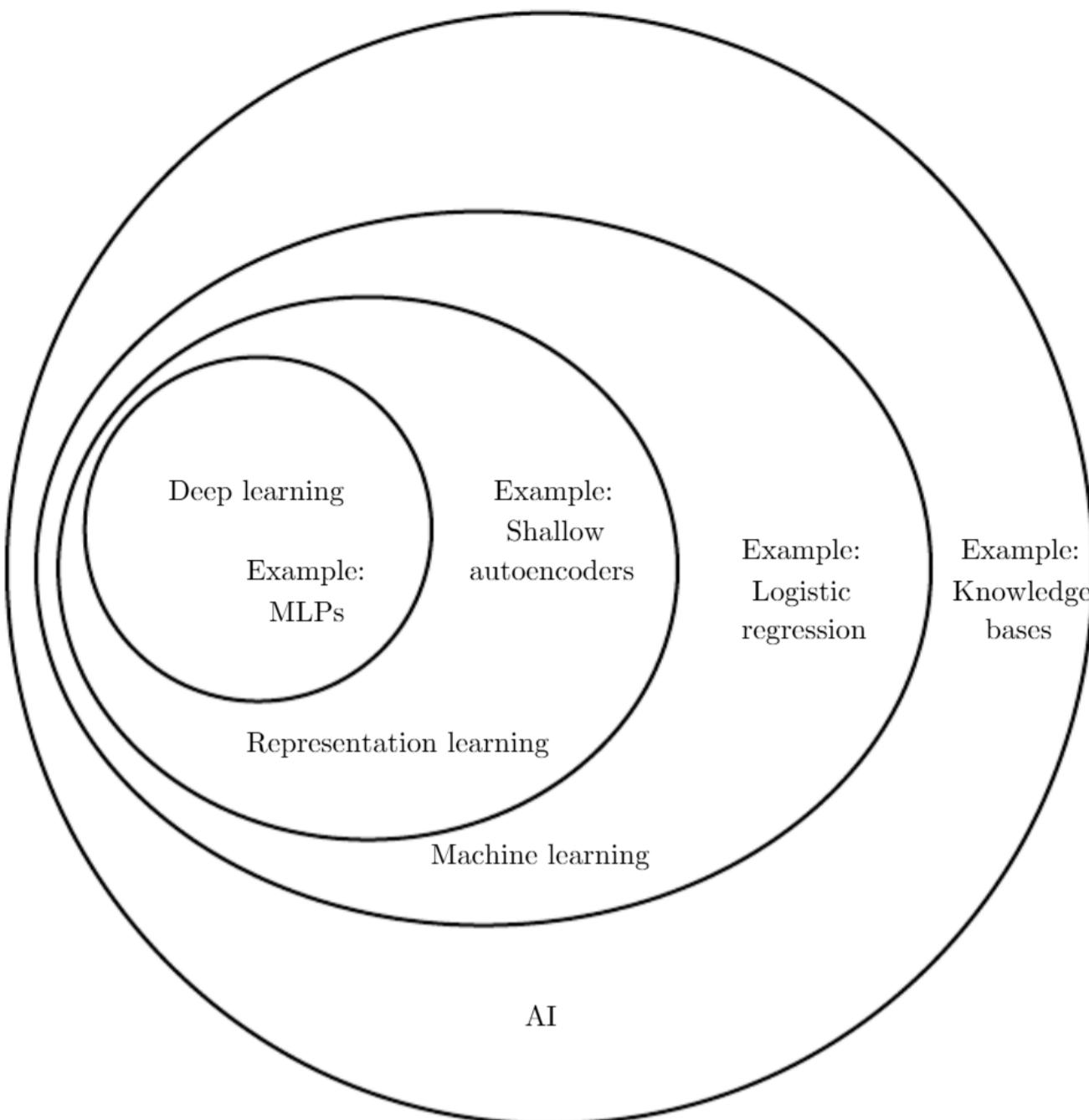


Figure 1.4: A Venn diagram showing how deep learning is a kind of representation learning, which is in turn a kind of machine learning, which is used for many but not all approaches to AI. Each section of the Venn diagram includes an example of an AI technology.

target audience is software engineers who do not have a machine learning or statistics background but want to rapidly acquire one and begin using deep learning in their product or platform. Deep learning has already proved useful in many software disciplines, including computer vision, speech and audio processing, natural language processing, robotics, bioinformatics and chemistry, videogames, search engines, online advertising and finance.

This book has been organized into three parts to best accommodate a variety of readers. Part I introduces basic mathematical tools and machine learning concepts. Part II describes the most established deep learning algorithms, which are essentially solved technologies. Part III describes more speculative ideas that are widely believed to be important for future research in deep learning.

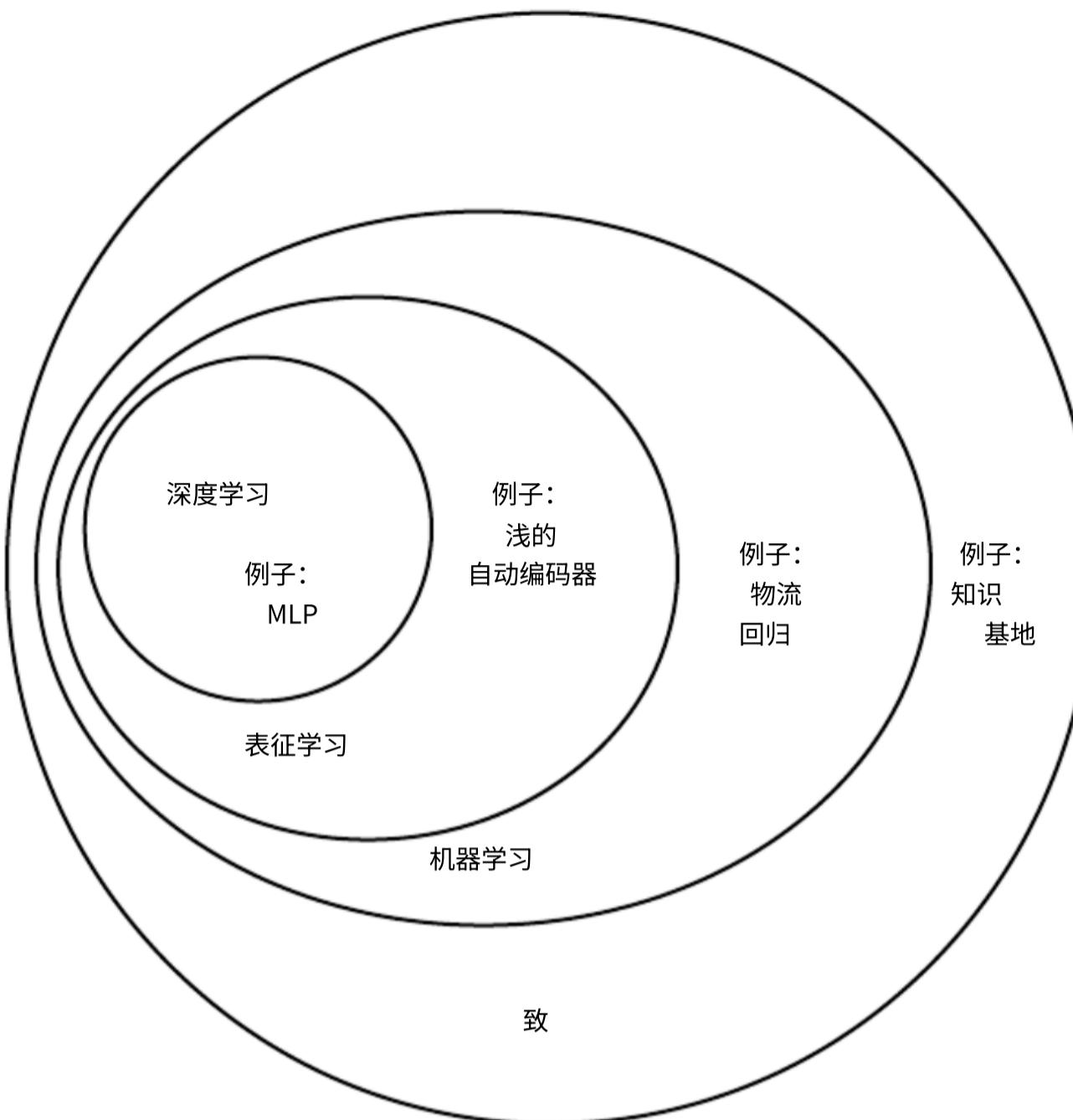


图 1.4：维恩图展示了深度学习是一种表征学习，而表征学习又是一种机器学习，用于许多但并非所有人工智能方法。维恩图的每一部分都包含一个人工智能技术的示例。

目标受众是那些没有机器学习或统计学背景，但想要快速掌握这些背景并开始在其产品或平台中使用深度学习的软件工程师。深度学习已经被证明在许多软件学科中有用，包括计算机视觉、语音和音频处理、自然语言处理、机器人技术、生物信息学和化学、视频游戏、搜索引擎、在线广告和金融。

本书分为三部分，以最好地适应各种读者。第一部分介绍基本的数学工具和机器学习概念。第二部分描述最成熟的深度学习算法，这些算法本质上是已解决的技术。第三部分描述更多推测性的想法，人们普遍认为这些想法对未来的深度学习研究很重要。

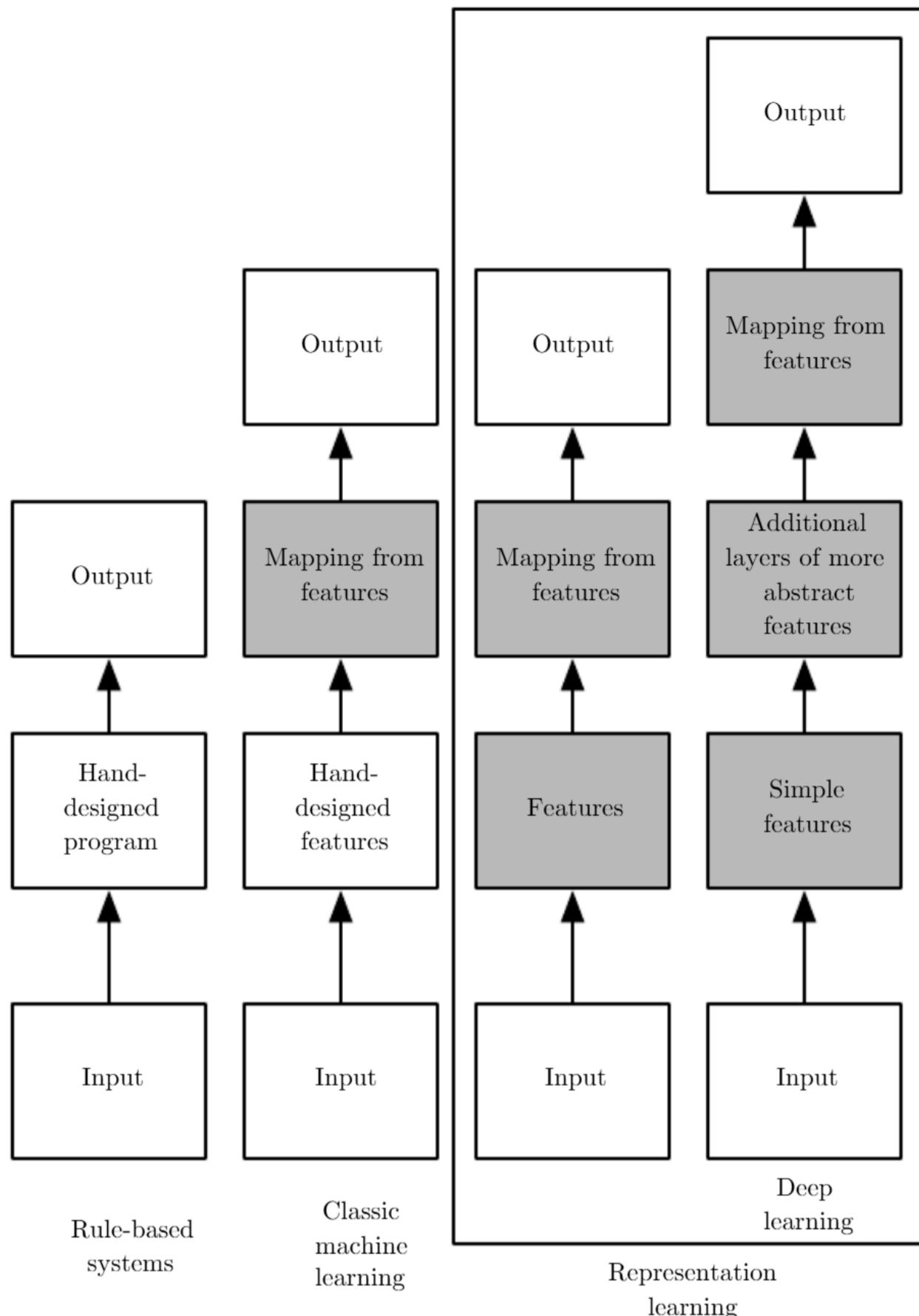


Figure1.5: Flowcharts showing how the different parts of an AI system relate to each other within different AI disciplines. Shaded boxes indicate components that are able to learn from data.

Readers should feel free to skip parts that are not relevant given their interests or background. Readers familiar with linear algebra, probability, and fundamental machine learning concepts can skip part **I**, for example, while those who just want to implement a working system need not read beyond part **II**. To help choose which

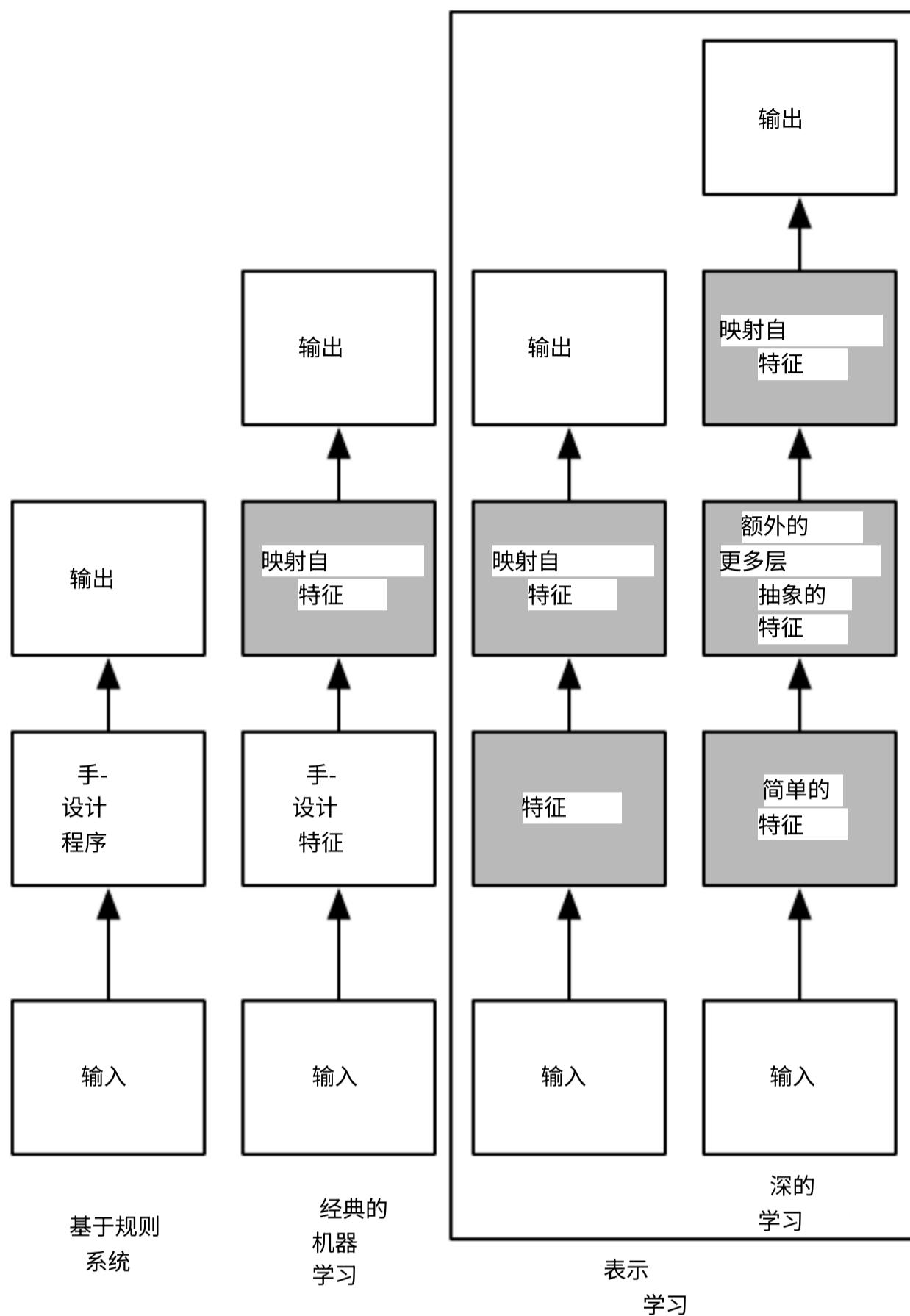


图 1.5：流程图展示了人工智能系统的不同部分在不同的人工智能学科中如何相互关联。阴影框表示能够从数据中学习的组件。

读者可以根据自己的兴趣或背景跳过不相关的部分。例如，熟悉线性代数、概率和基本机器学习概念的读者可以跳过此部分，而只想实现一个工作系统的读者则无需阅读此部分之后的内容。

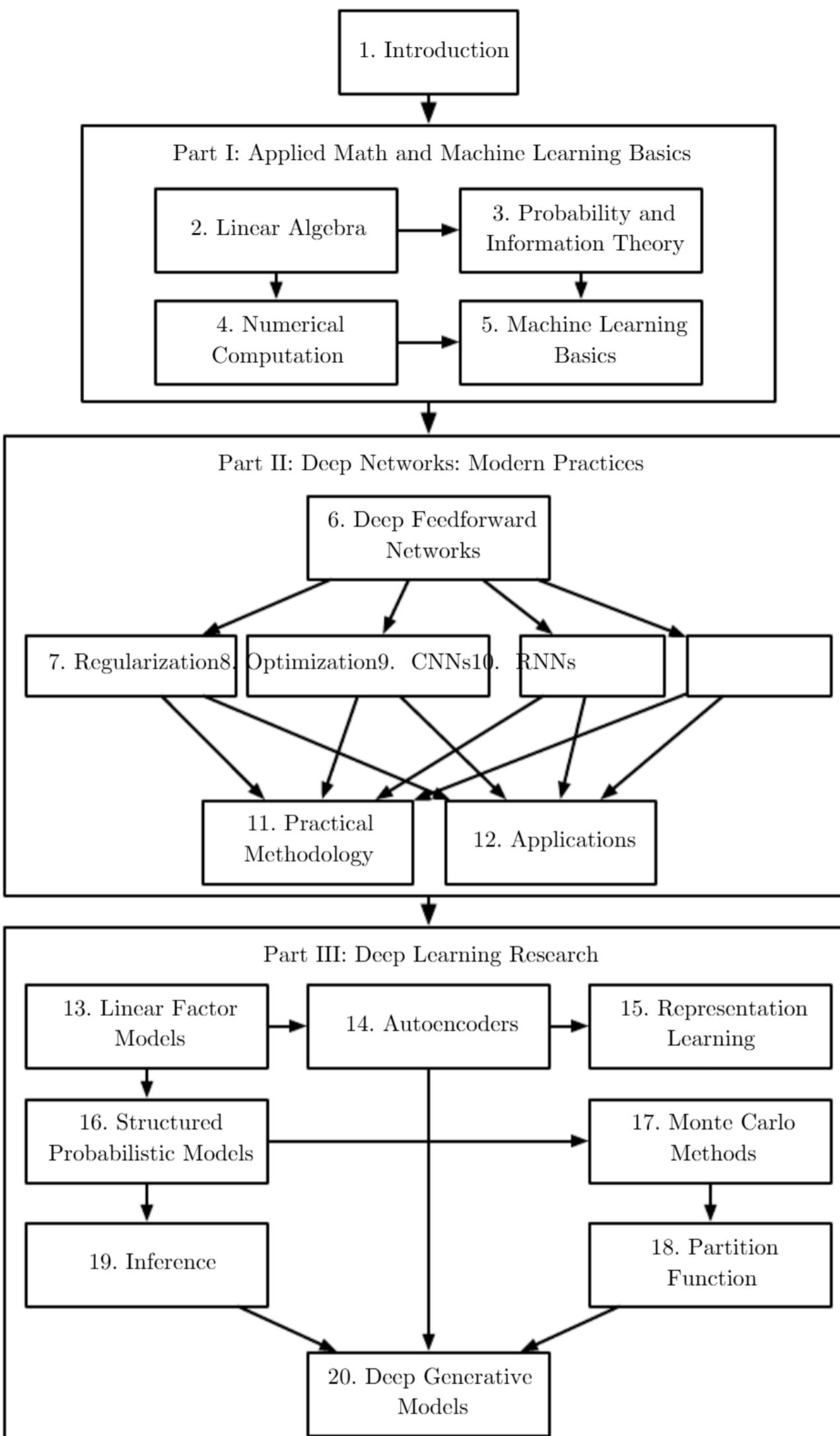


Figure1.6: Thehigh-levelorganizationofthebook. Anarrowfromonechaptertoanother indicatesthattheformerchapterisprerequisiteitematerialforunderstandingthelatter.

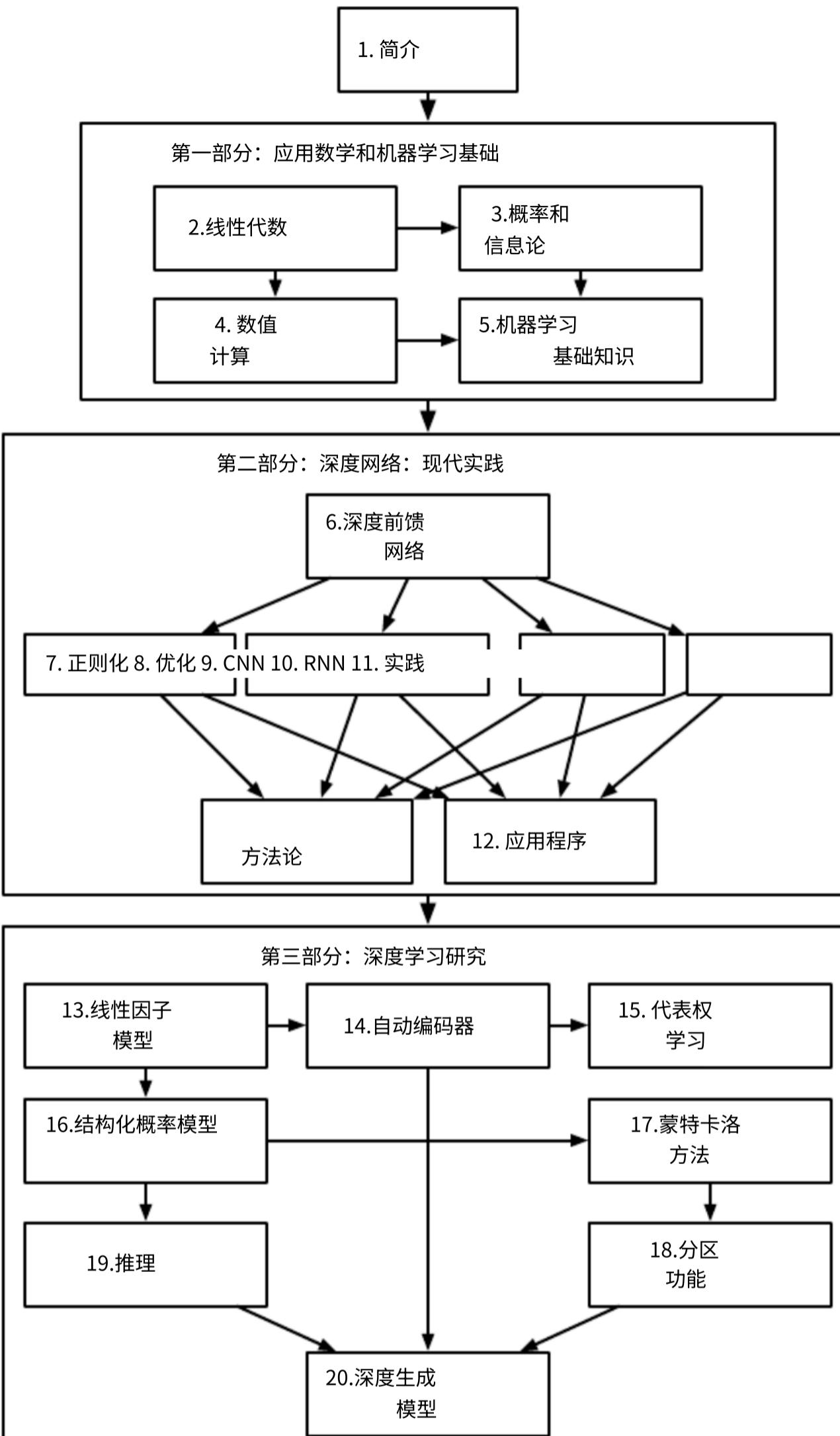


图 1.6：本书的高层组织结构。从一章到另一章的箭头表示前一章是理解后一章的先决条件材料。

chapter to read, figure 1.6 provides a flowchart showing the high-level organization of the book.

We do assume that all readers come from a computer science background. We assume familiarity with programming, a basic understanding of computational performance issues, complexity theory, introductory level calculus and some of the terminology of graph theory.

## 1.2 Historical Trends in Deep Learning

It is easiest to understand deep learning with some historical context. Rather than providing a detailed history of deep learning, we identify a few key trends:

- Deep learning has had a long and rich history, but has gone by many names, reflecting different philosophical viewpoints, and has waxed and waned in popularity.
- Deep learning has become more useful as the amount of available training data has increased.
- Deep learning models have grown in size over time as computer infrastructure (both hardware and software) for deep learning has improved.
- Deep learning has solved increasingly complicated applications with increasing accuracy over time.

### 1.2.1 The Many Names and Changing Fortunes of Neural Networks

We expect that many readers of this book have heard of deep learning as an exciting new technology, and are surprised to see a mention of ‘history’ in a book about an emerging field. In fact, deep learning dates back to the 1940s. Deep learning only appears to be new, because it was relatively unpopular for several years preceding its current popularity, and because it has gone through many different names, only recently being called “deep learning.” The field has been rebranded many times, reflecting the influence of different researchers and different perspectives.

A comprehensive history of deep learning is beyond the scope of this textbook. Some basic context, however, is useful for understanding deep learning. Broadly speaking, there have been three waves of development: deep learning known as **cybernetics** in the 1940s–1960s, deep learning known as **connectionism** in the

章节阅读，图提供了一个流程图，展示了本书的高级组织结构。假设所有读者都来自计算机科学背景。我们假设他们熟悉编程、对计算性能问题、复杂性理论、入门级微积分和一些图论术语有基本的了解。

## 1.2 深度学习的历史趋势

通过一些历史背景来理解深度学习是最容易的。我们不会提供深度学习的详细历史，而是确定一些关键趋势：

- 深度学习有着悠久而丰富的历史，但其名称众多，反映了不同的哲学观点，而且受欢迎程度时高时低。
- 随着可用训练数据量的增加，深度学习变得更加有用。
- 随着深度学习的计算机基础设施（硬件和软件）的改进，深度学习模型的规模不断扩大。
- 随着时间的推移，深度学习已经解决了越来越复杂的应用问题，并且准确性也越来越高。

### 1.2.1 神经网络的众多名称和变迁

我们预计本书的许多读者都听说过深度学习这一令人兴奋的新技术，并且很惊讶地看到在关于新兴领域的书中提到“历史”。事实上，深度学习可以追溯到 20 世纪 40 年代。深度学习只是看起来很新，因为在它现在流行之前的几年里它相对不受欢迎，并且因为它经历了许多不同的名称，直到最近才被称为“深度学习”。该领域已被多次重新命名，反映了不同研究人员和不同观点的影响。

深度学习的全面历史超出了本教科书的范围。然而，了解一些基本背景对理解深度学习还是有用的。广义上讲，深度学习经历了三次发展浪潮：20 世纪 40 年代至 60 年代被称为控制论的深度学习，20 世纪 60 年代被称为  
联结主义

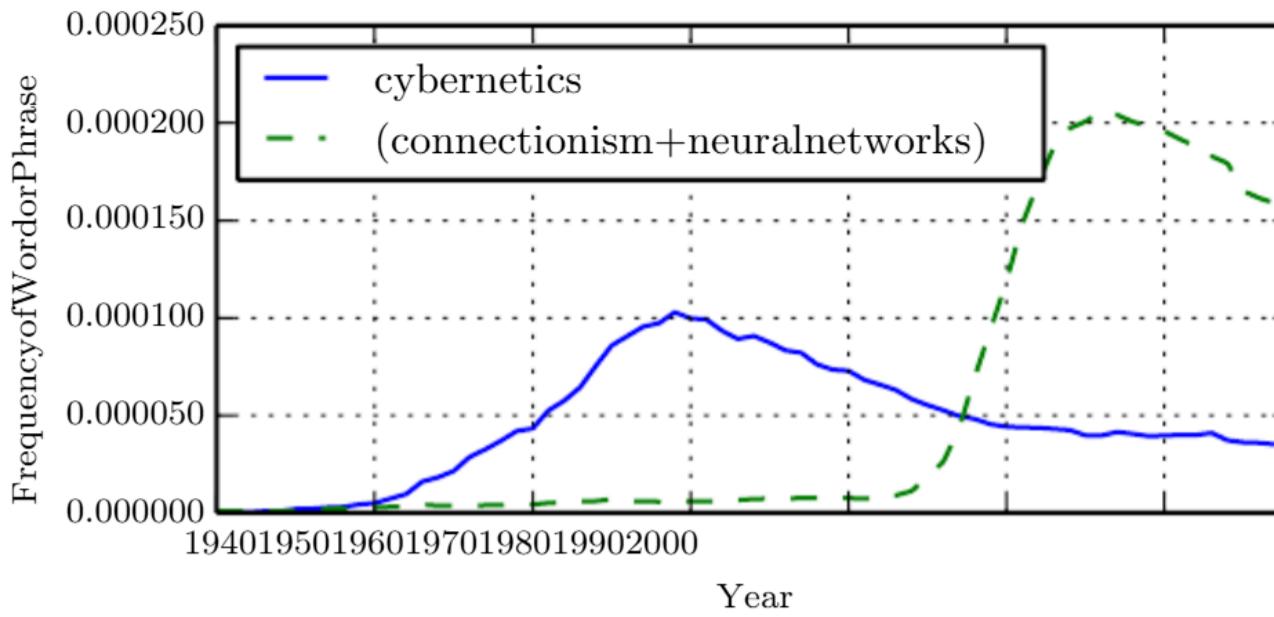


Figure 1.7: Two of the three historical waves of artificial neural nets research, as measured by the frequency of the phrases “cybernetics” and “connectionism” or “neural networks,” according to Google Books (the third wave is too recent to appear). The first wave started with cybernetics in the 1940s–1960s, with the development of theories of biological learning (McCulloch and Pitts , 1943 Hebb; , 1949) and implementations of the first models, such as the perceptron (Rosenblatt, 1958), enabling the training of a single neuron. The second wave started with the connectionist approach of the 1980–1995 period, with back-propagation (Rumelhart *et al.* , 1986a) to train a neural network with one or two hidden layers. The current and third wave, deep learning, started around 2006 (Hinton *et al.* , 2006 Bengio; *et al.* , 2007 Ranzato; *et al.* , 2007a) and is just now appearing in book form as of 2016. The other two waves similarly appeared in book form much later than the corresponding scientific activity occurred.

1980s–1990s, and the current resurgence under the name deep learning beginning in 2006. This is quantitatively illustrated in figure 1.7.

Some of the earliest learning algorithms were recognized today were intended to be computational models of biological learning, that is, models of how learning happens or could happen in the brain. As a result, one of the names that deep learning has gone by is **artificial neural networks** (ANNs). The corresponding perspective on deep learning models is that they are engineered systems inspired by the biological brain (whether the human brain or the brain of another animal). While the kinds of neural networks used for machine learning have sometimes been used to understand brain function (Hinton and Shallice , 1991), they are generally not designed to be realistic models of biological function. The neural perspective on deep learning is motivated by two main ideas. One idea is that the brain provides a proof by example that intelligent behavior is possible, and a conceptually straightforward path to building intelligence is to reverse engineer the computational principles behind the brain and duplicate its functionality. Another

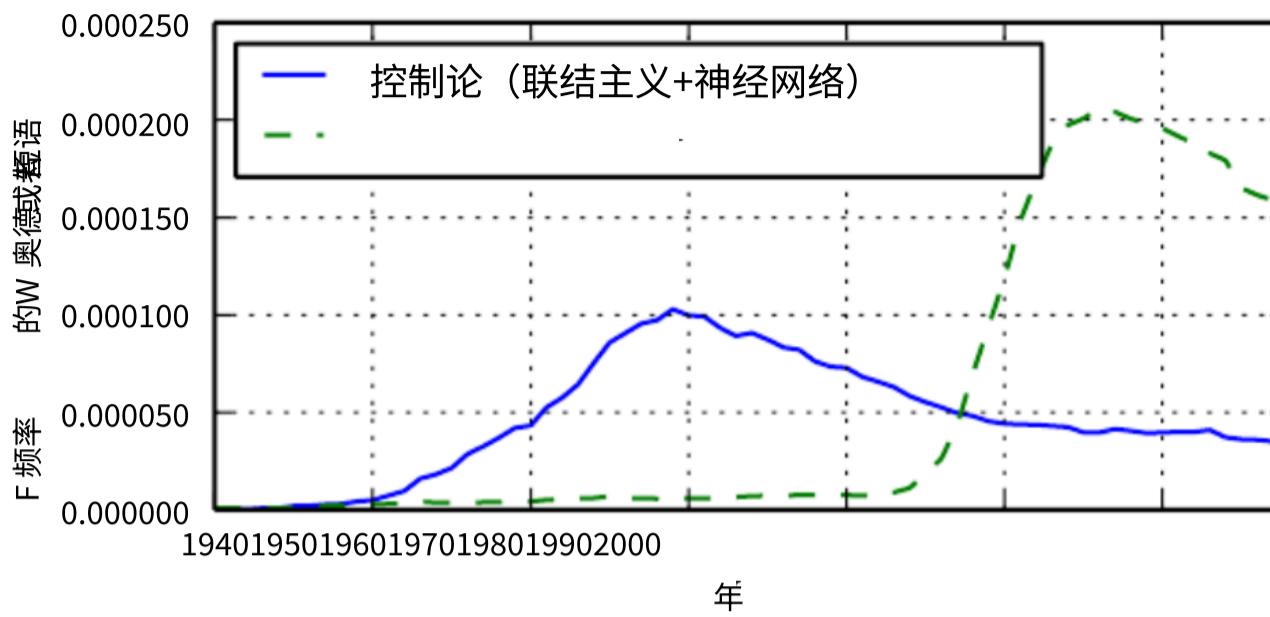


图 1.7：根据 Google 图书的数据，以“控制论”和“联结主义”或“神经网络”等词语出现的频率来衡量，人工神经网络研究经历了三次历史浪潮中的两次（第三次浪潮出现得较晚）。第一次浪潮始于 20 世纪 40 年代至 60 年代的控制论，当时出现了生物学习理论（McCulloch 和 Pitts 1943; Hebb 1949, 1951; 1952），并实现了第一批模型，如感知器（，），使单个神经元的训练成为可能。第二次浪潮始于 1958 年的 Rosenblatt 浪潮，始于 1980 年至 1995 年期间的联结主义方法，使用反向传播（，）来训练具有两层隐藏层的神经网络（Rumelhart et al. 1986）。当前的第三次浪潮，深度学习，始于 2006 年左右（Hinton 等人）。

2006Bengio 2007Ranzato 2007a; etal.; etal.,) 并且刚刚于 2016 年以书籍形式出现。另外两波浪潮同样以书籍形式出现，但比相应的科学活动发生的时间要晚得多。

深度学习最初在 20 世纪 80 年代至 90 年代兴起，并从 2006 年开始以深度学习的名义重新兴起。图 1.7 定量地说明了这一点。今天公认的一些最早的学习算法旨在成为生物学习的计算模型，即学习如何在大脑中发生或可能发生的模型。因此，深度学习的一个名称是人工神经网络 (ANN)。深度学习模型的相应观点是，它们是受生物大脑（无论是人类大脑还是其他动物的大脑）启发的工程系统。虽然用于机器学习的各种神经网络有时被用来理解大脑功能（，），但它们通常不是 Hinton 和 Shallice 1991 设计为生物功能的现实模型。深度学习的神经视角受到两个主要思想的推动。一个思想是，大脑通过例子证明了智能行为是可能的，而构建智能的概念上直接的途径是逆向工程大脑背后的计算原理并复制其功能。另一个

perspective is that it would be deeply interesting to understand the brain and the principles that underlie human intelligence, so machine learning models that shed light on these basic scientific questions are useful apart from their ability to solve engineering applications.

The modern term “deep learning” goes beyond the neuroscientific perspective on the current breed of machine learning models. It appeals to a more general principle of learning *multiple levels of composition*, which can be applied in machine learning frameworks that are not necessarily neurally inspired.

The earliest predecessors of modern deep learning were simple linear models motivated from a neuroscientific perspective. These models were designed to take a set of  $n$  input values  $x_1, \dots, x_n$  and associate them with an output  $y$ . These models would learn a set of weights  $w_1, \dots, w_n$  and compute their output  $f(\mathbf{x}, \mathbf{w}) = x_1 w_1 + \dots + x_n w_n$ . This first wave of neural networks research was known as cybernetics, as illustrated in figure 1.7.

The McCulloch-Pitts neuron ([McCulloch and Pitts, 1943](#)) was a nearly model of brain function. This linear model could recognize two different categories of inputs by testing whether  $f(\mathbf{x}, \mathbf{w})$  is positive or negative. Of course, for the model to correspond to the desired definition of the categories, the weights needed to be set correctly. These weights could be set by the human operator. In the 1950s, the perceptron ([Rosenblatt, 1958, 1962](#)) became the first model that could learn the weights that defined the categories given examples of inputs from each category. The **adaptive linear element** (ADALINE), which dates from about the same time, simply returned the value of  $f(\mathbf{x})$  itself to predict a real number ([Widrow and Hoff, 1960](#)) and could also learn to predict these numbers from data.

These simple learning algorithms greatly affected the modern landscape of machine learning. The training algorithm used to adapt the weights of the ADALINE was a special case of an algorithm called **stochastic gradient descent**. Slightly modified versions of the stochastic gradient descent algorithm remain the dominant training algorithms for deep learning models today.

Models based on the  $f(\mathbf{x}, \mathbf{w})$  used by the perceptron and ADALINE are recalled **linear models**. These models remain some of the most widely used machine learning models, though in many cases they are *trained* in different ways than the original models were retrained.

Linear models have many limitations. Most famously, they cannot learn the XOR function, where  $f([0, 1], \mathbf{w}) = 1$  and  $f([1, 0], \mathbf{w}) = 1$  but  $f([1, 1], \mathbf{w}) = 0$  and  $f([0, 0], \mathbf{w}) = 0$ . Critics who observed these flaws in linear models caused a backlash against biologically inspired learning in general ([Minsky and Papert, 1969](#)). This was the first major dip in the popularity of neural networks.

我的观点是，了解大脑和人类智能背后的原理将会非常有趣，因此，阐明这些基本科学问题的机器学习模型除了能够解决工程应用之外，还非常有用。

现代术语“深度学习”超越了神经科学对当前机器学习模型的视角。它诉诸于学习多层次组合的更普遍的原理，该原理可应用于不一定受神经启发的机器学习框架中。

现代深度学习的早期前身是从神经科学角度出发的简单线性模型。这些模型被设计用来接受一组输入值

$$x_1, \dots, x_n \text{ and } w \text{ 将它们与输出关联起来} \quad y.$$

深度学习模型的相应观点是，它们是受生物大脑启发而设计的系统（无论是人类大脑还是其他动物的大脑），这些模型将学习权重 第一波神经网络研究被称为控制论，如图 1.7 所示。McCulloch-Pitts 神经元（）是大脑功能的早期模型。该线性模型通过测试  $f(xw,)$  是正数还是负数，可以识别两种不同的输入类别。当然，为了使模型与所需的类别定义相对应，需要正确设置权重。这些权重可以由人类操作员设置。20 世纪 50 年代，感知器（Rosenblatt 1958 1962,）成为第一个可以通过给定每个类别的输入示例来学习定义类别的权重的模型。

自适应线性元素（ADALINE）大约在同一时期出现，它只是返回  $f(x)$  本身的值来预测一个实数（Widrow 和 Hoff 1960,），并且还可以学习从数据中预测这些数字。

这些简单的学习算法极大地影响了机器学习的现代格局。用于调整 ADALINE 权重的训练算法是称为随机梯度下降的算法的一个特例。稍微修改后的随机梯度下降算法版本仍然是当今深度学习模型的主要训练算法。

感知器和 ADALINE 使用的基于  $f(xw,)$  的模型称为线性模型。这些模型仍然是一些最广泛使用的机器学习模型，尽管在许多情况下，它们的训练方式与原始模型不同。

线性模型有许多局限性。最著名的是，它们不能学习 XOR 函数，其中  $f([0,1], w) = 1$  和  $f([1,0], w) = 1$ ，但  $f([1,1], w) = 0$  和  $f([0,0], w) = 0$ 。批评者观察到线性模型中的这些缺陷，引起了对生物启发式学习的强烈反对（Minsky and Papert, 1969）。这是神经网络流行的第一个主要原因。

Today, neuroscience is regarded as an important source of inspiration for deep learning researchers, but it is no longer the predominant guide for the field.

The main reason for the diminished role of neuroscience in deep learning research today is that we simply do not have enough information about the brain to use it as a guide. To obtain a deep understanding of the actual algorithms used by the brain, we would need to be able to monitor the activity of (at the very least) thousands of interconnected neurons simultaneously. Because we are not able to do this, we are far from understanding even some of the most simple and well-studied parts of the brain (Olshausen and Field , 2005).

Neuroscience has given us reasons to hope that a single deep learning algorithm can solve many different tasks. Neuroscientists have found that ferrets can learn to “see” with the auditory processing region of their brain if their brains are rewired to send visual signals to that area (Von Melchner *et al.* , 2000). This suggests that much of the mammalian brain might use a single algorithm to solve most of the different tasks that the brain solves. Before this hypothesis, machine learning research was more fragmented, with different communities of researchers studying natural language processing, vision, motion planning and speech recognition. Today, these application communities are still separate, but it is common for deep learning research groups to study many or even all these applications simultaneously.

We are able to draw some rough guidelines from neuroscience. The basic idea of having many computational units that become intelligent only via their interactions with each other is inspired by the brain. The neocognitron (Fukushima, 1980) introduced a powerful model architecture for processing images that was inspired by the structure of the mammalian visual system and later became the basis for the modern convolutional network (LeCun *et al.* , 1998b), as we will see in section 9.10. Most neural networks today are based on a model neuron called the **rectified linear unit**. The original cognitron (Fukushima, 1975) introduced a more complicated version that was highly inspired by our knowledge of brain function. The simplified modern version was developed incorporating ideas from many viewpoints, with Nair and Hinton (2010) and Glorot *et al.* (2011a) citing neuroscience as an influence, and Jarrett *et al.* (2009) citing more engineering-oriented influences. While neuroscience is an important source of inspiration, it need not be taken as a rigid guide. We know that actual neurons compute very different functions than modern rectified linear units, but greater neural realism has not yet led to an improvement in machine learning performance. Also, while neuroscience has successfully inspired several neural network architectures we do not yet know enough about biological learning for neuroscience to offer much guidance for the *learning algorithms* we use to train these architectures.

如今，神经科学被视为深度学习研究人员的重要灵感来源，但它不再是该领域的主要指南。

如今，神经科学在深度学习研究中的作用减弱的主要原因是我们在有足够的大脑信息来指导它。为了深入了解大脑使用的实际算法，我们需要能够同时监控（至少）数千个相互连接的神经元的活动。因为我们无法做到这一点，所以我们甚至无法理解大脑中一些最简单、研究最透彻的部分（，）。Olshausen and Field  
2005

神经科学让我们有理由希望单一的深度学习算法可以解决许多不同的任务。神经科学家发现，如果雪貂的大脑重新连接以向该区域发送视觉信号，它们就可以学会用大脑的听觉处理区域“看”（Von Melchner 2000 et al.，）。这表明大部分哺乳动物的大脑可能使用单一算法来解决大脑解决的大多数不同任务。在此假设之前，机器学习研究更加分散，不同的研究人员社区研究自然语言处理，视觉，运动规划和语音识别。今天，这些应用社区仍然是分开的，但深度学习研究小组同时研究许多甚至所有这些应用领域是很常见的。

我们可以从神经科学中得出一些粗略的指导方针。许多计算单元仅通过相互作用才能变得智能的基本思想是受到大脑的启发。新认知机（Fukushima, 1980）引入了一种强大的图像处理模型架构，其灵感来自哺乳动物视觉系统的结构，后来成为现代卷积网络（，）的基础，正如我们将在 1998b 节中看到的 LeCun 等人。当今大多数神经网络都基于称为 9.10 的模型神经元

整流线性单元。原始认知机（Fukushima 1975,）引入了一个更复杂的版本，极大地启发了您对大脑功能的知识。简化的现代版本结合了许多观点的想法而开发，其中（）和（）引用了 Nair and Hinton 2010 Glorot et al. 2011a 神经科学作为影响，并且（）引用了更多工程 - Jarrett et al. 2009 导向的影响。虽然神经科学是重要的灵感来源，但不必将其作为僵化的指南。我们知道实际的神经元计算的功能与现代整流线性单元完全不同，但更大的神经现实主义尚未导致机器学习性能的改进。此外，虽然神经科学已经成功启发了几种神经网络架构，但我们对生物学习的了解还不足以让神经科学为我们用来训练搜索结构的学习算法提供足够的指导。

Media accounts often emphasize the similarity of deep learning to the brain. While it is true that deep learning researchers are more likely to cite the brain as an influence than researchers working in other machine learning fields, such as kernel machines or Bayesian statistics, one should not view deep learning as an attempt to simulate the brain. Modern deep learning draws inspiration from many fields, especially applied math fundamentals like linear algebra, probability, information theory, and numerical optimization. While some deep learning researchers cite neuroscience as an important source of inspiration, others are not concerned with neuroscience at all.

It is worth noting that the effort to understand how the brain works on an algorithmic level is alive and well. This endeavor is primarily known as “computational neuroscience” and is a separate field of study from deep learning. It is common for researchers to move back and forth between both fields. The field of deep learning is primarily concerned with how to build computer systems that are able to successfully solve tasks requiring intelligence, while the field of computational neuroscience is primarily concerned with building more accurate models of how the brain actually works.

In the 1980s, the second wave of neural network research emerged, gaining great part via a movement called **connectionism**, or **parallel distributed processing** (Rumelhart *et al.*, 1986c; McClelland *et al.*, 1995). Connectionism arose in the context of cognitive science. Cognitive science is an interdisciplinary approach to understanding the mind, combining multiple different levels of analysis. During the early 1980s, most cognitive scientists studied models of symbolic reasoning. Despite their popularity, symbolic models were difficult to explain in terms of how the brain could actually implement them using neurons. The connectionists began to study models of cognition that could actually be grounded in neural implementations (Touretzky and Minton, 1985), reviving many ideas dating back to the work of psychologist Donald Hebb in the 1940s (Hebb, 1949).

The central idea in connectionism is that a large number of simple computational units can achieve intelligent behavior when networked together. This insight applies equally to neurons in biological nervous systems as it does to hidden units in computational models.

Several key concepts arose during the connectionism movement of the 1980s that remain central to today’s deep learning.

One of these concepts is that of **distributed representation** (Hinton *et al.*, 1986). This is the idea that each input to a system should be represented by many features, and each feature should be involved in the representation of many possible inputs. For example, suppose we have a vision system that can recognize

媒体报道常常强调深度学习与大脑的相似性。虽然深度学习研究人员比其他机器学习领域（例如核机器或贝叶斯统计）的研究人员更有可能引用大脑作为参考，但不应将深度学习视为模拟大脑的尝试。现代深度学习从许多领域汲取灵感，

尤其是应用数学基础知识，如线性代数、概率、信息论和数值优化。虽然一些深度学习研究人员将神经科学视为重要的灵感来源，但其他人根本不关心神经科学。

值得注意的是，理解大脑在算法层面如何运作的努力依然活跃。这项努力主要被称为“计算神经科学”，是一个独立于深度学习的研究领域。研究人员在两个领域之间来回切换是很常见的。深度学习领域主要关注如何构建能够成功解决需要智能的任务的计算机系统，而计算神经科学领域主要关注如何构建更精确的大脑实际运作模型。

20世纪80年代，神经网络研究的第二波浪潮兴起，很大程度上是通过一场被称为联结主义或并行分布式处理的运动（，；，）出现的。联结主义出现于 Rumelhart 等人 1986c McClelland 等人 1995 认知科学的背景下。认知科学是一种理解思维的跨学科方法，结合了多个不同层次的分析。在 20 世纪 80 年代早期，大多数认知科学家研究符号推理模型。尽管符号模型很受欢迎，但它们很难解释大脑如何真正实现那些令人沉思的神经元。联结主义者开始研究能够真正建立在神经实现基础上的认知模型（Touretzky and Minton 1985,），复兴了许多可以追溯到心理学家唐纳德·赫布在 20 世纪 40 年代的工作的想法（，）。赫布 1949

联结主义的核心思想是，大量简单的计算单元在联网时可以实现智能行为。这一见解同样适用于生物神经系统中的神经元，就像它适用于计算模型中的隐藏单元一样。

20世纪80年代的联结主义运动中出现了几个关键概念，这些概念至今仍是深度学习的核心。

其中一个概念是分布式表示（Hinton 等人，1986）。

cars, trucks, and birds, and these objects can each be red, green, or blue. One way of representing these inputs would be to have a separate neuron or hidden unit that activates for each of the nine possible combinations: red truck, red car, red bird, green truck, and so on. This requires nine different neurons, and each neuron must independently learn the concept of color and object identity. One way to improve on this situation is to use a distributed representation, with three neurons describing the color and three neurons describing the object identity. This requires only six neurons total instead of nine, and the neuron describing redness is able to learn about redness from images of cars, trucks and birds, not just from images of one specific category of objects. The concept of distributed representation is central to this book and is described in greater detail in chapter 15.

Another major accomplishment of the connectionist movement was the successful use of back-propagation to train deep neural networks with internal representations and the popularization of the back-propagation algorithm (Rumelhart *et al.*, 1986a; LeCun, 1987). This algorithm has waxed and waned in popularity but, as of this writing, is the dominant approach to training deep models.

During the 1990s, researchers made important advances in modeling sequences with neural networks. Hochreiter (1991) and Bengio *et al.* (1994) identified some of the fundamental mathematical difficulties in modeling long sequences, described in section 10.7. Hochreiter and Schmidhuber 1997 () introduced the long short-term memory (LSTM) network to resolve some of these difficulties. Today, the LSTM is widely used for many sequence modeling tasks, including many natural language processing tasks at Google.

The second wave of neural networks research lasted until the mid-1990s. Ventures based on neural networks and other AI technologies began to make unrealistically ambitious claims while seeking investments. When AI research did not fulfill these unreasonable expectations, investors were disappointed. Simultaneously, other fields of machine learning made advances. Kernel machines (Boser *et al.*, 1992 Cortes and Vapnik, 1995 Schölkopf *et al.*, 1999) and graphical models (Jordan, 1998) both achieved good results on many important tasks. These two factors led to a decline in the popularity of neural networks that lasted until 2007.

During this time, neural networks continued to obtain impressive performance on some tasks (LeCun *et al.*, 1998b; Bengio *et al.*, 2001). The Canadian Institute for Advanced Research (CIFAR) helped to keep neural networks research alive via its Neural Computation and Adaptive Perception (NCAP) research initiative. This program united machine learning research groups led by Geoffrey Hinton at University of Toronto, Yoshua Bengio at University of Montreal, and Yann LeCun at New York University. The multidisciplinary CIFAR NCAP research initiative

汽车、卡车和鸟类，这些物体可以是红色、绿色或蓝色。表示这些输入的一种方法是拥有一个单独的神经元或隐藏单元，可以激活九种可能的组合中的每一种：红色卡车、红色汽车、红色鸟、绿色卡车等等。这需要九个不同的神经元，每个神经元都必须独立学习颜色和物体身份的概念。改善这种情况的一种方法是使用分布式表示，其中三个神经元描述颜色，三个神经元描述物体身份。这只需要六个神经元而不是九个，并且描述红色的神经元能够从汽车、卡车和鸟类的图像中学习红色，而不仅仅是从某一特定类别的物体的图像中学习。分布式表示的概念是本书的核心，并将在第 15 章中详细描述。联结主义运动的另一项主要成就是成功使用反向传播来训练具有内部表示的深度神经网络，以及反向传播算法的普及 (Rumelhart 等人, ;, )。这种算法在 1986 年 LeCun 1987 年流行起来并逐渐衰落，但截至撰写本文时，它仍是训练深度模型的主要方法。

在 20 世纪 90 年代，研究人员在利用神经网络对序列进行建模方面取得了重要进展。() 和() 发现了长序列建模的一些基本数学难题，如第 10.7 节所述。Hochreiter 和 Schmidhuber1997() 引入了长短期记忆 (LSTM) 网络来解决其中的一些难题。如今，LSTM 被广泛用于许多序列建模任务，包括谷歌的许多自然语言处理任务。

神经网络研究的第二波浪潮一直持续到 20 世纪 90 年代中期。基于神经网络和其他人工智能技术的企业在寻求投资时开始做出不切实际的雄心勃勃的主张。当人工智能研究未能满足这些不合理的期望时，投资者感到失望。与此同时，机器学习的其他领域取得了进展。内核机器 (Boser 等人, 1992; CortesandVapnik, 1995; Schölkopf, 1999; Jor-; 等人,) 和图形模型 (dan 1998,) 都在许多重要任务上取得了良好的结果。这两个因素导致神经网络的流行度下降，直到 2007 年才有所下降。

在此期间，神经网络在某些任务上继续取得令人印象深刻的表现 (, ; , )。加拿大高等研究院 (CIFAR) 通过其神经计算与自适应感知 (NCAP) 研究计划，帮助神经网络研究保持活力。该计划联合了由多伦多大学的 Geoffrey Hinton、蒙特利尔大学的 Yoshua Bengio 和纽约大学的 Yann LeCun 领导的机器学习研究小组。这项多学科的 CIFAR 研究计划

also included neuroscientists and experts in human and computer vision.

At this point, deep networks were generally believed to be very difficult to train. We now know that algorithms that have existed since the 1980s work quite well, but this was not apparent circa 2006. The issue is perhaps simply that these algorithms were too computationally costly to allow much experimentation with the hardware available at the time.

The third wave of neural networks research began with a breakthrough in 2006. Geoffrey Hinton showed that a kind of neural network called a deep belief network could be efficiently trained using a strategy called greedy layer-wise pretraining (Hinton *et al.*, 2006), which we describe in more detail in section 15.1. The other CIFAR-affiliated research groups quickly showed that the same strategy could be used to train many other kinds of deep networks (Bengio *et al.*, 2007; Ranzato *et al.*, 2007a) and systematically helped to improve generalization on test examples. This wave of neural networks research popularized the use of the term “deep learning” to emphasize that researchers were now able to train deeper neural networks than had been possible before, and to focus attention on the theoretical importance of depth (Bengio and LeCun, 2007; Delalleau and Bengio, 2011; Pascanu; *et al.*, 2014a; Montufar; *et al.*, 2014). At this time, deep neural networks outperformed competing AI systems based on other machine learning technologies as well as hand-designed functionality. This third wave of popularity of neural networks continues to the time of this writing, though the focus of deep learning research has changed dramatically within the time of this wave. The third wave began with a focus on new unsupervised learning techniques and the ability of deep models to generalize well from small datasets, but today there is more interest in much older supervised learning algorithms and the ability of deep models to learn large labeled datasets.

### 1.2.2 Increasing Dataset Sizes

One may wonder why deep learning has only recently become recognized as a crucial technology even though the first experiments with artificial neural networks were conducted in the 1950s. Deep learning has been successfully used in commercial applications since the 1990s but was often regarded as being more of an art than a technology and something that only an expert could use, until recently. It is true that some skill is required to get good performance from a deep learning algorithm. Fortunately, the amount of skill required reduces as the amount of training data increases. The learning algorithms reaching human performance on complex tasks today are nearly identical to the learning algorithms that struggled to solve toy problems in the 1980s, though the models we train with these algorithms have

还包括神经科学家以及人类和计算机视觉专家。

此时，人们普遍认为深度网络很难训练。我们现在知道，自 20 世纪 80 年代以来存在的算法已经运行得很好，但这在 2006 年左右并不明显。问题可能只是这些算法的计算成本太高，无法利用当时可用的硬件进行大量实验。

第三次神经网络研究浪潮始于 2006 年的一项突破。Geoffrey Hinton 表明，一种称为深度信念网络的神经网络可以通过一种称为贪婪逐层预训练的策略进行有效训练 (, ), 我们将在第 3 节中对此进行更详细的描述。Hinton 等人。2006 15.1 其他 CIFAR 附属研究小组很快证明，同样的策略可以用于训练许多其他类型的深度网络 (,; Bengio 等人 2007; Ranzato 等人 2007,)，并系统地帮助提高测试样本的泛化能力。这波神经网络研究浪潮普及了“深度学习”一词的使用，强调研究人员现在能够训练比以前更深的神经网络，并将注意力集中在

深度的理论重要性 (Bengio and LeCun 2007; Delalleau and Bengio 2011; Pascanu 2014a; Montufar 2014; et al., ; et al.,)。目前，深度神经网络的表现优于基于其他机器学习技术以及手工设计功能的竞争人工智能系统。神经网络的第三次流行浪潮一直持续到本文撰写之时，尽管深度学习研究的重点在这一浪潮中发生了巨大变化。第三次浪潮开始时关注的是新的无监督学习技术和深度模型从小数据集很好地概括的能力，但今天人们对多持股监督学习算法和深度模型利用大型标记数据集的能力更感兴趣。

### 1.2.2 增加数据集大小

人们或许会疑惑，尽管深度学习的首次实验早在 20 世纪 50 年代就已开展，但为什么它直到最近才被公认为一项关键技术。深度学习自 20 世纪 90 年代以来就已成功应用于商业领域，但直到最近，它通常被认为更像是一门艺术而非技术，而且只有专家才能使用。的确，要使深度学习算法获得良好的性能，需要一些技能。幸运的是，随着训练数据量的增加，所需的技能数量会减少。如今，在复杂任务上达到人类水平的学习算法，与 20 世纪 80 年代难以解决玩具问题的学习算法几乎完全相同，尽管我们用这些算法训练的模型已经……

undergone changes that simplify the training of very deep architectures. The most important new development is that today we can provide these algorithms with the resources they need to succeed. Figure 1.8 shows how the size of benchmark datasets has expanded remarkably over time. This trend is driven by the increasing digitization of society. As more and more of our activities take place on computers, more and more of what we do is recorded. As our computers are increasingly networked together, it becomes easier to centralize these records and curate them into a dataset appropriate for machine learning applications. The age of ‘Big Data’

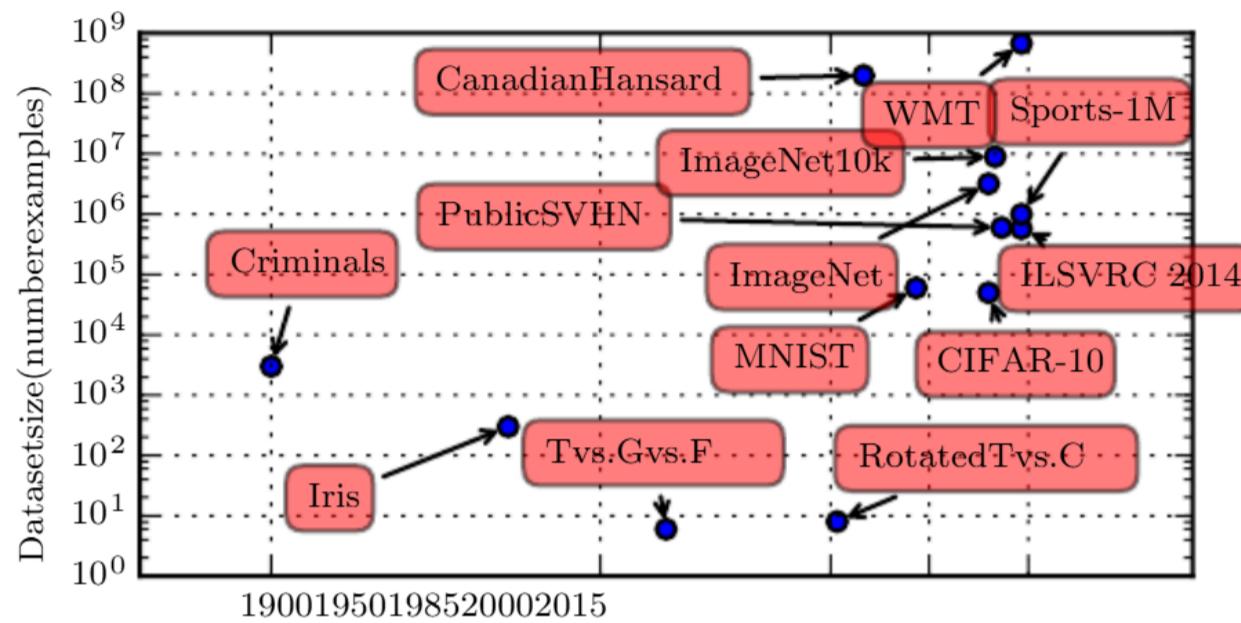


Figure 1.8: Increasing dataset size over time. In the early 1900s, statisticians studied datasets using hundreds or thousands of manually compiled measurements (Garson, 1900; Gosset, 1908; Anderson, 1935; Fisher, 1936). In the 1950s through the 1980s, the pioneers of biologically inspired machine learning often worked with small synthetic datasets, such as low-resolution bitmapsofletters, that were designed to incur low computational cost and demonstrate that neural networks were able to learn specific kinds of functions (Widrow and Hoff, 1960; Rumelhart, et al., 1986b). In the 1980s and 1990s, machine learning became more statistical and began to leverage larger datasets containing tens of thousands of examples, such as the MNIST dataset (shown in figure 1.9) of scans of handwritten numbers (LeCun, et al., 1998b). In the first decade of the 2000s, more sophisticated datasets of this same size, such as the CIFAR-10 dataset (Krizhevsky and Hinton, 2009), continued to be produced. Toward the end of that decade and throughout the first half of the 2010s, significantly larger datasets, containing hundreds of thousands to tens of millions of examples, completely changed what was possible with deep learning. These datasets included the public Street View House Numbers dataset (Netzer, et al., 2011), various versions of the ImageNet dataset (Deng, et al., 2009, 2010a; Russakovsky, et al., 2014a), and the Sports-1M dataset (Karpathy, et al., 2014). At the top of the graph, we see that datasets of translated sentences, such as IBM’s dataset constructed from the Canadian Hansard (Brown, et al., 1990) and the WMT2014 English to French dataset (Schwenk, 2014), are typically far ahead of other dataset sizes.

经历了简化非常深层架构的训练的改变。最重要的新发展是，今天我们可以为这些算法提供它们成功所需的资源。图显示，基准 1.8 数据集的大小随着时间的推移显著扩大。这一趋势是由社会日益数字化所驱动的。随着我们越来越多的活动在计算机上进行，我们所做的也越来越多地被记录下来。随着我们的计算机越来越多地联网在一起，将这些记录集中起来并将它们整理成适用于机器学习应用程序的数据集变得更加容易。“大数据”时代

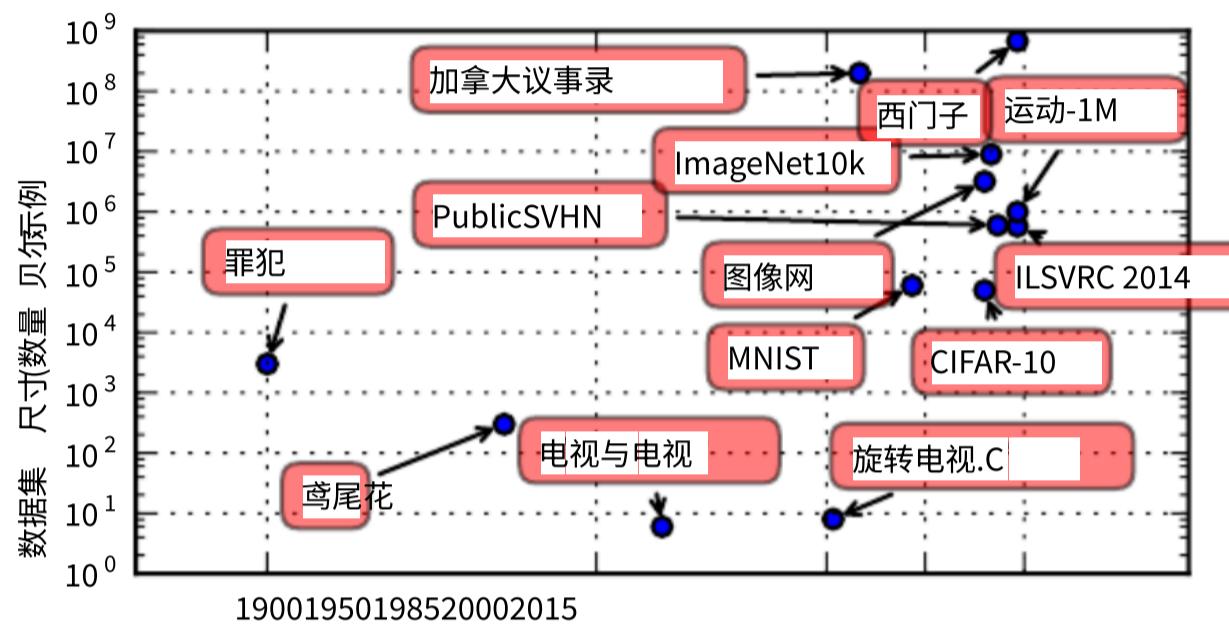


图 1.8：数据集大小随时间推移不断增加。20 世纪初，统计学家使用数百上千个手动编制的测量数据来研究数据集 (Garson 1900、Gosset 1908、Anderson 1935、Fisher 1936, 1991; 2001; 2002; 2003; 2004, 2005)。20 世纪 50 年代到 80 年代，受生物启发的机器学习的先驱者经常使用小型合成数据集，例如低分辨率字母位图，旨在降低计算成本并证明神经网络能够学习特定类型的函数 (Widrow and Hoff 1960、Rumelhart 1986b, 1992; 2003 年)。20 世纪 80 年代和 90 年代，机器学习变得更加具有统计意义，并开始利用包含数万个示例的更大数据集，例如包含 1.9 手写数字扫描图的 MNIST 数据集 (如图所示) (, )。在 21 世纪的第一个十年，同样大小的更复杂的 LeCun et al. 1998b 数据集，例如 CIFAR-10 数据集 (Krizhevsky and Hinton 2009,) 不断涌现。在 21 世纪第一个十年的末期以及整个 21 世纪第一个十年的前半段，包含数亿到数千万个示例的更大的数据集彻底改变了深度学习的可能性。这些数据集包括公共 StreetViewHouseNumbers 数据集 (, ) 以及各种 Netzer et al. 2011 版 ImageNet 数据集 (, , ; Deng et al. 2009 2010a Russakovsky 2014a et al. , ) 和 Sports-1M 数据集 (, )。在图的顶部，我们看到 Karpathy et al. 2014

翻译句子的数据集，例如 IBM 根据加拿大 Hansard 构建的数据集 (Brown 1990 等,) 和 WMT2014 英语到法语数据集 (Schwenk, 2014)，通常远远领先于其他数据集的规模。

8	9	0	1	2	3	4	7	8	9	0	1	2	3	4	5	6	7	8	6
4	2	6	4	7	5	5	4	7	8	9	2	9	3	9	3	8	2	0	5
0	1	0	4	2	6	5	3	5	3	8	0	0	3	4	1	5	3	0	8
3	0	6	2	7	1	1	8	1	7	1	3	8	9	7	6	7	4	1	6
7	5	1	7	1	9	8	0	6	9	4	9	9	3	7	1	9	2	2	5
3	7	8	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	0
1	2	3	4	5	6	7	8	9	8	1	0	5	5	1	9	0	4	1	9
3	8	4	7	7	8	5	0	6	5	5	3	3	3	9	8	1	4	0	6
1	0	0	6	2	1	1	3	2	8	8	7	8	4	6	0	2	0	3	6
8	7	1	5	9	9	3	2	4	9	4	6	5	3	2	8	5	9	4	1
6	5	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7
8	9	0	1	2	3	4	5	6	7	8	9	6	4	2	6	4	7	5	5
4	7	8	9	2	9	3	9	3	8	2	0	9	8	0	5	6	0	1	0
4	2	6	5	5	5	4	3	4	1	5	3	0	8	3	0	6	2	7	1
1	8	1	7	1	3	8	5	4	2	0	9	7	6	7	4	1	6	8	4
7	5	1	2	6	7	1	9	8	0	6	9	4	9	9	6	2	3	7	1
9	2	2	5	3	7	8	0	1	2	3	4	5	6	7	8	0	1	2	3
4	5	6	7	8	0	1	2	3	4	5	6	7	8	9	2	1	2	1	3
9	9	8	5	3	7	0	7	7	5	7	9	9	4	7	0	3	4	1	4
4	7	5	8	1	4	8	4	1	8	6	6	4	6	3	5	7	2	5	9

Figure 1.9: Example inputs from the MNIST dataset. The “NIST” stands for National Institute of Standards and Technology, the agency that originally collected this data. The “M” stands for “modified,” since the data has been preprocessed for easier use with machine learning algorithms. The MNIST dataset consists of scans of handwritten digits and associated labels describing which digit 0–9 is contained in each image. This simple classification problem is one of the simplest and most widely used tests in deep learning research. It remains popular despite being quite easy for modern techniques to solve. Geoffrey Hinton has described it as “the *drosophila* of machine learning,” meaning that it enables machine learning researchers to study their algorithms in controlled laboratory conditions, much as biologists often study fruit flies.

has made machine learning much easier because the key burden of statistical estimation—generalizing well to new data after observing only a small amount of data—has been considerably lightened. As of 2016, a rough rule of thumb is that a supervised deep learning algorithm will generally achieve acceptable performance with around 5,000 labeled examples per category and will match or

8	9	0	1	2	3	4	7	8	9	0	1	2	3	4	5	6	7	8	6
4	2	6	4	7	5	5	4	7	8	9	2	9	3	9	3	8	2	0	5
0	1	0	4	2	6	5	3	5	3	8	0	0	3	4	1	5	3	0	8
3	0	6	2	7	1	1	8	1	7	1	3	8	9	7	6	7	4	1	6
7	5	1	7	1	9	8	0	6	9	4	9	9	3	7	1	9	2	2	5
3	7	8	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	0
1	2	3	4	5	6	7	8	9	8	1	0	5	5	1	9	0	4	1	9
3	8	4	7	7	8	5	0	6	5	5	3	3	3	9	8	1	4	0	6
1	0	0	6	2	1	1	3	2	8	8	7	8	4	6	0	2	0	3	6
8	7	1	5	9	9	3	2	4	9	4	6	5	3	2	5	5	9	4	1
6	5	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7
8	9	0	1	2	3	4	5	6	7	8	9	6	4	2	6	4	7	5	5
4	7	8	9	2	9	3	9	3	8	2	0	9	8	0	5	6	0	1	0
4	2	6	5	5	5	4	3	4	1	5	3	0	8	3	0	6	2	7	1
1	8	1	7	1	3	8	5	4	2	0	9	7	6	7	4	1	6	8	4
7	5	1	2	6	7	1	9	8	0	6	9	4	9	9	6	2	3	7	1
9	2	2	5	3	7	8	0	1	2	3	4	5	6	7	8	0	1	2	3
4	5	6	7	8	0	1	2	3	4	5	6	7	8	9	2	1	2	1	3
9	9	8	5	3	7	0	7	7	5	7	9	9	4	7	0	3	4	1	4
4	7	5	8	1	4	8	4	1	8	6	6	4	6	3	5	7	2	5	9

图 1.9：MNIST 数据集的示例输入。“NIST” 代表美国国家标准与技术研究院 (National Institute of Standards and Technology)，该机构最初收集了这些数据。“M” 代表 “已修改”，因为数据已进行预处理，以便于机器学习算法使用。MNIST 数据集包含手写数字扫描图以及描述每幅图像中包含哪个数字 0-9 的相关标签。这个简单的分类问题是深度学习研究中简单且最广泛的测试之一。尽管对于现代技术来说，它很容易解决，但仍然很受欢迎。

杰弗里 · 辛顿 (Geoffrey Hinton) 将其描述为 “机器学习中的果蝇”，这意味着它使机器学习研究人员能够在受控的实验室条件下研究他们的算法，就像生物学家经常研究果蝇一样。

使得机器学习变得更容易，因为统计估计的关键负担——在仅观察少量数据后很好地推广到新数据——已经大大减轻。截至 2016 年，一个粗略的经验法则是，监督深度学习算法通常会在每个类别中大约 5,000 个带标签的示例中获得可接受的性能，并且将匹配或

exceed human performance when trained with a dataset containing at least 10 million labeled examples. Working successfully with datasets smaller than this is an important research area, focusing in particular on how we can take advantage of large quantities of unlabeled examples, with unsupervised or semi-supervised learning.

### 1.2.3 Increasing Model Sizes

Another key reason that neural networks are wildly successful today after enjoying comparatively little success since the 1980s is that we have the computational resources to run much larger models today. One of the main insights of connectionism is that animals become intelligent when many of their neurons work together. An individual neuron or small collection of neurons is not particularly useful.

Biological neurons are notes especially densely connected. As seen in figure 1.10, our machine learning models have had a number of connections per neuron within an order of magnitude of even mammalian brains for decades.

In terms of the total number of neurons, neural networks have been astonishingly small until quite recently, as shown in figure 1.11. Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. This growth is driven by faster computers with larger memory and by the availability of larger datasets. Larger networks are able to achieve higher accuracy on more complex tasks. This trend looks set to continue for decades. Unless new technologies enable faster scaling, artificial neural networks will not have the same number of neurons as the human brain until at least the 2050s. Biological neurons may represent more complicated functions than current artificial neurons, so biological neural networks may be even larger than this plot portrays.

In retrospect, it is not particularly surprising that neural networks with fewer neurons than a leech were unable to solve sophisticated artificial intelligence problems. Eventoday's networks, which we consider quite large from a computational systems point of view, are smaller than the nervous system of even relatively primitive vertebrate animals like frogs.

The increase in model size over time, due to the availability of faster CPUs, the advent of general purpose GPUs (described in section 12.1.2), faster network connectivity and better software infrastructure for distributed computing, is one of the most important trends in the history of deep learning. This trend is generally expected to continue well into the future.

在使用至少包含 1000 万个标记示例的数据集进行训练时，其表现将超越人类。成功处理小于此的数据集是一个重要的研究领域，尤其关注如何利用无监督或半监督学习来利用大量未标记的示例。

### 1.2.3 增加模型尺寸

神经网络自 1980 年代以来一直没有取得太大成功，而如今却异常成功的另一个关键原因是，我们如今拥有了处理更大模型的计算资源。联结主义的主要见解之一是，当许多神经元协同工作时，动物就会变得聪明。

单个神经元或小群神经元并不是特别有用。

生物神经元的连接并不是特别密集。如图 1 所示，几十年来，我们的机器学习模型每个神经元的连接数量甚至与哺乳动物大脑的连接数量相当。

就神经元的总数而言，神经网络直到最近才变得非常小，如图所示。自从引入 1.11 个隐藏单元以来，人工神经网络的规模大约每 2.4 年翻一番。这种增长是由速度更快、内存更大的计算机以及更大数据集的可用性驱动的。更大的网络能够在更复杂的任务上实现更高的准确率。这种趋势看起来将持续几十年。除非新技术能够实现更快的扩展，否则至少在 2050 年代之前，人工神经网络的神经元数量都不会与人脑相同。生物神经元可能比当前的人工神经元代表更复杂的功能，因此生物神经网络可能比该图描绘的还要大。

回想起来，神经元比水蛭少的神经网络无法解决复杂的人工智能问题，这并不特别令人惊讶。即使从计算系统的角度来看，今天的网络已经相当大了，但它甚至比青蛙等相对原始的脊椎动物的神经系统还要小。

由于更快的 CPU 的出现、通用 GPU 的出现（如第节所述）、更快的网络连接以及更好的分布式计算软件基础设施，模型大小随着时间的推移而增加，这是深度学习历史上最重要的趋势之一。人们普遍预计这种趋势将持续到未来。

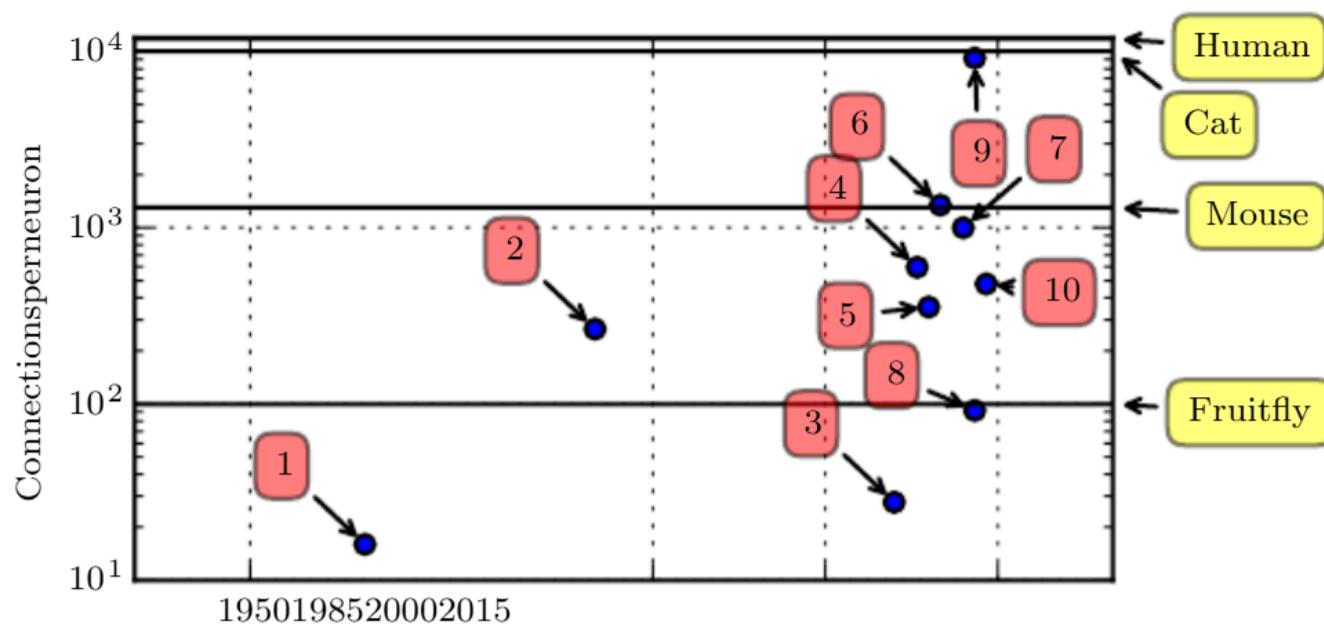


Figure1.10: Numberofconnectionspерneuronovertime.Initially,thenumberofconnectionsbetweenneuronsinartificialneuralnetworkswaslimitedbyhardwarecapabilities.Today,thenumberofconnectionsbetweenneuronsismostlyadesignconsideration.Somewhatificialneuralnetworkshavenearlyasmaymanyconnectionspерneuronasacat,andonit isquitecommonforotherneuralnetworkstohaveasmaymanyconnectionspерneuronassmallermammalslikemice.Eventhehumanbraindoesnothaveanexorbitantamount ofconnectionspерneuron.Biologicalneuralnetworkssizesfrom [Wikipedia \(2015\)](#).

1. Adaptive line element ([Widrow and Hoff , 1960](#))
2. Neocognitron ([Fukushima, 1980](#))
3. GPU-accelerated convolutional network ([Chellapilla et al. , 2006](#))
4. Deep Boltzmann machine ([Salakhutdinov and Hinton , 2009a](#))
5. Unsupervised convolutional network ([Jarrett et al. , 2009](#))
6. GPU-accelerated multilayer perceptron ([Ciresan et al. , 2010](#))
7. Distributed autoencoder ([Le et al. , 2012](#))
8. Multi-GPU convolutional network ([Krizhevsky et al. , 2012](#))
9. COTS HPC unsupervised convolutional network ([Coates et al. , 2013](#))
10. GoogLeNet ([Szegedy et al. , 2014a](#))

## 1.2.4 Increasing Accuracy, Complexity and Real-World Impact

Since the 1980s, deep learning has consistently improved its ability to provide accurate recognition and prediction. Moreover, deep learning has consistently been applied with success to broader and broader sets of applications.

The earliest deep models were used to recognize individual objects in tightly cropped, extremely small images ([Rumelhart et al. , 1986a](#)). Since then there has been a gradual increase in the size of images neural networks could process. Modern object recognition networks process rich high-resolution photographs and don't

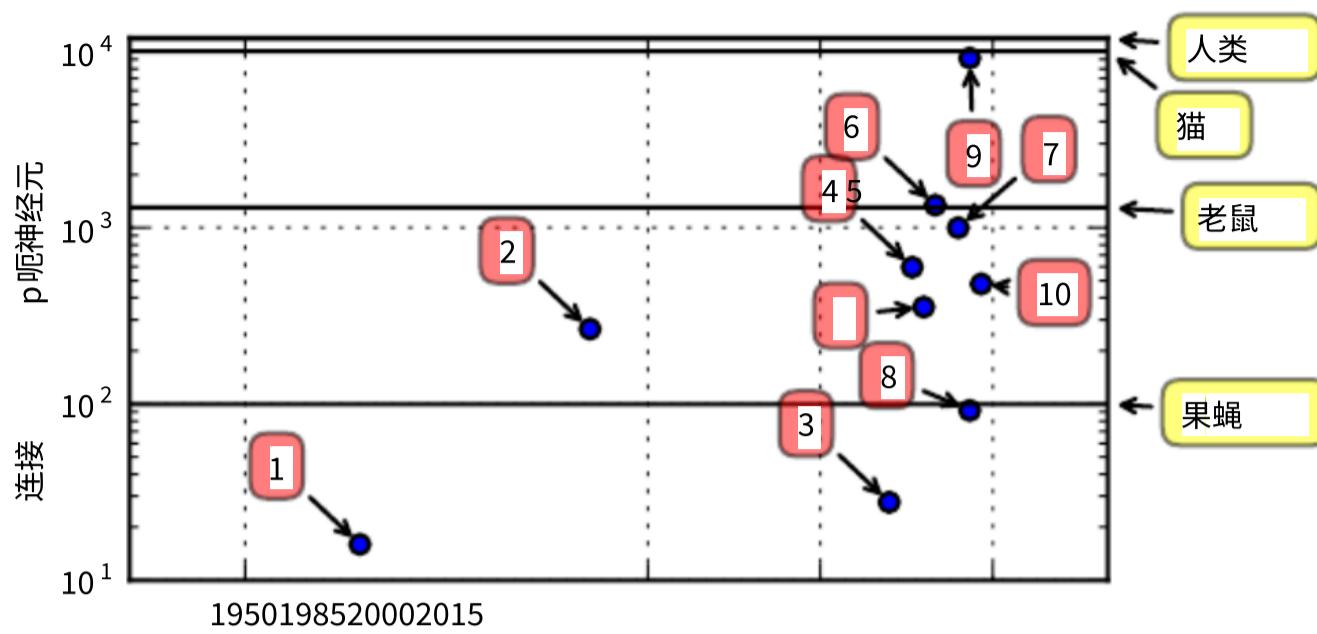


图 1.10：每个神经元的连接数量随时间变化。最初，人工神经网络中神经元之间的连接数量受限于硬件能力。如今，神经元之间的连接数量主要取决于设计。一些人工神经网络早期每个神经元的连接数量与猫一样多，其他神经网络每个神经元的连接数量也相当常见。

体型较小的哺乳动物，如老鼠。即使人类大脑，每个神经元的连接数量也远不及人类。生物神经网络的大小来自 [\(\)](#)。维基百科 2015

- 1.自适应线性单元 (, ) Widrow and Hoff 1960
- 2.Neocognitron (福岛 1980,)
- 3.GPU 加速卷积网络 (, ) Chellapilla 等人, 2006
- 4.深玻尔兹曼机 (Salakhutdinov and Hinton 2009a,)
- 5.无监督卷积网络 (, ) Jarrett 等人, 2009
- 6.GPU 加速多层感知器 (, ) Ciresan 等人, 2010
- 7.分布式自动编码器 (Le et al. 2012)
- 8.多 GPU 卷积网络 (, ) Krizhevsky 等人, 2012
- 9.COTS HPC unsupervised convolutional network (,) Coates 等人, 2013
- 10.GoogLeNet (,) Szegedy 等人。2014a

#### 1.2.4 提高准确性、复杂性和现实世界影响

自 20 世纪 80 年代以来，深度学习不断提高其提供准确识别和预测的能力。此外，深度学习已成功应用于越来越广泛的应用领域。

最早的深度模型用于识别紧密裁剪、极小图像中的单个物体 (, )。自 1986 年 Rumelhart 等人提出以来，神经网络可以处理的图像大小逐渐增加。现代物体识别网络可以处理丰富的高分辨率照片，并且不会

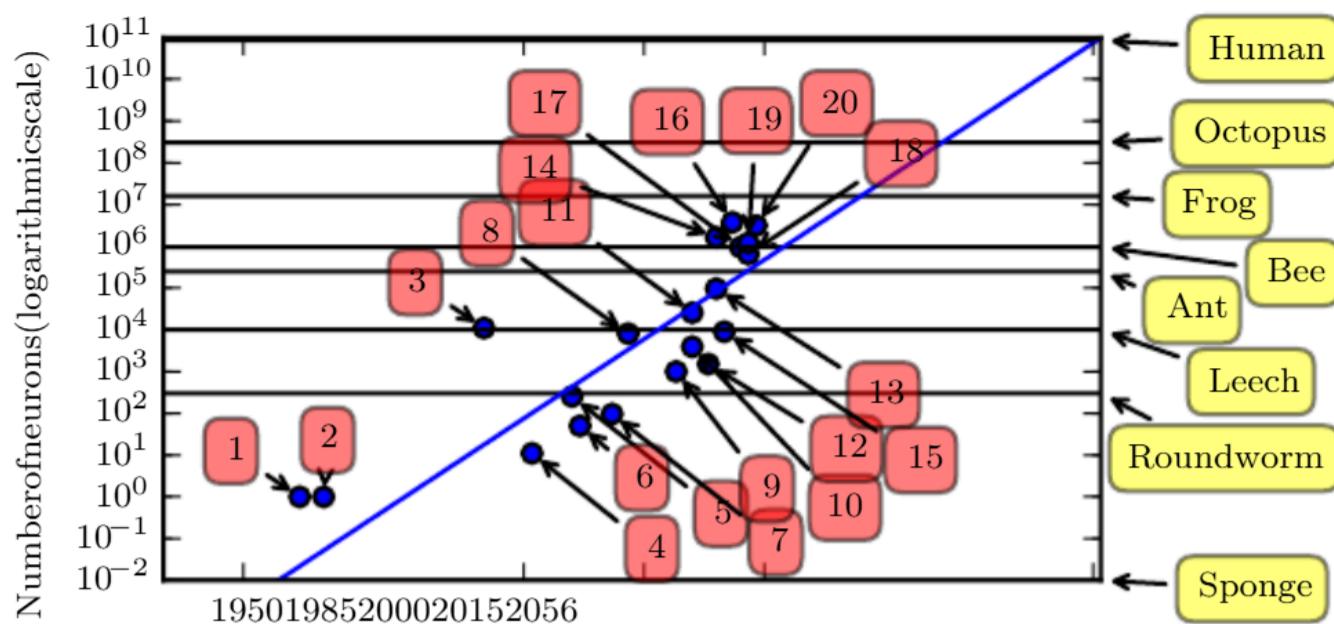


Figure 1.11: Increasing neural network size over time. Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from [Wikipedia \(2015\)](#).

1. Perceptron (Rosenblatt, 1958, 1962)
2. Adaptive linear element (Widrow and Hoff, 1960)
3. Neocognitron (Fukushima, 1980)
4. Early back-propagation network (Rumelhart et al., 1986b)
5. Recurrent neural network for speech recognition (Robinson and Fallside, 1991)
6. Multilayer perceptron for speech recognition (Bengio et al., 1991)
7. Meanfield sigmoid belief network (Saul et al., 1996)
8. LeNet-5 (LeCun et al., 1998b)
9. Echo state network (Jaeger and Haas, 2004)
10. Deep belief network (Hinton et al., 2006)
11. GPU-accelerated convolutional network (Chellapilla et al., 2006)
12. Deep Boltzmann machine (Salakhutdinov and Hinton, 2009a)
13. GPU-accelerated deep belief network (Raina et al., 2009)
14. Unsupervised convolutional network (Jarrett et al., 2009)
15. GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
16. OMP-1 network (Coates and Ng, 2011)
17. Distributed autoencoder (Le et al., 2012)
18. Multi-GPU convolutional network (Krizhevsky et al., 2012)
19. COTS HPC unsupervised convolutional network (Coates et al., 2013)
20. GoogLeNet (Szegedy et al., 2014a)

have a requirement that the photo be cropped near the object to be recognized (Krizhevsky et al., 2012). Similarly, the earliest networks could recognize only two kinds of objects (or in some cases, the absence or presence of a single kind of object), while these modern networks typically recognize at least 1,000 different categories of objects. The largest contest in object recognition is the ImageNet

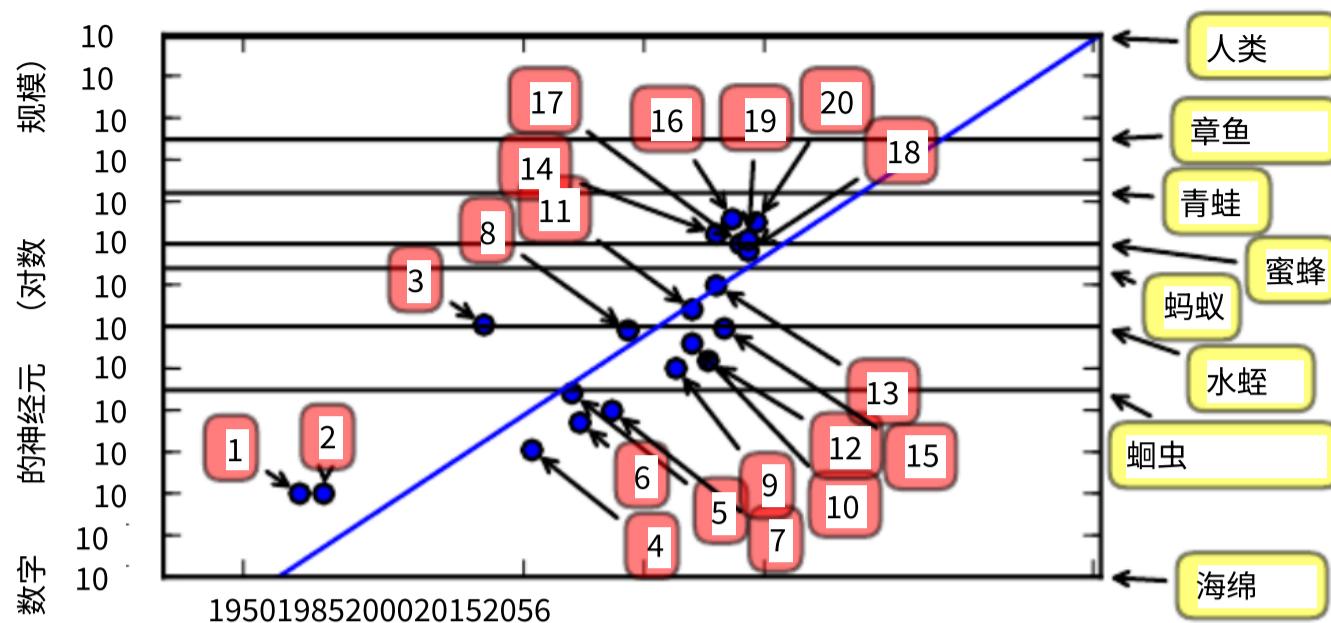


图 1.11：神经网络规模随时间推移而增加。自引入隐藏单元以来，人工神经网络的规模大约每 2.4 年翻一番。生物神经网络的规模从 ()。维基百科 2015

1. 感知器 (,,) Rosenblatt 1958 1962
2. 自适应线性单元 (,,) Widrow and Hoff 1960
3. Neocognitron (福岛 1980,)
4. 早期反向传播网络 (,,) Rumelhart 等人, 1986b
5. Recurrent neural network for speech recognition (Robinson and Fallside 1991,)
6. Multilayer perceptron for speech recognition (,,) Bengio 等人, 1991
7. Meanfield sigmoid belief network (,) Saul 等人, 1996 年
8. LeNet-5 (,) LeCun 等人。1998b
9. Echo state network (,) Jaeger and Haas 2004
10. 深度信念网络 (Hinton 2006 等,)
11. GPU 加速卷积网络 (,) Chellapilla 等人, 2006 年
12. 深玻尔兹曼机 (Salakhutdinov and Hinton 2009a,)
13. GPU 加速的深度信念网络 (Raina 等人, 2009 年)
14. 无监督卷积网络 (,,) Jarrett 等人, 2009
15. GPU 加速多层感知器 (,,) Ciresan 等人, 2010
16. OMP-1 网络 (,) Coates and Ng 2011
17. 分布式自动编码器 (Le et al. 2012)
18. 多 GPU 卷积网络 (,,) Krizhevsky 等人, 2012
19. COTS HPC unsupervised convolutional network (,) Coates 等人, 2013
20. GoogLeNet (,) Szegedy 等人。2014a

要求将照片裁剪到物体附近才能识别 (,,)。同样，最早的网络只能识别 Krizhevsky et al. 2012 两种物体（或者在某些情况下，识别单一类型物体的存在或不存在），而这些现代网络通常可以识别至少 1,000 种不同类别的物体。物体识别领域最大的竞赛是 ImageNet

Large Scale Visual Recognition Challenge (ILSVRC) held each year. A dramatic moment in the meteoric rise of deep learning came when a convolutional network won this challenge for the first time and by a wide margin, bringing down the state-of-the-art top-5 error rate from 26.1 percent to 15.3 percent (Krizhevsky et al., 2012), meaning that the convolutional network produces a ranked list of possible categories for each image, and the correct category appeared in the first five entries of this list for all but 15.3 percent of the test examples. Since then, these competitions are consistently won by deep convolutional nets, and as of this writing, advances in deep learning have brought the latest top-5 error rate in this contest down to 3.6 percent, as shown in figure 1.12.

Deep learning has also had a dramatic impact on speech recognition. After improving throughout the 1990s, the error rates for speech recognition stagnated starting in about 2000. The introduction of deep learning (Dahl et al., 2010; Deng et al., 2010b; Seide et al., 2011; Hinton et al., 2012a) to speech recognition resulted in a sudden drop in error rates, with some error rates cut in half. We explore this history in more detail in section 12.3.

Deep networks have also had spectacular successes for pedestrian detection and image segmentation (Sermanet et al., 2013; Farabet et al., 2013; Couprie et al., 2013) and yielded superhuman performance in traffic sign classification (Ciresan et al., 2012).

At the same time that the scale and accuracy of deep networks have increased,

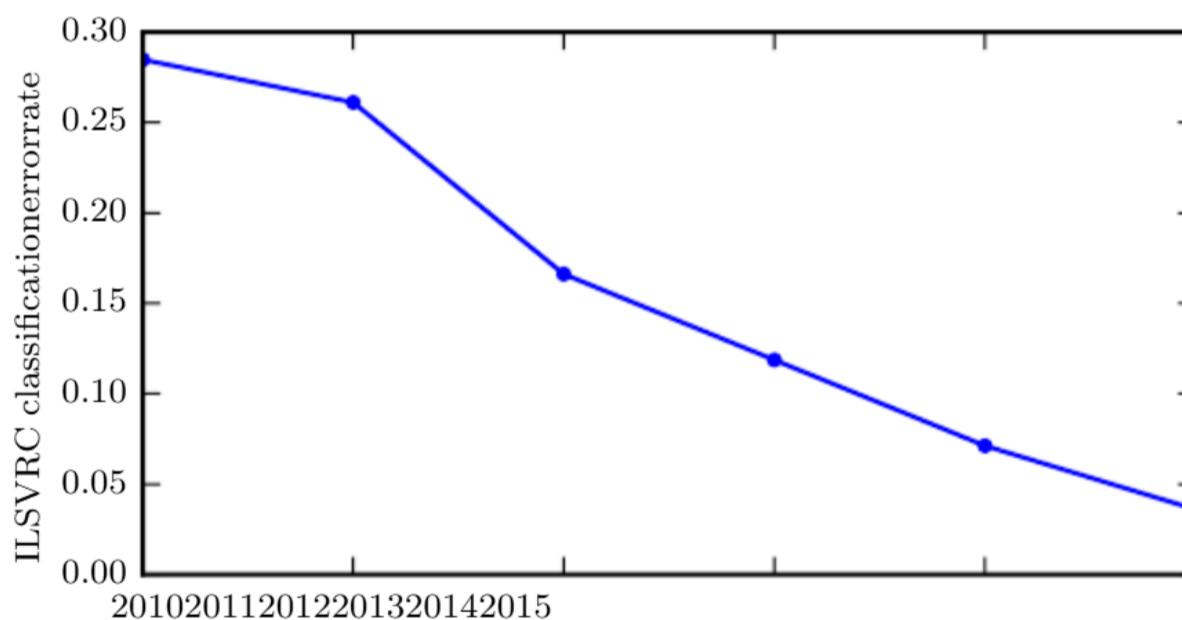


Figure 1.12: Decreasing error rate over time. Since deep networks reached the scale necessary to compete in the ImageNet Large Scale Visual Recognition Challenge, they have consistently won the competition every year, yielding lower and lower error rates each time. Data from Russakovsky et al. (2014b) and He et al. (2015).

大规模视觉识别挑战赛 (ILSVRC) 每年举办一次。深度学习飞速发展的戏剧性时刻是，一个卷积网络首次以大幅优势赢得该挑战，将当时最先进的 top-5 错误率从 26.1% 降至 15.3% (Krizhevsky et al., )，这意味着卷积网络为每个图像生成一个包含 2012 个可能类别的排序列表，并且除了 15.3% 的测试示例外，所有测试示例的正确类别都出现在此列表的前五个条目中。自那以后，这些比赛一直由深度卷积网络赢得，截至本文撰写时，深度学习的进步已将本次竞赛的最新 top-5 错误率降至 3.6%，如图 1.12 所示深度学习也对语音识别产生了巨大的影响。在经历了整个 20 世纪 90 年代的改进之后，语音识别的错误率从 2000 年左右开始停滞不前。深度学习 (,; Dahl et al. 2010 Deng et al. et al. 2010b Seide ,; 2011 Hinton ,) 引入语音识别导致错误率突然下降，有些错误率减少了一半。我们将在第 12.3 节中更详细地探讨这段历史。深度网络在行人检测和图像分割 (,; Sermanet et al. 2013 Farabet 2013Couprie et al. ,; et al. , 2013) 并在交通标志分类中取得了超越人类的表现 (Ciresan et al. , )。2012 在深度网络规模和准确性不断提高的同时，

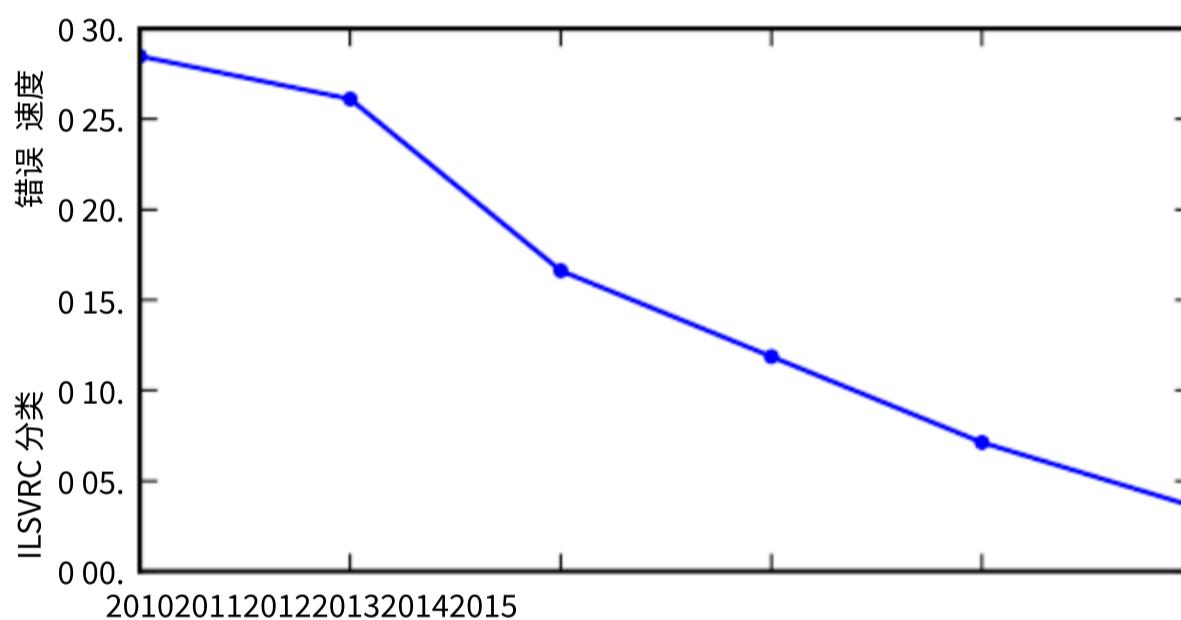


图 1.12：随着时间的推移，错误率不断降低。自从深度网络达到参加 ImageNet 大规模视觉识别挑战赛所需的规模以来，它们每年都连续赢得比赛，并且每次的错误率都越来越低。数据来自 Russakovsky 2014b、He 2015 等 ( ) 和等 ( )。

so has the complexity of the tasks that they can solve. **Goodfellow et al.** (2014d) showed that neural networks could learn to output an entire sequence of characters transcribed from an image, rather than just identifying a single object. Previously, it was widely believed that this kind of learning required labeling of the individual elements of the sequence (Gülçehre and Bengio , 2013). Recurrent neural networks, such as the LSTM sequence model mentioned above, are now used to model relationships between *sequences* and other *sequences* rather than just fixed inputs. This sequence-to-sequence learning seems to be on the cusp of revolutionizing another application: machine translation (Sutskever et al. , 2014 Bahdanau; et al. , 2015).

This trend of increasing complexity has been pushed to its logical conclusion with the introduction of neural Turing machines (Graves et al. , 2014) that learn to read from memory cells and write arbitrary content to memory cells. Such neural networks can learn simple programs from examples of desired behavior. For example, they can learn to sort lists of numbers given examples of scrambled and sorted sequences. This self-programming technology is in its infancy, but in the future it could in principle be applied to nearly any task.

Another crowning achievement of deep learning is its extension to the domain of **reinforcement learning**. In the context of reinforcement learning, an autonomous agent must learn to perform a task by trial and error, without any guidance from the human operator. DeepMind demonstrated that a reinforcement learning system based on deep learning is capable of learning to play Atari video games, reaching human-level performance on many tasks (Mnih et al. , 2015). Deep learning has also significantly improved the performance of reinforcement learning for robotics (Finn et al. , 2015).

Many of these applications of deep learning are highly profitable. Deep learning is now used by many top technology companies, including Google, Microsoft, Facebook, IBM, Baidu, Apple, Adobe, Netflix, NVIDIA, and NEC.

Advances in deep learning have also depended heavily on advances in software infrastructure. Software libraries such as Theano (Bergstra et al. , 2010; Bastien et al. , 2012), PyLearn2 (Goodfellow et al. , 2013c), Torch (Collobert et al. , 2011b), DistBelief (Dean et al. , 2012), Caffe (Jia, 2013), MXNet (Chen et al. , 2015), and TensorFlow (Abadi et al. , 2015) have all supported important research projects or commercial products.

Deep learning has also made contributions to other sciences. Modern convolutional networks for object recognition provide a model of visual processing that neuroscientists can study (DiCarlo, 2013). Deep learning also provides useful tools for processing massive amounts of data and making useful predictions in scientific

因此，它们所能解决的任务也很复杂。（）Goodfellow 等人，2014 年发现，神经网络可以学习输出从图像转录的整个字符序列，而不仅仅是识别单个对象。此前，人们普遍认为这种学习需要标记序列中的各个元素（，）。循环神经网络，例如上面提到的 LSTM 序列模型，现在用于模拟序列之间的关系，而不仅仅是固定的输入。这种序列到序列的学习似乎正处于另一个应用革命的风口浪尖：机器翻译（Sutskever, 2014；Bahdanau et al., ; 等人，2015）。

随着神经图灵机（Graves 2014 等）的引入，这种复杂性不断增加的趋势被推向了逻辑结论，神经图灵机可以学习从记忆单元读取数据并向记忆单元写入任意内容。这种神经网络可以从所需行为的示例中学习简单的程序。例如，它们可以学习根据给出的混乱和排序序列的示例对数字列表进行排序。这种自编程技术尚处于起步阶段，但在未来，它原则上可以应用于任何任务。

深度学习的另一项伟大成就是其向强化学习领域的扩展。在强化学习的背景下，自主智能体必须通过反复试验来学习特定格式的任务，而无需人类操作员的任何指导。DeepMind 证明，基于深度学习的强化学习系统能够学习玩 Atari 电子游戏，在许多任务上达到人类水平的表现（，）。深度学习还显著提高了机器人强化学习的性能（，）。Finn 等人，2015

深度学习的许多应用都利润丰厚。目前，许多顶尖科技公司都在使用深度学习，包括谷歌、微软、Facebook、IBM、百度、苹果、Adobe、Netflix、NVIDIA 和 NEC。

深度学习的进步也在很大程度上依赖于软件基础设施的进步。诸如 Theano (2011, 2012; Bergstra 等人, 2010 Bastien 等人, 2011)、PyLearn2 (2012 Goodfellow 等人, 2012)、Torch (2013、2014)、DistBelief (2013、2014)、Caffe (2013、2014)、MXNet (2013、2014) 和 TensorFlow (2013、2014) 等软件库都支持重要的研究项目或商业产品。

深度学习也对其他科学做出了贡献。用于物体识别的现代卷积网络为神经科学家提供了研究视觉处理的理想模型（，）。深度学习还提供了有用的工具 DiCarlo 2013，用于处理大量数据并在科学的研究中做出有用的预测

fields. It has been successfully used to predict how molecules will interact in order to help pharmaceutical companies design new drugs (Dahl *et al.*, 2014), to search for subatomic particles (Baldi *et al.*, 2014), and to automatically parse microscope images used to construct a 3-D map of the human brain (Knowles-Barley *et al.*, 2014). We expect deep learning to appear in more and more scientific fields in the future.

In summary, deep learning is an approach to machine learning that has drawn heavily on our knowledge of the human brain, statistics and applied math as it developed over the past several decades. In recent years, deep learning has seen tremendous growth in its popularity and usefulness, largely as the result of more powerful computers, larger datasets and techniques to train deeper networks. The years ahead are full of challenges and opportunities to improve deep learning even further and to bring it to new frontiers.

它已被成功用于预测分子如何相互作用，以帮助制药公司设计新药（，），搜索 Dahl et al. 2014 中的亚原子粒子（，），以及自动解析用于构建人脑 3-D 地图的显微镜 Baldi et al. 2014 图像（，Knowles-Barley et al. 2014）。我们期待深度学习在未来出现在越来越多的科学领域。

总而言之，深度学习是一种机器学习方法，它大量借鉴了我们对人脑、统计学和应用数学的知识，并在过去几十年中得到了发展。近年来，深度学习的普及度和实用性得到了巨大增长，这主要是由于更强大的计算机、更大的数据集和训练更深层网络的技术。未来几年充满了挑战和机遇，可以进一步改进深度学习并将其带入新的领域。