

# Deep Learning

Ian Goodfellow  
Yoshua Bengio  
Aaron Courville

# Contents

<b>Website</b>	<b>vii</b>
<b>Acknowledgments</b>	<b>viii</b>
<b>Notation</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Who Should Read This Book? . . . . .	8
1.2 Historical Trends in Deep Learning . . . . .	11
<b>I Applied Math and Machine Learning Basics</b>	<b>29</b>
<b>2 Linear Algebra</b>	<b>31</b>
2.1 Scalars, Vectors, Matrices and Tensors . . . . .	31
2.2 Multiplying Matrices and Vectors . . . . .	34
2.3 Identity and Inverse Matrices . . . . .	36
2.4 Linear Dependence and Span . . . . .	37
2.5 Norms . . . . .	39
2.6 Special Kinds of Matrices and Vectors . . . . .	40
2.7 Eigendecomposition . . . . .	42
2.8 Singular Value Decomposition . . . . .	44
2.9 The Moore-Penrose Pseudoinverse . . . . .	45
2.10 The Trace Operator . . . . .	46
2.11 The Determinant . . . . .	47
2.12 Example: Principal Components Analysis . . . . .	48
<b>3 Probability and Information Theory</b>	<b>53</b>
3.1 Why Probability? . . . . .	54

3.2	Random Variables . . . . .	56
3.3	Probability Distributions . . . . .	56
3.4	Marginal Probability . . . . .	58
3.5	Conditional Probability . . . . .	59
3.6	The Chain Rule of Conditional Probabilities . . . . .	59
3.7	Independence and Conditional Independence . . . . .	60
3.8	Expectation, Variance and Covariance . . . . .	60
3.9	Common Probability Distributions . . . . .	62
3.10	Useful Properties of Common Functions . . . . .	67
3.11	Bayes' Rule . . . . .	70
3.12	Technical Details of Continuous Variables . . . . .	71
3.13	Information Theory . . . . .	73
3.14	Structured Probabilistic Models . . . . .	75
<b>4</b>	<b>Numerical Computation</b>	<b>80</b>
4.1	Overflow and Underflow . . . . .	80
4.2	Poor Conditioning . . . . .	82
4.3	Gradient-Based Optimization . . . . .	82
4.4	Constrained Optimization . . . . .	93
4.5	Example: Linear Least Squares . . . . .	96
<b>5</b>	<b>Machine Learning Basics</b>	<b>98</b>
5.1	Learning Algorithms . . . . .	99
5.2	Capacity, Overfitting and Underfitting . . . . .	110
5.3	Hyperparameters and Validation Sets . . . . .	120
5.4	Estimators, Bias and Variance . . . . .	122
5.5	Maximum Likelihood Estimation . . . . .	131
5.6	Bayesian Statistics . . . . .	135
5.7	Supervised Learning Algorithms . . . . .	140
5.8	Unsupervised Learning Algorithms . . . . .	146
5.9	Stochastic Gradient Descent . . . . .	151
5.10	Building a Machine Learning Algorithm . . . . .	153
5.11	Challenges Motivating Deep Learning . . . . .	155
<b>II</b>	<b>Deep Networks: Modern Practices</b>	<b>166</b>
<b>6</b>	<b>Deep Feedforward Networks</b>	<b>168</b>
6.1	Example: Learning XOR . . . . .	171
6.2	Gradient-Based Learning . . . . .	177

6.3	Hidden Units . . . . .	191
6.4	Architecture Design . . . . .	197
6.5	Back-Propagation and Other Differentiation Algorithms . . . . .	204
6.6	Historical Notes . . . . .	224
<b>7</b>	<b>Regularization for Deep Learning</b>	<b>228</b>
7.1	Parameter Norm Penalties . . . . .	230
7.2	Norm Penalties as Constrained Optimization . . . . .	237
7.3	Regularization and Under-Constrained Problems . . . . .	239
7.4	Dataset Augmentation . . . . .	240
7.5	Noise Robustness . . . . .	242
7.6	Semi-Supervised Learning . . . . .	243
7.7	Multi-Task Learning . . . . .	244
7.8	Early Stopping . . . . .	246
7.9	Parameter Tying and Parameter Sharing . . . . .	253
7.10	Sparse Representations . . . . .	254
7.11	Bagging and Other Ensemble Methods . . . . .	256
7.12	Dropout . . . . .	258
7.13	Adversarial Training . . . . .	268
7.14	Tangent Distance, Tangent Prop, and Manifold Tangent Classifier	270
<b>8</b>	<b>Optimization for Training Deep Models</b>	<b>274</b>
8.1	How Learning Differs from Pure Optimization . . . . .	275
8.2	Challenges in Neural Network Optimization . . . . .	282
8.3	Basic Algorithms . . . . .	294
8.4	Parameter Initialization Strategies . . . . .	301
8.5	Algorithms with Adaptive Learning Rates . . . . .	306
8.6	Approximate Second-Order Methods . . . . .	310
8.7	Optimization Strategies and Meta-Algorithms . . . . .	317
<b>9</b>	<b>Convolutional Networks</b>	<b>330</b>
9.1	The Convolution Operation . . . . .	331
9.2	Motivation . . . . .	335
9.3	Pooling . . . . .	339
9.4	Convolution and Pooling as an Infinitely Strong Prior . . . . .	345
9.5	Variants of the Basic Convolution Function . . . . .	347
9.6	Structured Outputs . . . . .	358
9.7	Data Types . . . . .	360
9.8	Efficient Convolution Algorithms . . . . .	362
9.9	Random or Unsupervised Features . . . . .	363

9.10	The Neuroscientific Basis for Convolutional Networks . . . . .	364
9.11	Convolutional Networks and the History of Deep Learning . . . . .	371
<b>10</b>	<b>Sequence Modeling: Recurrent and Recursive Nets</b>	<b>373</b>
10.1	Unfolding Computational Graphs . . . . .	375
10.2	Recurrent Neural Networks . . . . .	378
10.3	Bidirectional RNNs . . . . .	394
10.4	Encoder-Decoder Sequence-to-Sequence Architectures . . . . .	396
10.5	Deep Recurrent Networks . . . . .	398
10.6	Recursive Neural Networks . . . . .	400
10.7	The Challenge of Long-Term Dependencies . . . . .	401
10.8	Echo State Networks . . . . .	404
10.9	Leaky Units and Other Strategies for Multiple Time Scales . . . . .	406
10.10	The Long Short-Term Memory and Other Gated RNNs . . . . .	408
10.11	Optimization for Long-Term Dependencies . . . . .	413
10.12	Explicit Memory . . . . .	416
<b>11</b>	<b>Practical Methodology</b>	<b>421</b>
11.1	Performance Metrics . . . . .	422
11.2	Default Baseline Models . . . . .	425
11.3	Determining Whether to Gather More Data . . . . .	426
11.4	Selecting Hyperparameters . . . . .	427
11.5	Debugging Strategies . . . . .	436
11.6	Example: Multi-Digit Number Recognition . . . . .	440
<b>12</b>	<b>Applications</b>	<b>443</b>
12.1	Large-Scale Deep Learning . . . . .	443
12.2	Computer Vision . . . . .	452
12.3	Speech Recognition . . . . .	458
12.4	Natural Language Processing . . . . .	461
12.5	Other Applications . . . . .	478
<b>III</b>	<b>Deep Learning Research</b>	<b>486</b>
<b>13</b>	<b>Linear Factor Models</b>	<b>489</b>
13.1	Probabilistic PCA and Factor Analysis . . . . .	490
13.2	Independent Component Analysis (ICA) . . . . .	491
13.3	Slow Feature Analysis . . . . .	493
13.4	Sparse Coding . . . . .	496

13.5	Manifold Interpretation of PCA . . . . .	499
<b>14</b>	<b>Autoencoders</b>	<b>502</b>
14.1	Undercomplete Autoencoders . . . . .	503
14.2	Regularized Autoencoders . . . . .	504
14.3	Representational Power, Layer Size and Depth . . . . .	508
14.4	Stochastic Encoders and Decoders . . . . .	509
14.5	Denoising Autoencoders . . . . .	510
14.6	Learning Manifolds with Autoencoders . . . . .	515
14.7	Contractive Autoencoders . . . . .	521
14.8	Predictive Sparse Decomposition . . . . .	523
14.9	Applications of Autoencoders . . . . .	524
<b>15</b>	<b>Representation Learning</b>	<b>526</b>
15.1	Greedy Layer-Wise Unsupervised Pretraining . . . . .	528
15.2	Transfer Learning and Domain Adaptation . . . . .	536
15.3	Semi-Supervised Disentangling of Causal Factors . . . . .	541
15.4	Distributed Representation . . . . .	546
15.5	Exponential Gains from Depth . . . . .	553
15.6	Providing Clues to Discover Underlying Causes . . . . .	554
<b>16</b>	<b>Structured Probabilistic Models for Deep Learning</b>	<b>558</b>
16.1	The Challenge of Unstructured Modeling . . . . .	559
16.2	Using Graphs to Describe Model Structure . . . . .	563
16.3	Sampling from Graphical Models . . . . .	580
16.4	Advantages of Structured Modeling . . . . .	582
16.5	Learning about Dependencies . . . . .	582
16.6	Inference and Approximate Inference . . . . .	584
16.7	The Deep Learning Approach to Structured Probabilistic Models	585
<b>17</b>	<b>Monte Carlo Methods</b>	<b>590</b>
17.1	Sampling and Monte Carlo Methods . . . . .	590
17.2	Importance Sampling . . . . .	592
17.3	Markov Chain Monte Carlo Methods . . . . .	595
17.4	Gibbs Sampling . . . . .	599
17.5	The Challenge of Mixing between Separated Modes . . . . .	599
<b>18</b>	<b>Confronting the Partition Function</b>	<b>605</b>
18.1	The Log-Likelihood Gradient . . . . .	606
18.2	Stochastic Maximum Likelihood and Contrastive Divergence . . .	607

18.3	Pseudolikelihood . . . . .	615
18.4	Score Matching and Ratio Matching . . . . .	617
18.5	Denoising Score Matching . . . . .	619
18.6	Noise-Contrastive Estimation . . . . .	620
18.7	Estimating the Partition Function . . . . .	623
<b>19</b>	<b>Approximate Inference</b>	<b>631</b>
19.1	Inference as Optimization . . . . .	633
19.2	Expectation Maximization . . . . .	634
19.3	MAP Inference and Sparse Coding . . . . .	635
19.4	Variational Inference and Learning . . . . .	638
19.5	Learned Approximate Inference . . . . .	651
<b>20</b>	<b>Deep Generative Models</b>	<b>654</b>
20.1	Boltzmann Machines . . . . .	654
20.2	Restricted Boltzmann Machines . . . . .	656
20.3	Deep Belief Networks . . . . .	660
20.4	Deep Boltzmann Machines . . . . .	663
20.5	Boltzmann Machines for Real-Valued Data . . . . .	676
20.6	Convolutional Boltzmann Machines . . . . .	683
20.7	Boltzmann Machines for Structured or Sequential Outputs . . . . .	685
20.8	Other Boltzmann Machines . . . . .	686
20.9	Back-Propagation through Random Operations . . . . .	687
20.10	Directed Generative Nets . . . . .	692
20.11	Drawing Samples from Autoencoders . . . . .	711
20.12	Generative Stochastic Networks . . . . .	714
20.13	Other Generation Schemes . . . . .	716
20.14	Evaluating Generative Models . . . . .	717
20.15	Conclusion . . . . .	720
	<b>Bibliography</b>	<b>721</b>
	<b>Index</b>	<b>777</b>

# Website

[www.deeplearningbook.org](http://www.deeplearningbook.org)

This book is accompanied by the above website. The website provides a variety of supplementary material, including exercises, lecture slides, corrections of mistakes, and other resources that should be useful to both readers and instructors.



# Acknowledgments

This book would not have been possible without the contributions of many people.

We would like to thank those who commented on our proposal for the book and helped plan its contents and organization: Guillaume Alain, Kyunghyun Cho, Çağlar Gülçehre, David Krueger, Hugo Larochelle, Razvan Pascanu and Thomas Rohée.

We would like to thank the people who offered feedback on the content of the book itself. Some offered feedback on many chapters: Martín Abadi, Guillaume Alain, Ion Androutsopoulos, Fred Bertsch, Olexa Bilaniuk, Ufuk Can Biçici, Matko Bošnjak, John Boersma, Greg Brockman, Alexandre de Brébisson, Pierre Luc Carrier, Sarath Chandar, Pawel Chilinski, Mark Daoust, Oleg Dashevskii, Laurent Dinh, Stephan Dreseidl, Jim Fan, Miao Fan, Meire Fortunato, Frédéric Francis, Nando de Freitas, Çağlar Gülçehre, Jurgen Van Gael, Javier Alonso García, Jonathan Hunt, Gopi Jeyaram, Chingiz Kabayev, Lukasz Kaiser, Varun Kanade, Asifullah Khan, Akiel Khan, John King, Diederik P. Kingma, Yann LeCun, Rudolf Mathey, Matías Mattamala, Abhinav Maurya, Kevin Murphy, Oleg Mürk, Roman Novak, Augustus Q. Odena, Simon Pavlik, Karl Pichotta, Eddie Pierce, Kari Pulli, Roussel Rahman, Tapani Raiko, Anurag Ranjan, Johannes Roith, Mihaela Rosca, Halis Sak, César Salgado, Grigory Sapunov, Yoshinori Sasaki, Mike Schuster, Julian Serban, Nir Shabat, Ken Shirriff, Andre Simpelo, Scott Stanley, David Sussillo, Ilya Sutskever, Carles Gelada Sáez, Graham Taylor, Valentin Tolmer, Massimiliano Tomassoli, An Tran, Shubhendu Trivedi, Alexey Umnov, Vincent Vanhoucke, Marco Visentini-Scarzanella, Martin Vita, David Warde-Farley, Dustin Webb, Kelvin Xu, Wei Xue, Ke Yang, Li Yao, Zygmunt Zajac and Ozan Çağlayan.

We would also like to thank those who provided us with useful feedback on individual chapters:

- **Notation:** Zhang Yuanhang.
- Chapter 1, **Introduction:** Yusuf Akgul, Sebastien Bratieres, Samira Ebrahimi,

Charlie Gorichanaz, Brendan Loudermilk, Eric Morris, Cosmin Pârvulescu and Alfredo Solano.

- Chapter 2, **Linear Algebra**: Amjad Almahairi, Nikola Banić, Kevin Bennett, Philippe Castonguay, Oscar Chang, Eric Fosler-Lussier, Andrey Khalyavin, Sergey Oreshkov, István Petrás, Dennis Prangle, Thomas Rohée, Gitanjali Gulve Sehgal, Colby Toland, Alessandro Vitale and Bob Welland.
- Chapter 3, **Probability and Information Theory**: John Philip Anderson, Kai Arulkumaran, Vincent Dumoulin, Rui Fa, Stephan Gouws, Artem Oboturov, Antti Rasmus, Alexey Surkov and Volker Tresp.
- Chapter 4, **Numerical Computation**: Tran Lam An, Ian Fischer and Hu Yuhuang.
- Chapter 5, **Machine Learning Basics**: Dzmitry Bahdanau, Justin Domingue, Nikhil Garg, Makoto Otsuka, Bob Pepin, Philip Popien, Emmanuel Rayner, Peter Shepard, Kee-Bong Song, Zheng Sun and Andy Wu.
- Chapter 6, **Deep Feedforward Networks**: Uriel Berdugo, Fabrizio Bottarel, Elizabeth Burl, Ishan Durugkar, Jeff Hlywa, Jong Wook Kim, David Krueger and Aditya Kumar Praharaj.
- Chapter 7, **Regularization for Deep Learning**: Morten Kolbæk, Kshitij Lauria, Inkyu Lee, Sunil Mohan, Hai Phong Phan and Joshua Salisbury.
- Chapter 8, **Optimization for Training Deep Models**: Marcel Ackermann, Peter Armitage, Rowel Atienza, Andrew Brock, Tegan Maharaj, James Martens, Kashif Rasul, Klaus Strobl and Nicholas Turner.
- Chapter 9, **Convolutional Networks**: Martín Arjovsky, Eugene Brevdo, Konstantin Divilov, Eric Jensen, Mehdi Mirza, Alex Paino, Marjorie Sayer, Ryan Stout and Wentao Wu.
- Chapter 10, **Sequence Modeling: Recurrent and Recursive Nets**: Gökçen Eraslan, Steven Hickson, Razvan Pascanu, Lorenzo von Ritter, Rui Rodrigues, Dmitriy Serdyuk, Dongyu Shi and Kaiyu Yang.
- Chapter 11, **Practical Methodology**: Daniel Beckstein.
- Chapter 12, **Applications**: George Dahl, Vladimir Nekrasov and Ribana Roscher.
- Chapter 13, **Linear Factor Models**: Jayanth Koushik.

- Chapter 15, **Representation Learning**: Kunal Ghosh.
- Chapter 16, **Structured Probabilistic Models for Deep Learning**: Minh Lê and Anton Varfolom.
- Chapter 18, **Confronting the Partition Function**: Sam Bowman.
- Chapter 19, **Approximate Inference**: Yujia Bao.
- Chapter 20, **Deep Generative Models**: Nicolas Chapados, Daniel Galvez, Wenming Ma, Fady Medhat, Shakir Mohamed and Grégoire Montavon.
- Bibliography: Lukas Michelbacher and Leslie N. Smith.

We also want to thank those who allowed us to reproduce images, figures or data from their publications. We indicate their contributions in the figure captions throughout the text.

We would like to thank Lu Wang for writing pdf2htmlEX, which we used to make the web version of the book, and for offering support to improve the quality of the resulting HTML.

We would like to thank Ian's wife Daniela Flori Goodfellow for patiently supporting Ian during the writing of the book as well as for help with proofreading.

We would like to thank the Google Brain team for providing an intellectual environment where Ian could devote a tremendous amount of time to writing this book and receive feedback and guidance from colleagues. We would especially like to thank Ian's former manager, Greg Corrado, and his current manager, Samy Bengio, for their support of this project. Finally, we would like to thank Geoffrey Hinton for encouragement when writing was difficult.

# Notation

This section provides a concise reference describing the notation used throughout this book. If you are unfamiliar with any of the corresponding mathematical concepts, we describe most of these ideas in chapters 2–4.

## Numbers and Arrays

$a$	A scalar (integer or real)
$\mathbf{a}$	A vector
$\mathbf{A}$	A matrix
$\mathbf{A}$	A tensor
$\mathbf{I}_n$	Identity matrix with $n$ rows and $n$ columns
$\mathbf{I}$	Identity matrix with dimensionality implied by context
$\mathbf{e}^{(i)}$	Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at position $i$
$\text{diag}(\mathbf{a})$	A square, diagonal matrix with diagonal entries given by $\mathbf{a}$
$a$	A scalar random variable
$\mathbf{a}$	A vector-valued random variable
$\mathbf{A}$	A matrix-valued random variable

## Sets and Graphs

$\mathbb{A}$	A set
$\mathbb{R}$	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and $n$
$[a, b]$	The real interval including $a$ and $b$
$(a, b]$	The real interval excluding $a$ but including $b$
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of $\mathbb{A}$ that are not in $\mathbb{B}$
$\mathcal{G}$	A graph
$Pa_{\mathcal{G}}(x_i)$	The parents of $x_i$ in $\mathcal{G}$

## Indexing

$a_i$	Element $i$ of vector $\mathbf{a}$ , with indexing starting at 1
$a_{-i}$	All elements of vector $\mathbf{a}$ except for element $i$
$A_{i,j}$	Element $i, j$ of matrix $\mathbf{A}$
$\mathbf{A}_{i,:}$	Row $i$ of matrix $\mathbf{A}$
$\mathbf{A}_{:,i}$	Column $i$ of matrix $\mathbf{A}$
$A_{i,j,k}$	Element $(i, j, k)$ of a 3-D tensor $\mathbf{A}$
$\mathbf{A}_{::,i}$	2-D slice of a 3-D tensor
$a_i$	Element $i$ of the random vector $\mathbf{a}$

## Linear Algebra Operations

$\mathbf{A}^\top$	Transpose of matrix $\mathbf{A}$
$\mathbf{A}^+$	Moore-Penrose pseudoinverse of $\mathbf{A}$
$\mathbf{A} \odot \mathbf{B}$	Element-wise (Hadamard) product of $\mathbf{A}$ and $\mathbf{B}$
$\det(\mathbf{A})$	Determinant of $\mathbf{A}$

## Calculus

$\frac{dy}{dx}$	Derivative of $y$ with respect to $x$
$\frac{\partial y}{\partial x}$	Partial derivative of $y$ with respect to $x$
$\nabla_{\mathbf{x}} y$	Gradient of $y$ with respect to $\mathbf{x}$
$\nabla_{\mathbf{X}} y$	Matrix derivatives of $y$ with respect to $\mathbf{X}$
$\nabla_{\mathbf{x}} y$	Tensor containing derivatives of $y$ with respect to $\mathbf{X}$
$\frac{\partial f}{\partial \mathbf{x}}$	Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or $\mathbf{H}(f)(\mathbf{x})$	The Hessian matrix of $f$ at input point $\mathbf{x}$
$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of $\mathbf{x}$
$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to $\mathbf{x}$ over the set $\mathbb{S}$

## Probability and Information Theory

$a \perp b$	The random variables $a$ and $b$ are independent
$a \perp b \mid c$	They are conditionally independent given $c$
$P(a)$	A probability distribution over a discrete variable
$p(a)$	A probability distribution over a continuous variable, or over a variable whose type has not been specified
$a \sim P$	Random variable $a$ has distribution $P$
$\mathbb{E}_{\mathbf{x} \sim P}[f(\mathbf{x})]$ or $\mathbb{E}f(\mathbf{x})$	Expectation of $f(\mathbf{x})$ with respect to $P(\mathbf{x})$
$\text{Var}(f(\mathbf{x}))$	Variance of $f(\mathbf{x})$ under $P(\mathbf{x})$
$\text{Cov}(f(\mathbf{x}), g(\mathbf{x}))$	Covariance of $f(\mathbf{x})$ and $g(\mathbf{x})$ under $P(\mathbf{x})$
$H(\mathbf{x})$	Shannon entropy of the random variable $\mathbf{x}$
$D_{\text{KL}}(P \parallel Q)$	Kullback-Leibler divergence of $P$ and $Q$
$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian distribution over $\mathbf{x}$ with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$

### Functions

$f : \mathbb{A} \rightarrow \mathbb{B}$	The function $f$ with domain $\mathbb{A}$ and range $\mathbb{B}$
$f \circ g$	Composition of the functions $f$ and $g$
$f(\mathbf{x}; \boldsymbol{\theta})$	A function of $\mathbf{x}$ parametrized by $\boldsymbol{\theta}$ . (Sometimes we write $f(\mathbf{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation)
$\log x$	Natural logarithm of $x$
$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
$\ \mathbf{x}\ _p$	$L^p$ norm of $\mathbf{x}$
$\ \mathbf{x}\ $	$L^2$ norm of $\mathbf{x}$
$x^+$	Positive part of $x$ , i.e., $\max(0, x)$
$\mathbf{1}_{\text{condition}}$	is 1 if the condition is true, 0 otherwise

Sometimes we use a function  $f$  whose argument is a scalar but apply it to a vector, matrix, or tensor:  $f(\mathbf{x})$ ,  $f(\mathbf{X})$ , or  $f(\mathbf{X})$ . This denotes the application of  $f$  to the array element-wise. For example, if  $\mathbf{C} = \sigma(\mathbf{X})$ , then  $C_{i,j,k} = \sigma(X_{i,j,k})$  for all valid values of  $i$ ,  $j$  and  $k$ .

### Datasets and Distributions

$p_{\text{data}}$	The data generating distribution
$\hat{p}_{\text{data}}$	The empirical distribution defined by the training set
$\mathbb{X}$	A set of training examples
$\mathbf{x}^{(i)}$	The $i$ -th example (input) from a dataset
$y^{(i)}$ or $\mathbf{y}^{(i)}$	The target associated with $\mathbf{x}^{(i)}$ for supervised learning
$\mathbf{X}$	The $m \times n$ matrix with input example $\mathbf{x}^{(i)}$ in row $\mathbf{X}_{i,:}$

# 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 in terms of its relation to simpler concepts. By gathering knowledge from experience, this approach avoids the need for human operators to formally specify all of the knowledge that the computer needs. The hierarchy of concepts allows the computer to learn complicated concepts by building them out of simpler ones. If we draw a graph showing how these



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

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 are only recently matching some of the abilities of average human beings 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 about statements in these formal languages automatically using logical inference rules. This is known as the **knowledge base** approach to artificial intelligence. None of these projects has led to a 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 "FredWhileShaving" 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 patterns from raw data. This capability is known as **machine learning**. The

introduction of machine learning allowed computers to tackle problems involving knowledge of the real world and make 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 the way that the features are defined in any way. If logistic regression was 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 these features to a simple machine learning algorithm. For example, a useful feature for speaker identification from sound is an estimate of the size of speaker’s vocal tract. It therefore gives a strong clue as to whether the speaker is a man, woman, or child.

However, for many tasks, 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 foreground obscuring part of the wheel, and so on.

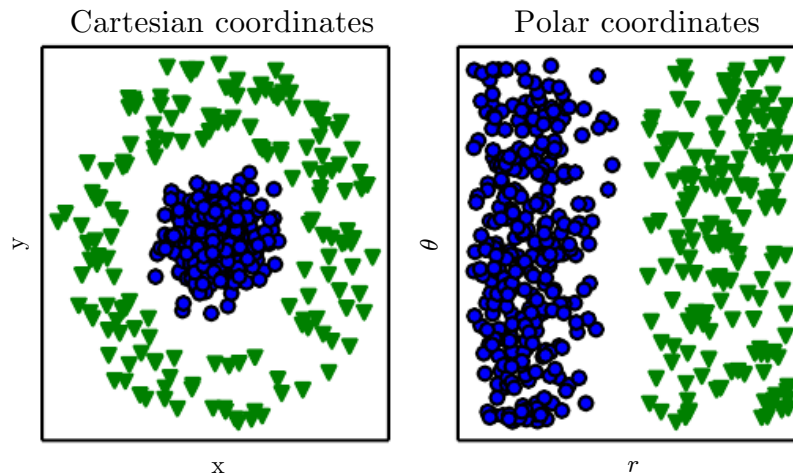


Figure 1.1: Example of different representations: suppose we want to separate two categories of data by drawing a line between them in a scatterplot. In the plot on the left, we represent some data using Cartesian coordinates, and the task is impossible. In the plot on the right, we represent the data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

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 allow 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 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 that converts the input data into a different representation, and a **decoder** function that converts the new representation back into the original format. Autoencoders are trained to preserve as much information as possible when an input is run through the encoder and then the decoder, but 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 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 either as 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 that 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 allows 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 allows 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 executing another set of instructions in parallel. Networks with greater depth can execute more instructions in sequence. Sequential instructions offer great power because later

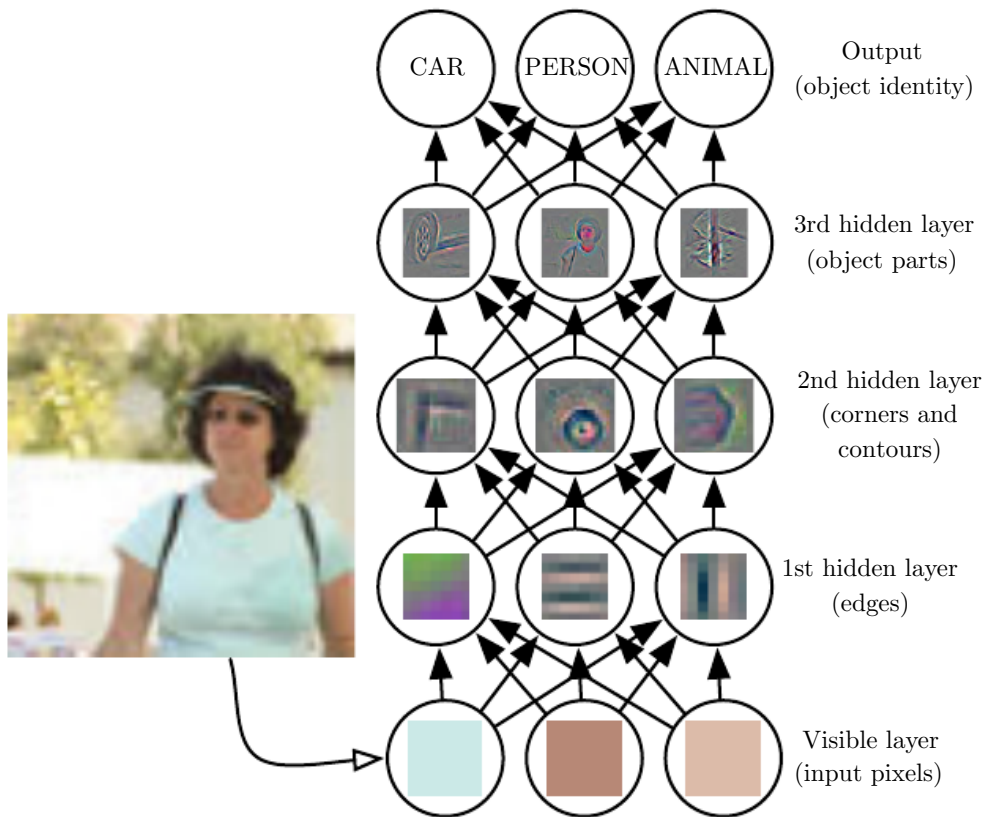


Figure 1.2: Illustration of a deep learning model. It is difficult for a computer to understand the meaning of raw sensory input data, such as this image represented as a collection of pixel values. The function mapping from a set of pixels to an object identity is very complicated. Learning or evaluating 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 \(2014\)](#).

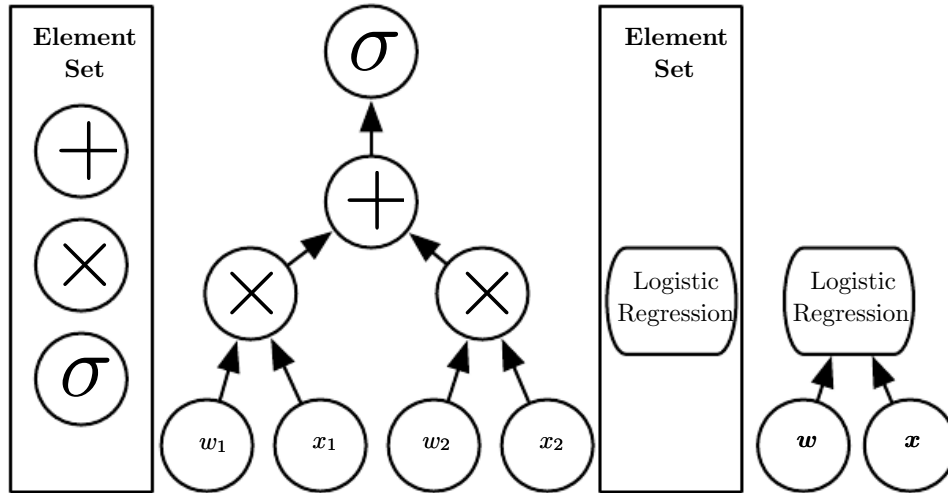


Figure 1.3: Illustration of computational graphs mapping an input to an output where each node performs an operation. Depth is the length of the longest path from input to output but depends on the definition of what constitutes a possible computational step. The computation depicted in these graphs 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 sigmoids 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.

instructions can refer back to the results of earlier instructions. According to this view of deep learning, not all of 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 input 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 flow chart 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.

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 only see one eye. After detecting that a face is present, it can then infer that a second eye is probably present as well. In this case, the graph of concepts only includes 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 smallest 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 safely be regarded as the study of models that either involve a greater amount of composition of 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 allows computer systems to improve with experience and data. According to the authors of this book, 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 learning to represent 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 main target audiences in mind. One of these 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 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 proven useful in

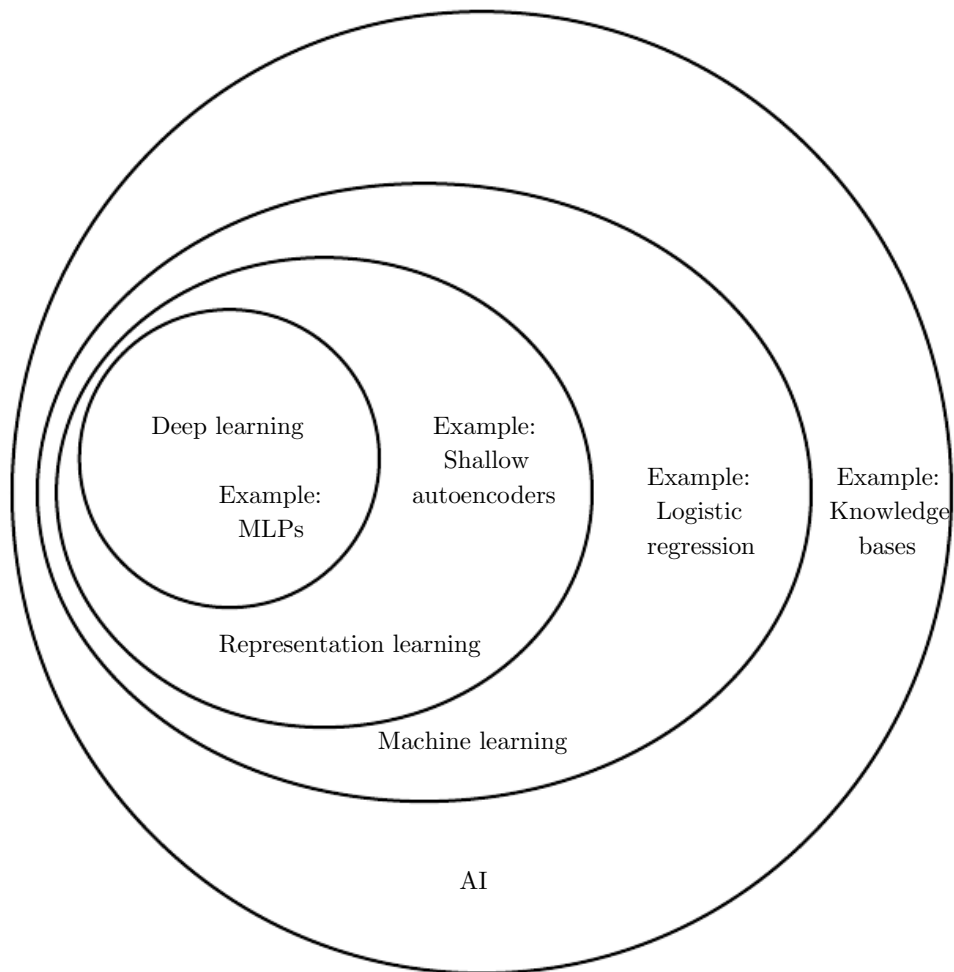


Figure 1.4: A Venn diagram showing how deep learning is a kind of representation learning, which is in turn a kind of machine learning, which is used for many but not all approaches to AI. Each section of the Venn diagram includes an example of an AI technology.



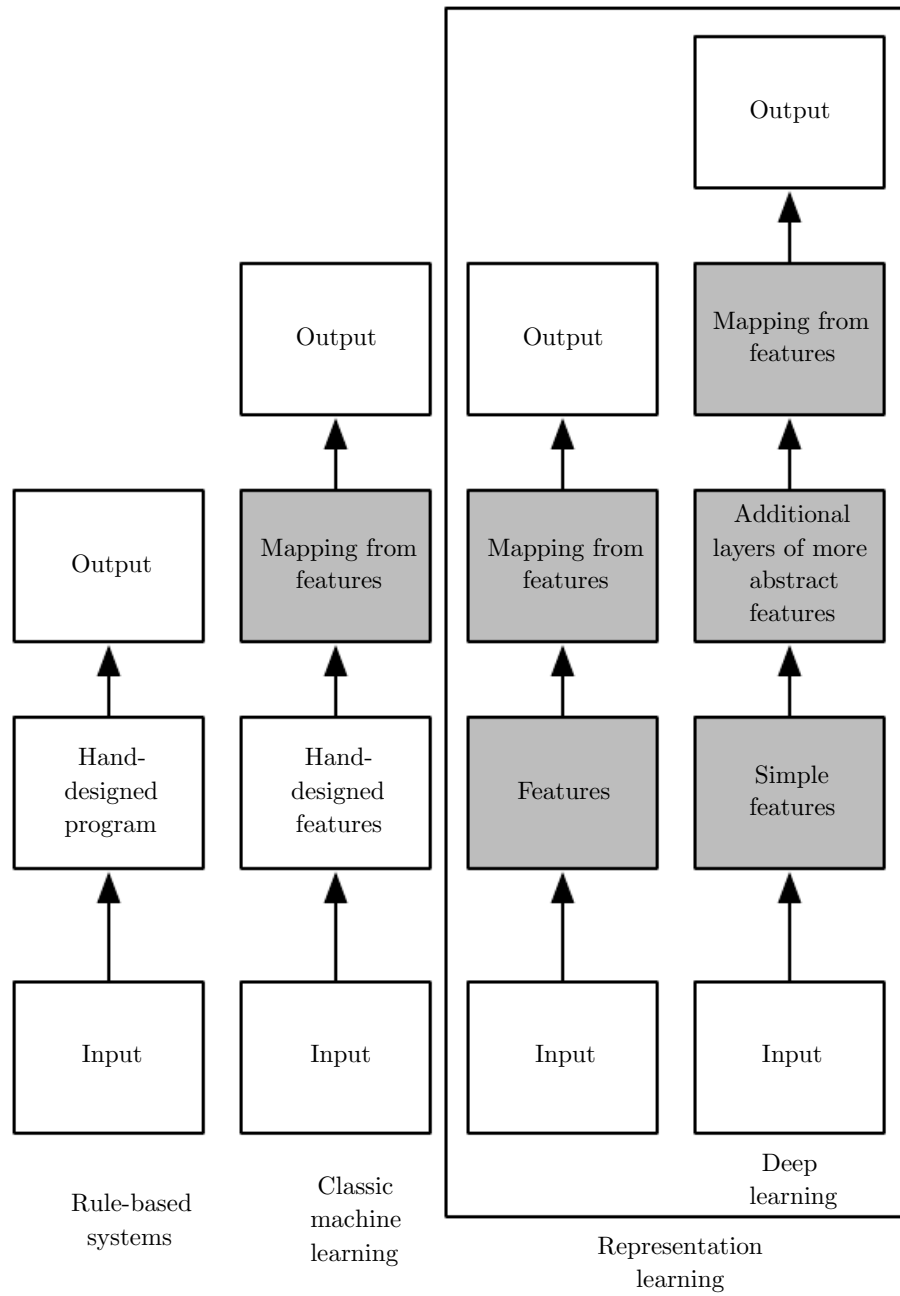


Figure 1.5: Flowcharts showing how the different parts of an AI system relate to each other within different AI disciplines. Shaded boxes indicate components that are able to learn from data.

many software disciplines including computer vision, speech and audio processing, natural language processing, robotics, bioinformatics and chemistry, video games, search engines, online advertising and finance.

This book has been organized into three parts in order 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 that are essentially solved technologies. Part **III** describes more speculative ideas that are widely believed to be important for future research in deep learning.

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 readers who just want to implement a working system need not read beyond part **II**. To help choose which chapters 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.

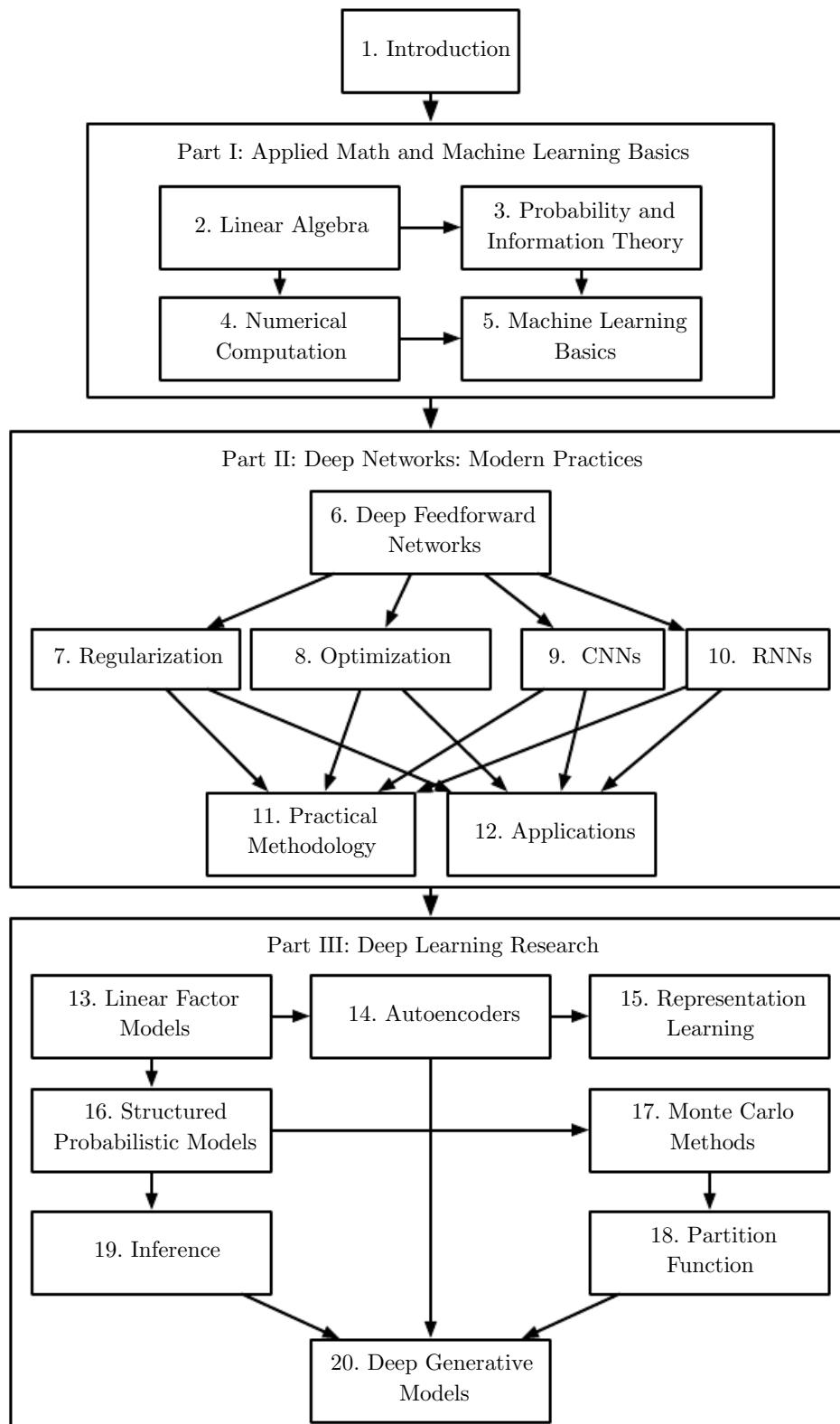


Figure 1.6: The high-level organization of the book. An arrow from one chapter to another indicates that the former chapter is prerequisite material for understanding the latter.

### 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, and has only recently become 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. However, some basic context is useful for understanding deep learning. Broadly speaking, there have been three waves of development of deep learning: deep learning known as **cybernetics** in the 1940s–1960s, deep learning known as **connectionism** in the 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 we recognize today were intended to be computational models of biological learning, i.e. 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 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.

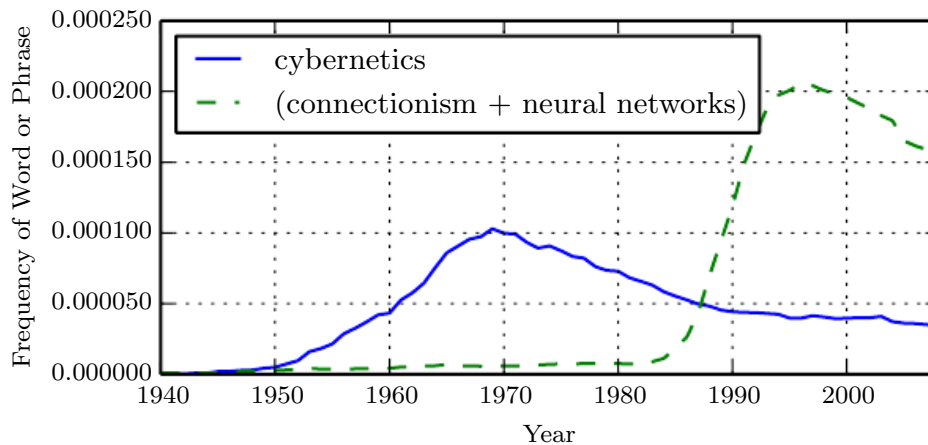


Figure 1.7: The figure shows two of the three historical waves of artificial neural nets research, as measured by the frequency of the phrases “cybernetics” and “connectionism” or “neural networks” according to Google Books (the third wave is too recent to appear). The first wave started with cybernetics 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) allowing 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.

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 an early 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 defining 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 called **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 trained.

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.

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 a reason 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 of these application areas 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.

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 in 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 and 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 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 only from images of one specific category of objects. The concept of distributed representation is central to this book, and will be 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 currently 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 or 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 CIFAR NCAP research initiative had a multi-disciplinary nature that also included neuroscientists and experts in human and computer vision.

At this point in time, 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 pre-training (Hinton *et al.*, 2006), which will be described 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 leverage large labeled datasets.

### 1.2.2 Increasing Dataset Sizes

One may wonder why deep learning has only recently become recognized as a crucial technology 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 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 increased 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” 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 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 not especially densely connected. As seen in figure 1.10, our machine learning models have had a number of connections per neuron that was 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 allow 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. Even today’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,

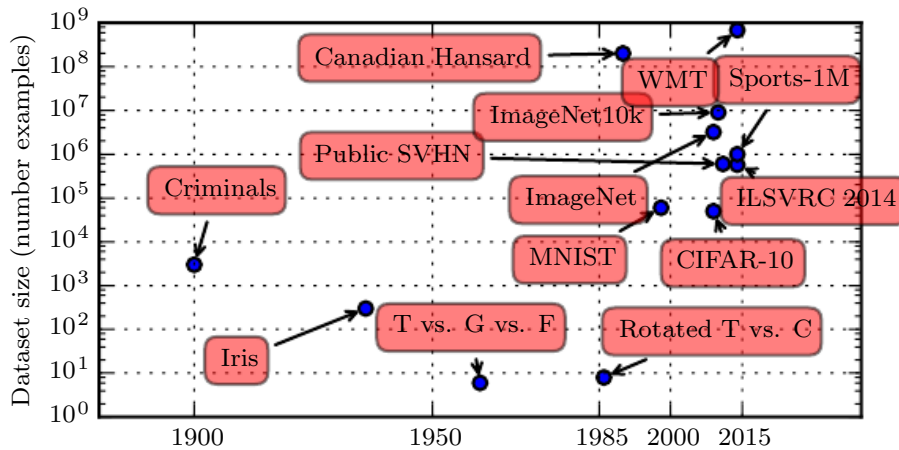


Figure 1.8: Dataset sizes have increased greatly over time. In the early 1900s, statisticians studied datasets using hundreds or thousands of manually compiled measurements (Garson, 1900; Gosset, 1908; Anderson, 1935; Fisher, 1936). In the 1950s through 1980s, the pioneers of biologically inspired machine learning often worked with small, synthetic datasets, such as low-resolution bitmaps of letters, 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 in nature 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 (Karpthy *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 WMT 2014 English to French dataset (Schwenk, 2014) are typically far ahead of other dataset sizes.

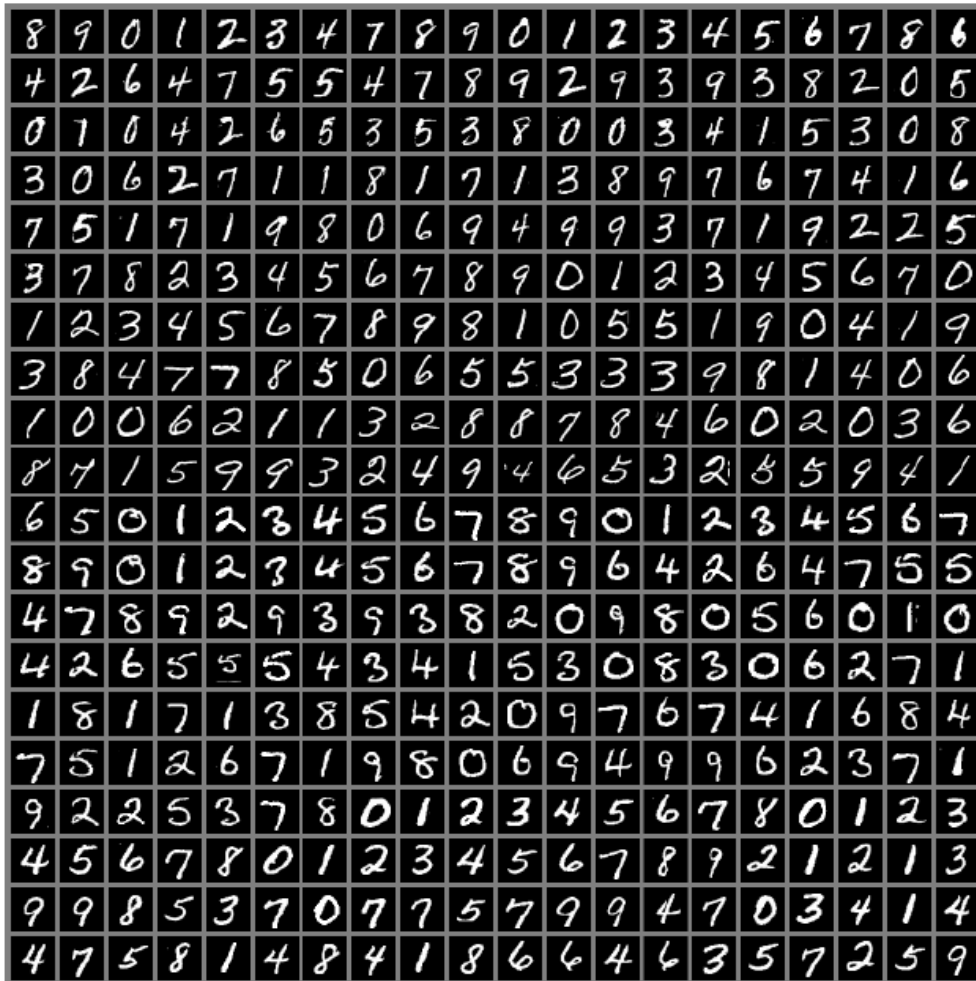


Figure 1.9: Example inputs from the MNIST dataset. The “NIST” stands for National Institute of Standards and Technology, the agency that originally collected this data. The “M” stands for “modified,” since the data has been preprocessed for easier use with machine learning algorithms. The MNIST dataset 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 allows machine learning researchers to study their algorithms in controlled laboratory conditions, much as biologists often study fruit flies.

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.

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

Since the 1980s, deep learning has consistently improved in its ability to provide accurate recognition or 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 do not have a requirement that the photo be cropped near the object to be recognized (Krizhevsky *et al.*, 2012). Similarly, the earliest networks could only recognize 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 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% to 15.3% (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% 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%, 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 of error rates, with some error rates cut in half. We will 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

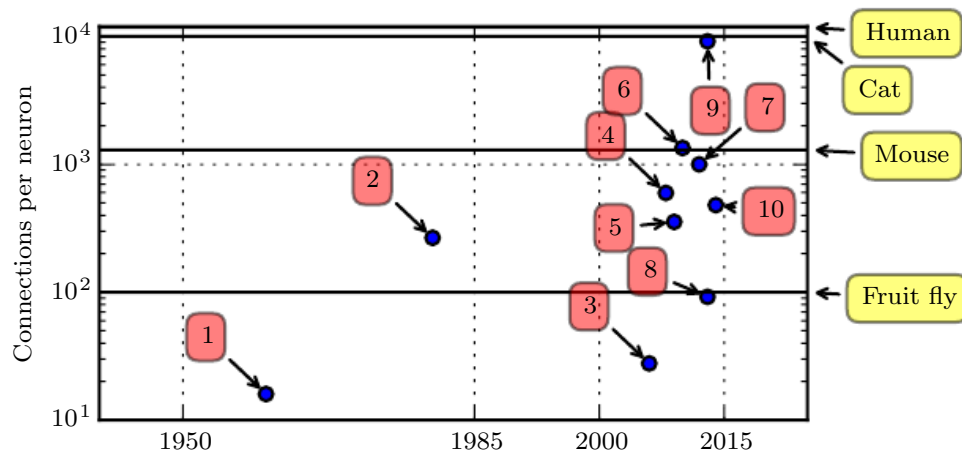


Figure 1.10: Initially, the number of connections between neurons in artificial neural networks was limited by hardware capabilities. Today, the number of connections between neurons is mostly a design consideration. Some artificial neural networks have nearly as many connections per neuron as a cat, and it is quite common for other neural networks to have as many connections per neuron as smaller mammals like mice. Even the human brain does not have an exorbitant amount of connections per neuron. Biological neural network sizes from [Wikipedia \(2015\)](#).

1. Adaptive linear 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](#))



*et al.*, 2012).

At the same time that the scale and accuracy of deep networks has increased, 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.*, 2014a) 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 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 back 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 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, it has seen tremendous growth in its popularity and usefulness, due in large part to 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 bring it to new frontiers.

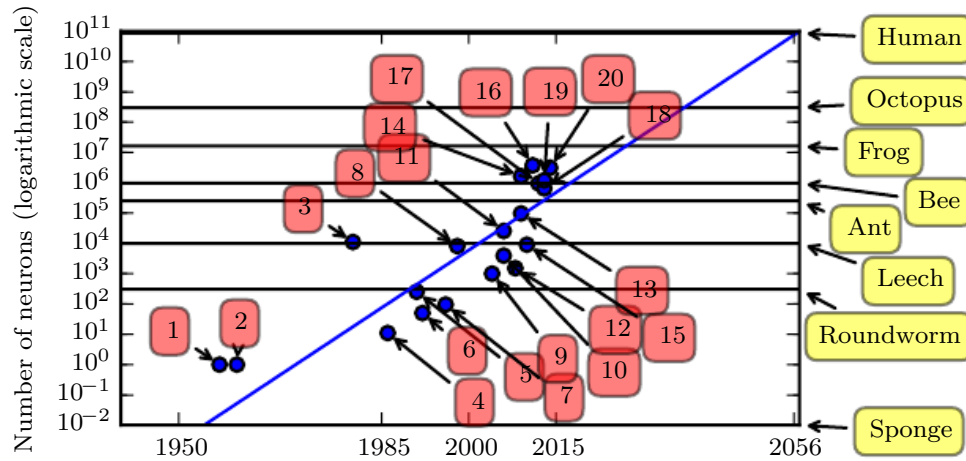


Figure 1.11: Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from [Wikipedia \(2015\)](#).

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. Mean field 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](#))

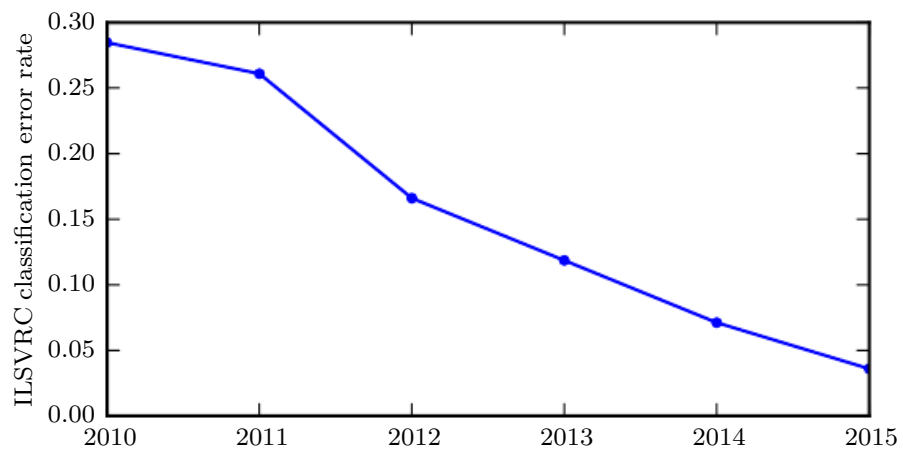


Figure 1.12: Since deep networks reached the scale necessary to compete in the ImageNet Large Scale Visual Recognition Challenge, they have consistently won the competition every year, and yielded lower and lower error rates each time. Data from [Russakovsky \*et al.\* \(2014b\)](#) and [He \*et al.\*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

## Part I

# Applied Math and Machine Learning Basics

---

This part of the book introduces the basic mathematical concepts needed to understand deep learning. We begin with general ideas from applied math that allow us to define functions of many variables, find the highest and lowest points on these functions and quantify degrees of belief.

Next, we describe the fundamental goals of machine learning. We describe how to accomplish these goals by specifying a model that represents certain beliefs, designing a cost function that measures how well those beliefs correspond with reality and using a training algorithm to minimize that cost function.

This elementary framework is the basis for a broad variety of machine learning algorithms, including approaches to machine learning that are not deep. In the subsequent parts of the book, we develop deep learning algorithms within this framework.

## Chapter 2

# Linear Algebra

Linear algebra is a branch of mathematics that is widely used throughout science and engineering. However, because linear algebra is a form of continuous rather than discrete mathematics, many computer scientists have little experience with it. A good understanding of linear algebra is essential for understanding and working with many machine learning algorithms, especially deep learning algorithms. We therefore precede our introduction to deep learning with a focused presentation of the key linear algebra prerequisites.

If you are already familiar with linear algebra, feel free to skip this chapter. If you have previous experience with these concepts but need a detailed reference sheet to review key formulas, we recommend *The Matrix Cookbook* ([Petersen and Pedersen, 2006](#)). If you have no exposure at all to linear algebra, this chapter will teach you enough to read this book, but we highly recommend that you also consult another resource focused exclusively on teaching linear algebra, such as [Shilov \(1977\)](#). This chapter will completely omit many important linear algebra topics that are not essential for understanding deep learning.

## 2.1 Scalars, Vectors, Matrices and Tensors

The study of linear algebra involves several types of mathematical objects:

- **Scalars:** A scalar is just a single number, in contrast to most of the other objects studied in linear algebra, which are usually arrays of multiple numbers. We write scalars in italics. We usually give scalars lower-case variable names. When we introduce them, we specify what kind of number they are. For

example, we might say “Let  $s \in \mathbb{R}$  be the slope of the line,” while defining a real-valued scalar, or “Let  $n \in \mathbb{N}$  be the number of units,” while defining a natural number scalar.

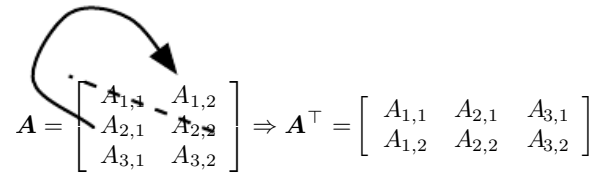
- **Vectors:** A vector is an array of numbers. The numbers are arranged in order. We can identify each individual number by its index in that ordering. Typically we give vectors lower case names written in bold typeface, such as  $\mathbf{x}$ . The elements of the vector are identified by writing its name in italic typeface, with a subscript. The first element of  $\mathbf{x}$  is  $x_1$ , the second element is  $x_2$  and so on. We also need to say what kind of numbers are stored in the vector. If each element is in  $\mathbb{R}$ , and the vector has  $n$  elements, then the vector lies in the set formed by taking the Cartesian product of  $\mathbb{R}$   $n$  times, denoted as  $\mathbb{R}^n$ . When we need to explicitly identify the elements of a vector, we write them as a column enclosed in square brackets:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}. \quad (2.1)$$

We can think of vectors as identifying points in space, with each element giving the coordinate along a different axis.

Sometimes we need to index a set of elements of a vector. In this case, we define a set containing the indices and write the set as a subscript. For example, to access  $x_1$ ,  $x_3$  and  $x_6$ , we define the set  $S = \{1, 3, 6\}$  and write  $\mathbf{x}_S$ . We use the  $-$  sign to index the complement of a set. For example  $\mathbf{x}_{-1}$  is the vector containing all elements of  $\mathbf{x}$  except for  $x_1$ , and  $\mathbf{x}_{-S}$  is the vector containing all of the elements of  $\mathbf{x}$  except for  $x_1$ ,  $x_3$  and  $x_6$ .

- **Matrices:** A matrix is a 2-D array of numbers, so each element is identified by two indices instead of just one. We usually give matrices upper-case variable names with bold typeface, such as  $\mathbf{A}$ . If a real-valued matrix  $\mathbf{A}$  has a height of  $m$  and a width of  $n$ , then we say that  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . We usually identify the elements of a matrix using its name in italic but not bold font, and the indices are listed with separating commas. For example,  $A_{1,1}$  is the upper left entry of  $\mathbf{A}$  and  $A_{m,n}$  is the bottom right entry. We can identify all of the numbers with vertical coordinate  $i$  by writing a “:” for the horizontal coordinate. For example,  $\mathbf{A}_{i,:}$  denotes the horizontal cross section of  $\mathbf{A}$  with vertical coordinate  $i$ . This is known as the  $i$ -th **row** of  $\mathbf{A}$ . Likewise,  $\mathbf{A}_{:,i}$  is



$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow \mathbf{A}^\top = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

Figure 2.1: The transpose of the matrix can be thought of as a mirror image across the main diagonal.

the  $i$ -th **column** of  $\mathbf{A}$ . When we need to explicitly identify the elements of a matrix, we write them as an array enclosed in square brackets:

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}. \quad (2.2)$$

Sometimes we may need to index matrix-valued expressions that are not just a single letter. In this case, we use subscripts after the expression, but do not convert anything to lower case. For example,  $f(\mathbf{A})_{i,j}$  gives element  $(i, j)$  of the matrix computed by applying the function  $f$  to  $\mathbf{A}$ .

- **Tensors:** In some cases we will need an array with more than two axes. In the general case, an array of numbers arranged on a regular grid with a variable number of axes is known as a tensor. We denote a tensor named “A” with this typeface:  $\mathbf{A}$ . We identify the element of  $\mathbf{A}$  at coordinates  $(i, j, k)$  by writing  $A_{i,j,k}$ .

One important operation on matrices is the **transpose**. The transpose of a matrix is the mirror image of the matrix across a diagonal line, called the **main diagonal**, running down and to the right, starting from its upper left corner. See figure 2.1 for a graphical depiction of this operation. We denote the transpose of a matrix  $\mathbf{A}$  as  $\mathbf{A}^\top$ , and it is defined such that

$$(\mathbf{A}^\top)_{i,j} = A_{j,i}. \quad (2.3)$$

Vectors can be thought of as matrices that contain only one column. The transpose of a vector is therefore a matrix with only one row. Sometimes we



define a vector by writing out its elements in the text inline as a row matrix, then using the transpose operator to turn it into a standard column vector, e.g.,  $\mathbf{x} = [x_1, x_2, x_3]^\top$ .

A scalar can be thought of as a matrix with only a single entry. From this, we can see that a scalar is its own transpose:  $a = a^\top$ .

We can add matrices to each other, as long as they have the same shape, just by adding their corresponding elements:  $\mathbf{C} = \mathbf{A} + \mathbf{B}$  where  $C_{i,j} = A_{i,j} + B_{i,j}$ .

We can also add a scalar to a matrix or multiply a matrix by a scalar, just by performing that operation on each element of a matrix:  $\mathbf{D} = a \cdot \mathbf{B} + c$  where  $D_{i,j} = a \cdot B_{i,j} + c$ .

In the context of deep learning, we also use some less conventional notation. We allow the addition of matrix and a vector, yielding another matrix:  $\mathbf{C} = \mathbf{A} + \mathbf{b}$ , where  $C_{i,j} = A_{i,j} + b_j$ . In other words, the vector  $\mathbf{b}$  is added to each row of the matrix. This shorthand eliminates the need to define a matrix with  $\mathbf{b}$  copied into each row before doing the addition. This implicit copying of  $\mathbf{b}$  to many locations is called **broadcasting**.

## 2.2 Multiplying Matrices and Vectors

One of the most important operations involving matrices is multiplication of two matrices. The **matrix product** of matrices  $\mathbf{A}$  and  $\mathbf{B}$  is a third matrix  $\mathbf{C}$ . In order for this product to be defined,  $\mathbf{A}$  must have the same number of columns as  $\mathbf{B}$  has rows. If  $\mathbf{A}$  is of shape  $m \times n$  and  $\mathbf{B}$  is of shape  $n \times p$ , then  $\mathbf{C}$  is of shape  $m \times p$ . We can write the matrix product just by placing two or more matrices together, e.g.

$$\mathbf{C} = \mathbf{AB}. \quad (2.4)$$

The product operation is defined by

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}. \quad (2.5)$$

Note that the standard product of two matrices is *not* just a matrix containing the product of the individual elements. Such an operation exists and is called the **element-wise product** or **Hadamard product**, and is denoted as  $\mathbf{A} \odot \mathbf{B}$ .

The **dot product** between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  of the same dimensionality is the matrix product  $\mathbf{x}^\top \mathbf{y}$ . We can think of the matrix product  $\mathbf{C} = \mathbf{AB}$  as computing  $C_{i,j}$  as the dot product between row  $i$  of  $\mathbf{A}$  and column  $j$  of  $\mathbf{B}$ .

Matrix product operations have many useful properties that make mathematical analysis of matrices more convenient. For example, matrix multiplication is distributive:

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}. \quad (2.6)$$

It is also associative:

$$\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}. \quad (2.7)$$

Matrix multiplication is *not* commutative (the condition  $\mathbf{AB} = \mathbf{BA}$  does not always hold), unlike scalar multiplication. However, the dot product between two vectors is commutative:

$$\mathbf{x}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{x}. \quad (2.8)$$

The transpose of a matrix product has a simple form:

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top. \quad (2.9)$$

This allows us to demonstrate equation 2.8, by exploiting the fact that the value of such a product is a scalar and therefore equal to its own transpose:

$$\mathbf{x}^\top \mathbf{y} = \left( \mathbf{x}^\top \mathbf{y} \right)^\top = \mathbf{y}^\top \mathbf{x}. \quad (2.10)$$

Since the focus of this textbook is not linear algebra, we do not attempt to develop a comprehensive list of useful properties of the matrix product here, but the reader should be aware that many more exist.

We now know enough linear algebra notation to write down a system of linear equations:

$$\mathbf{Ax} = \mathbf{b} \quad (2.11)$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a known matrix,  $\mathbf{b} \in \mathbb{R}^m$  is a known vector, and  $\mathbf{x} \in \mathbb{R}^n$  is a vector of unknown variables we would like to solve for. Each element  $x_i$  of  $\mathbf{x}$  is one of these unknown variables. Each row of  $\mathbf{A}$  and each element of  $\mathbf{b}$  provide another constraint. We can rewrite equation 2.11 as:

$$\mathbf{A}_{1,:}\mathbf{x} = b_1 \quad (2.12)$$

$$\mathbf{A}_{2,:}\mathbf{x} = b_2 \quad (2.13)$$

$$\dots \quad (2.14)$$

$$\mathbf{A}_{m,:}\mathbf{x} = b_m \quad (2.15)$$

or, even more explicitly, as:

$$\mathbf{A}_{1,1}x_1 + \mathbf{A}_{1,2}x_2 + \dots + \mathbf{A}_{1,n}x_n = b_1 \quad (2.16)$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Figure 2.2: **Example identity matrix:** This is  $\mathbf{I}_3$ .

$$\mathbf{A}_{2,1}x_1 + \mathbf{A}_{2,2}x_2 + \cdots + \mathbf{A}_{2,n}x_n = b_2 \quad (2.17)$$

$$\dots \quad (2.18)$$

$$\mathbf{A}_{m,1}x_1 + \mathbf{A}_{m,2}x_2 + \cdots + \mathbf{A}_{m,n}x_n = b_m. \quad (2.19)$$

Matrix-vector product notation provides a more compact representation for equations of this form.

## 2.3 Identity and Inverse Matrices

Linear algebra offers a powerful tool called **matrix inversion** that allows us to analytically solve equation 2.11 for many values of  $\mathbf{A}$ .

To describe matrix inversion, we first need to define the concept of an **identity matrix**. An identity matrix is a matrix that does not change any vector when we multiply that vector by that matrix. We denote the identity matrix that preserves  $n$ -dimensional vectors as  $\mathbf{I}_n$ . Formally,  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ , and

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{I}_n \mathbf{x} = \mathbf{x}. \quad (2.20)$$

The structure of the identity matrix is simple: all of the entries along the main diagonal are 1, while all of the other entries are zero. See figure 2.2 for an example.

The **matrix inverse** of  $\mathbf{A}$  is denoted as  $\mathbf{A}^{-1}$ , and it is defined as the matrix such that

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}_n. \quad (2.21)$$

We can now solve equation 2.11 by the following steps:

$$\mathbf{A} \mathbf{x} = \mathbf{b} \quad (2.22)$$

$$\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.23)$$

$$\mathbf{I}_n \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.24)$$

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}. \quad (2.25)$$

Of course, this process depends on it being possible to find  $\mathbf{A}^{-1}$ . We discuss the conditions for the existence of  $\mathbf{A}^{-1}$  in the following section.

When  $\mathbf{A}^{-1}$  exists, several different algorithms exist for finding it in closed form. In theory, the same inverse matrix can then be used to solve the equation many times for different values of  $\mathbf{b}$ . However,  $\mathbf{A}^{-1}$  is primarily useful as a theoretical tool, and should not actually be used in practice for most software applications. Because  $\mathbf{A}^{-1}$  can be represented with only limited precision on a digital computer, algorithms that make use of the value of  $\mathbf{b}$  can usually obtain more accurate estimates of  $\mathbf{x}$ .

## 2.4 Linear Dependence and Span

In order for  $\mathbf{A}^{-1}$  to exist, equation 2.11 must have exactly one solution for every value of  $\mathbf{b}$ . However, it is also possible for the system of equations to have no solutions or infinitely many solutions for some values of  $\mathbf{b}$ . It is not possible to have more than one but less than infinitely many solutions for a particular  $\mathbf{b}$ ; if both  $\mathbf{x}$  and  $\mathbf{y}$  are solutions then

$$\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \quad (2.26)$$

is also a solution for any real  $\alpha$ .

To analyze how many solutions the equation has, we can think of the columns of  $\mathbf{A}$  as specifying different directions we can travel from the **origin** (the point specified by the vector of all zeros), and determine how many ways there are of reaching  $\mathbf{b}$ . In this view, each element of  $\mathbf{x}$  specifies how far we should travel in each of these directions, with  $x_i$  specifying how far to move in the direction of column  $i$ :

$$\mathbf{Ax} = \sum_i x_i \mathbf{A}_{:,i}. \quad (2.27)$$

In general, this kind of operation is called a **linear combination**. Formally, a linear combination of some set of vectors  $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$  is given by multiplying each vector  $\mathbf{v}^{(i)}$  by a corresponding scalar coefficient and adding the results:

$$\sum_i c_i \mathbf{v}^{(i)}. \quad (2.28)$$

The **span** of a set of vectors is the set of all points obtainable by linear combination of the original vectors.

Determining whether  $\mathbf{Ax} = \mathbf{b}$  has a solution thus amounts to testing whether  $\mathbf{b}$  is in the span of the columns of  $\mathbf{A}$ . This particular span is known as the **column space** or the **range** of  $\mathbf{A}$ .

In order for the system  $\mathbf{Ax} = \mathbf{b}$  to have a solution for all values of  $\mathbf{b} \in \mathbb{R}^m$ , we therefore require that the column space of  $\mathbf{A}$  be all of  $\mathbb{R}^m$ . If any point in  $\mathbb{R}^m$  is excluded from the column space, that point is a potential value of  $\mathbf{b}$  that has no solution. The requirement that the column space of  $\mathbf{A}$  be all of  $\mathbb{R}^m$  implies immediately that  $\mathbf{A}$  must have at least  $m$  columns, i.e.,  $n \geq m$ . Otherwise, the dimensionality of the column space would be less than  $m$ . For example, consider a  $3 \times 2$  matrix. The target  $\mathbf{b}$  is 3-D, but  $\mathbf{x}$  is only 2-D, so modifying the value of  $\mathbf{x}$  at best allows us to trace out a 2-D plane within  $\mathbb{R}^3$ . The equation has a solution if and only if  $\mathbf{b}$  lies on that plane.

Having  $n \geq m$  is only a necessary condition for every point to have a solution. It is not a sufficient condition, because it is possible for some of the columns to be redundant. Consider a  $2 \times 2$  matrix where both of the columns are identical. This has the same column space as a  $2 \times 1$  matrix containing only one copy of the replicated column. In other words, the column space is still just a line, and fails to encompass all of  $\mathbb{R}^2$ , even though there are two columns.

Formally, this kind of redundancy is known as **linear dependence**. A set of vectors is **linearly independent** if no vector in the set is a linear combination of the other vectors. If we add a vector to a set that is a linear combination of the other vectors in the set, the new vector does not add any points to the set's span. This means that for the column space of the matrix to encompass all of  $\mathbb{R}^m$ , the matrix must contain at least one set of  $m$  linearly independent columns. This condition is both necessary and sufficient for equation 2.11 to have a solution for every value of  $\mathbf{b}$ . Note that the requirement is for a set to have exactly  $m$  linear independent columns, not at least  $m$ . No set of  $m$ -dimensional vectors can have more than  $m$  mutually linearly independent columns, but a matrix with more than  $m$  columns may have more than one such set.

In order for the matrix to have an inverse, we additionally need to ensure that equation 2.11 has *at most* one solution for each value of  $\mathbf{b}$ . To do so, we need to ensure that the matrix has at most  $m$  columns. Otherwise there is more than one way of parametrizing each solution.

Together, this means that the matrix must be **square**, that is, we require that  $m = n$  and that all of the columns must be linearly independent. A square matrix with linearly dependent columns is known as **singular**.

If  $\mathbf{A}$  is not square or is square but singular, it can still be possible to solve the equation. However, we can not use the method of matrix inversion to find the

solution.

So far we have discussed matrix inverses as being multiplied on the left. It is also possible to define an inverse that is multiplied on the right:

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}. \quad (2.29)$$

For square matrices, the left inverse and right inverse are equal.

## 2.5 Norms

Sometimes we need to measure the size of a vector. In machine learning, we usually measure the size of vectors using a function called a **norm**. Formally, the  $L^p$  norm is given by

$$\|\mathbf{x}\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}} \quad (2.30)$$

for  $p \in \mathbb{R}, p \geq 1$ .

Norms, including the  $L^p$  norm, are functions mapping vectors to non-negative values. On an intuitive level, the norm of a vector  $\mathbf{x}$  measures the distance from the origin to the point  $\mathbf{x}$ . More rigorously, a norm is any function  $f$  that satisfies the following properties:

- $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$  (the **triangle inequality**)
- $\forall \alpha \in \mathbb{R}, f(\alpha\mathbf{x}) = |\alpha|f(\mathbf{x})$

The  $L^2$  norm, with  $p = 2$ , is known as the **Euclidean norm**. It is simply the Euclidean distance from the origin to the point identified by  $\mathbf{x}$ . The  $L^2$  norm is used so frequently in machine learning that it is often denoted simply as  $\|\mathbf{x}\|$ , with the subscript 2 omitted. It is also common to measure the size of a vector using the squared  $L^2$  norm, which can be calculated simply as  $\mathbf{x}^\top \mathbf{x}$ .

The squared  $L^2$  norm is more convenient to work with mathematically and computationally than the  $L^2$  norm itself. For example, the derivatives of the squared  $L^2$  norm with respect to each element of  $\mathbf{x}$  each depend only on the corresponding element of  $\mathbf{x}$ , while all of the derivatives of the  $L^2$  norm depend on the entire vector. In many contexts, the squared  $L^2$  norm may be undesirable because it increases very slowly near the origin. In several machine learning

applications, it is important to discriminate between elements that are exactly zero and elements that are small but nonzero. In these cases, we turn to a function that grows at the same rate in all locations, but retains mathematical simplicity: the  $L^1$  norm. The  $L^1$  norm may be simplified to

$$\|\mathbf{x}\|_1 = \sum_i |x_i|. \quad (2.31)$$

The  $L^1$  norm is commonly used in machine learning when the difference between zero and nonzero elements is very important. Every time an element of  $\mathbf{x}$  moves away from 0 by  $\epsilon$ , the  $L^1$  norm increases by  $\epsilon$ .

We sometimes measure the size of the vector by counting its number of nonzero elements. Some authors refer to this function as the “ $L^0$  norm,” but this is incorrect terminology. The number of non-zero entries in a vector is not a norm, because scaling the vector by  $\alpha$  does not change the number of nonzero entries. The  $L^1$  norm is often used as a substitute for the number of nonzero entries.

One other norm that commonly arises in machine learning is the  $L^\infty$  norm, also known as the **max norm**. This norm simplifies to the absolute value of the element with the largest magnitude in the vector,

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (2.32)$$

Sometimes we may also wish to measure the size of a matrix. In the context of deep learning, the most common way to do this is with the otherwise obscure **Frobenius norm**:

$$\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}, \quad (2.33)$$

which is analogous to the  $L^2$  norm of a vector.

The dot product of two vectors can be rewritten in terms of norms. Specifically,

$$\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta \quad (2.34)$$

where  $\theta$  is the angle between  $\mathbf{x}$  and  $\mathbf{y}$ .

## 2.6 Special Kinds of Matrices and Vectors

Some special kinds of matrices and vectors are particularly useful.

**Diagonal** matrices consist mostly of zeros and have non-zero entries only along the main diagonal. Formally, a matrix  $\mathbf{D}$  is diagonal if and only if  $D_{i,j} = 0$  for

all  $i \neq j$ . We have already seen one example of a diagonal matrix: the identity matrix, where all of the diagonal entries are 1. We write  $\text{diag}(\mathbf{v})$  to denote a square diagonal matrix whose diagonal entries are given by the entries of the vector  $\mathbf{v}$ . Diagonal matrices are of interest in part because multiplying by a diagonal matrix is very computationally efficient. To compute  $\text{diag}(\mathbf{v})\mathbf{x}$ , we only need to scale each element  $x_i$  by  $v_i$ . In other words,  $\text{diag}(\mathbf{v})\mathbf{x} = \mathbf{v} \odot \mathbf{x}$ . Inverting a square diagonal matrix is also efficient. The inverse exists only if every diagonal entry is nonzero, and in that case,  $\text{diag}(\mathbf{v})^{-1} = \text{diag}([1/v_1, \dots, 1/v_n]^\top)$ . In many cases, we may derive some very general machine learning algorithm in terms of arbitrary matrices, but obtain a less expensive (and less descriptive) algorithm by restricting some matrices to be diagonal.

Not all diagonal matrices need be square. It is possible to construct a rectangular diagonal matrix. Non-square diagonal matrices do not have inverses but it is still possible to multiply by them cheaply. For a non-square diagonal matrix  $\mathbf{D}$ , the product  $\mathbf{D}\mathbf{x}$  will involve scaling each element of  $\mathbf{x}$ , and either concatenating some zeros to the result if  $\mathbf{D}$  is taller than it is wide, or discarding some of the last elements of the vector if  $\mathbf{D}$  is wider than it is tall.

A **symmetric** matrix is any matrix that is equal to its own transpose:

$$\mathbf{A} = \mathbf{A}^\top. \quad (2.35)$$

Symmetric matrices often arise when the entries are generated by some function of two arguments that does not depend on the order of the arguments. For example, if  $\mathbf{A}$  is a matrix of distance measurements, with  $\mathbf{A}_{i,j}$  giving the distance from point  $i$  to point  $j$ , then  $\mathbf{A}_{i,j} = \mathbf{A}_{j,i}$  because distance functions are symmetric.

A **unit vector** is a vector with **unit norm**:

$$\|\mathbf{x}\|_2 = 1. \quad (2.36)$$

A vector  $\mathbf{x}$  and a vector  $\mathbf{y}$  are **orthogonal** to each other if  $\mathbf{x}^\top \mathbf{y} = 0$ . If both vectors have nonzero norm, this means that they are at a 90 degree angle to each other. In  $\mathbb{R}^n$ , at most  $n$  vectors may be mutually orthogonal with nonzero norm. If the vectors are not only orthogonal but also have unit norm, we call them **orthonormal**.

An **orthogonal matrix** is a square matrix whose rows are mutually orthonormal and whose columns are mutually orthonormal:

$$\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I}. \quad (2.37)$$



This implies that

$$\mathbf{A}^{-1} = \mathbf{A}^\top, \quad (2.38)$$

so orthogonal matrices are of interest because their inverse is very cheap to compute. Pay careful attention to the definition of orthogonal matrices. Counterintuitively, their rows are not merely orthogonal but fully orthonormal. There is no special term for a matrix whose rows or columns are orthogonal but not orthonormal.

## 2.7 Eigendecomposition

Many mathematical objects can be understood better by breaking them into constituent parts, or finding some properties of them that are universal, not caused by the way we choose to represent them.

For example, integers can be decomposed into prime factors. The way we represent the number 12 will change depending on whether we write it in base ten or in binary, but it will always be true that  $12 = 2 \times 2 \times 3$ . From this representation we can conclude useful properties, such as that 12 is not divisible by 5, or that any integer multiple of 12 will be divisible by 3.

Much as we can discover something about the true nature of an integer by decomposing it into prime factors, we can also decompose matrices in ways that show us information about their functional properties that is not obvious from the representation of the matrix as an array of elements.

One of the most widely used kinds of matrix decomposition is called **eigendecomposition**, in which we decompose a matrix into a set of eigenvectors and eigenvalues.

An **eigenvector** of a square matrix  $\mathbf{A}$  is a non-zero vector  $\mathbf{v}$  such that multiplication by  $\mathbf{A}$  alters only the scale of  $\mathbf{v}$ :

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}. \quad (2.39)$$

The scalar  $\lambda$  is known as the **eigenvalue** corresponding to this eigenvector. (One can also find a **left eigenvector** such that  $\mathbf{v}^\top \mathbf{A} = \lambda \mathbf{v}^\top$ , but we are usually concerned with right eigenvectors).

If  $\mathbf{v}$  is an eigenvector of  $\mathbf{A}$ , then so is any rescaled vector  $s\mathbf{v}$  for  $s \in \mathbb{R}, s \neq 0$ . Moreover,  $s\mathbf{v}$  still has the same eigenvalue. For this reason, we usually only look for unit eigenvectors.

Suppose that a matrix  $\mathbf{A}$  has  $n$  linearly independent eigenvectors,  $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$ , with corresponding eigenvalues  $\{\lambda_1, \dots, \lambda_n\}$ . We may concatenate all of the

## Effect of eigenvectors and eigenvalues

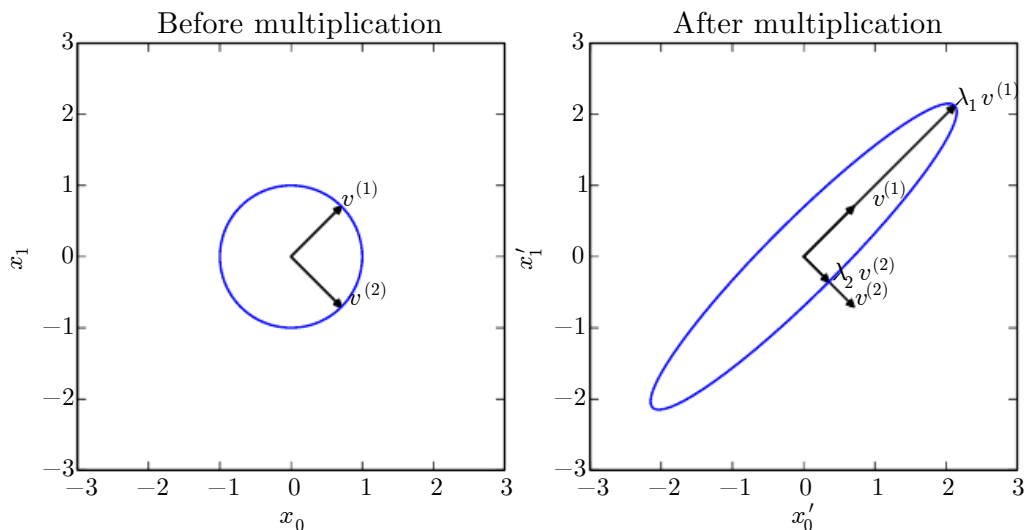


Figure 2.3: An example of the effect of eigenvectors and eigenvalues. Here, we have a matrix  $\mathbf{A}$  with two orthonormal eigenvectors,  $\mathbf{v}^{(1)}$  with eigenvalue  $\lambda_1$  and  $\mathbf{v}^{(2)}$  with eigenvalue  $\lambda_2$ . (Left) We plot the set of all unit vectors  $\mathbf{u} \in \mathbb{R}^2$  as a unit circle. (Right) We plot the set of all points  $\mathbf{A}\mathbf{u}$ . By observing the way that  $\mathbf{A}$  distorts the unit circle, we can see that it scales space in direction  $\mathbf{v}^{(i)}$  by  $\lambda_i$ .

eigenvectors to form a matrix  $\mathbf{V}$  with one eigenvector per column:  $\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}]$ . Likewise, we can concatenate the eigenvalues to form a vector  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^\top$ . The **eigendecomposition** of  $\mathbf{A}$  is then given by

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (2.40)$$

We have seen that *constructing* matrices with specific eigenvalues and eigenvectors allows us to stretch space in desired directions. However, we often want to **decompose** matrices into their eigenvalues and eigenvectors. Doing so can help us to analyze certain properties of the matrix, much as decomposing an integer into its prime factors can help us understand the behavior of that integer.

Not every matrix can be decomposed into eigenvalues and eigenvectors. In some

cases, the decomposition exists, but may involve complex rather than real numbers. Fortunately, in this book, we usually need to decompose only a specific class of matrices that have a simple decomposition. Specifically, every real symmetric matrix can be decomposed into an expression using only real-valued eigenvectors and eigenvalues:

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top, \quad (2.41)$$

where  $\mathbf{Q}$  is an orthogonal matrix composed of eigenvectors of  $\mathbf{A}$ , and  $\mathbf{\Lambda}$  is a diagonal matrix. The eigenvalue  $\Lambda_{i,i}$  is associated with the eigenvector in column  $i$  of  $\mathbf{Q}$ , denoted as  $\mathbf{Q}_{:,i}$ . Because  $\mathbf{Q}$  is an orthogonal matrix, we can think of  $\mathbf{A}$  as scaling space by  $\lambda_i$  in direction  $\mathbf{v}^{(i)}$ . See figure 2.3 for an example.

While any real symmetric matrix  $\mathbf{A}$  is guaranteed to have an eigendecomposition, the eigendecomposition may not be unique. If any two or more eigenvectors share the same eigenvalue, then any set of orthogonal vectors lying in their span are also eigenvectors with that eigenvalue, and we could equivalently choose a  $\mathbf{Q}$  using those eigenvectors instead. By convention, we usually sort the entries of  $\mathbf{\Lambda}$  in descending order. Under this convention, the eigendecomposition is unique only if all of the eigenvalues are unique.

The eigendecomposition of a matrix tells us many useful facts about the matrix. The matrix is singular if and only if any of the eigenvalues are zero. The eigendecomposition of a real symmetric matrix can also be used to optimize quadratic expressions of the form  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A}\mathbf{x}$  subject to  $\|\mathbf{x}\|_2 = 1$ . Whenever  $\mathbf{x}$  is equal to an eigenvector of  $\mathbf{A}$ ,  $f$  takes on the value of the corresponding eigenvalue. The maximum value of  $f$  within the constraint region is the maximum eigenvalue and its minimum value within the constraint region is the minimum eigenvalue.

A matrix whose eigenvalues are all positive is called **positive definite**. A matrix whose eigenvalues are all positive or zero-valued is called **positive semidefinite**. Likewise, if all eigenvalues are negative, the matrix is **negative definite**, and if all eigenvalues are negative or zero-valued, it is **negative semidefinite**. Positive semidefinite matrices are interesting because they guarantee that  $\forall \mathbf{x}, \mathbf{x}^\top \mathbf{A}\mathbf{x} \geq 0$ . Positive definite matrices additionally guarantee that  $\mathbf{x}^\top \mathbf{A}\mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$ .

## 2.8 Singular Value Decomposition

In section 2.7, we saw how to decompose a matrix into eigenvectors and eigenvalues. The **singular value decomposition** (SVD) provides another way to factorize a matrix, into **singular vectors** and **singular values**. The SVD allows us to discover some of the same kind of information as the eigendecomposition. However,

the SVD is more generally applicable. Every real matrix has a singular value decomposition, but the same is not true of the eigenvalue decomposition. For example, if a matrix is not square, the eigendecomposition is not defined, and we must use a singular value decomposition instead.

Recall that the eigendecomposition involves analyzing a matrix  $\mathbf{A}$  to discover a matrix  $\mathbf{V}$  of eigenvectors and a vector of eigenvalues  $\boldsymbol{\lambda}$  such that we can rewrite  $\mathbf{A}$  as

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (2.42)$$

The singular value decomposition is similar, except this time we will write  $\mathbf{A}$  as a product of three matrices:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top. \quad (2.43)$$

Suppose that  $\mathbf{A}$  is an  $m \times n$  matrix. Then  $\mathbf{U}$  is defined to be an  $m \times m$  matrix,  $\mathbf{D}$  to be an  $m \times n$  matrix, and  $\mathbf{V}$  to be an  $n \times n$  matrix.

Each of these matrices is defined to have a special structure. The matrices  $\mathbf{U}$  and  $\mathbf{V}$  are both defined to be orthogonal matrices. The matrix  $\mathbf{D}$  is defined to be a diagonal matrix. Note that  $\mathbf{D}$  is not necessarily square.

The elements along the diagonal of  $\mathbf{D}$  are known as the **singular values** of the matrix  $\mathbf{A}$ . The columns of  $\mathbf{U}$  are known as the **left-singular vectors**. The columns of  $\mathbf{V}$  are known as the **right-singular vectors**.

We can actually interpret the singular value decomposition of  $\mathbf{A}$  in terms of the eigendecomposition of functions of  $\mathbf{A}$ . The left-singular vectors of  $\mathbf{A}$  are the eigenvectors of  $\mathbf{A} \mathbf{A}^\top$ . The right-singular vectors of  $\mathbf{A}$  are the eigenvectors of  $\mathbf{A}^\top \mathbf{A}$ . The non-zero singular values of  $\mathbf{A}$  are the square roots of the eigenvalues of  $\mathbf{A}^\top \mathbf{A}$ . The same is true for  $\mathbf{A} \mathbf{A}^\top$ .

Perhaps the most useful feature of the SVD is that we can use it to partially generalize matrix inversion to non-square matrices, as we will see in the next section.

## 2.9 The Moore-Penrose Pseudoinverse

Matrix inversion is not defined for matrices that are not square. Suppose we want to make a left-inverse  $\mathbf{B}$  of a matrix  $\mathbf{A}$ , so that we can solve a linear equation

$$\mathbf{A} \mathbf{x} = \mathbf{y} \quad (2.44)$$

by left-multiplying each side to obtain

$$\mathbf{x} = \mathbf{B}\mathbf{y}. \quad (2.45)$$

Depending on the structure of the problem, it may not be possible to design a unique mapping from  $\mathbf{A}$  to  $\mathbf{B}$ .

If  $\mathbf{A}$  is taller than it is wide, then it is possible for this equation to have no solution. If  $\mathbf{A}$  is wider than it is tall, then there could be multiple possible solutions.

The **Moore-Penrose pseudoinverse** allows us to make some headway in these cases. The pseudoinverse of  $\mathbf{A}$  is defined as a matrix

$$\mathbf{A}^+ = \lim_{\alpha \searrow 0} (\mathbf{A}^\top \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^\top. \quad (2.46)$$

Practical algorithms for computing the pseudoinverse are not based on this definition, but rather the formula

$$\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{U}^\top, \quad (2.47)$$

where  $\mathbf{U}$ ,  $\mathbf{D}$  and  $\mathbf{V}$  are the singular value decomposition of  $\mathbf{A}$ , and the pseudoinverse  $\mathbf{D}^+$  of a diagonal matrix  $\mathbf{D}$  is obtained by taking the reciprocal of its non-zero elements then taking the transpose of the resulting matrix.

When  $\mathbf{A}$  has more columns than rows, then solving a linear equation using the pseudoinverse provides one of the many possible solutions. Specifically, it provides the solution  $\mathbf{x} = \mathbf{A}^+ \mathbf{y}$  with minimal Euclidean norm  $\|\mathbf{x}\|_2$  among all possible solutions.

When  $\mathbf{A}$  has more rows than columns, it is possible for there to be no solution. In this case, using the pseudoinverse gives us the  $\mathbf{x}$  for which  $\mathbf{A}\mathbf{x}$  is as close as possible to  $\mathbf{y}$  in terms of Euclidean norm  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$ .

## 2.10 The Trace Operator

The trace operator gives the sum of all of the diagonal entries of a matrix:

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}. \quad (2.48)$$

The trace operator is useful for a variety of reasons. Some operations that are difficult to specify without resorting to summation notation can be specified using

matrix products and the trace operator. For example, the trace operator provides an alternative way of writing the Frobenius norm of a matrix:

$$\|A\|_F = \sqrt{\text{Tr}(\mathbf{A}\mathbf{A}^\top)}. \quad (2.49)$$

Writing an expression in terms of the trace operator opens up opportunities to manipulate the expression using many useful identities. For example, the trace operator is invariant to the transpose operator:

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^\top). \quad (2.50)$$

The trace of a square matrix composed of many factors is also invariant to moving the last factor into the first position, if the shapes of the corresponding matrices allow the resulting product to be defined:

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA}) \quad (2.51)$$

or more generally,

$$\text{Tr}\left(\prod_{i=1}^n \mathbf{F}^{(i)}\right) = \text{Tr}\left(\mathbf{F}^{(n)} \prod_{i=1}^{n-1} \mathbf{F}^{(i)}\right). \quad (2.52)$$

This invariance to cyclic permutation holds even if the resulting product has a different shape. For example, for  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times m}$ , we have

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}) \quad (2.53)$$

even though  $\mathbf{AB} \in \mathbb{R}^{m \times m}$  and  $\mathbf{BA} \in \mathbb{R}^{n \times n}$ .

Another useful fact to keep in mind is that a scalar is its own trace:  $a = \text{Tr}(a)$ .

## 2.11 The Determinant

The determinant of a square matrix, denoted  $\det(\mathbf{A})$ , is a function mapping matrices to real scalars. The determinant is equal to the product of all the eigenvalues of the matrix. The absolute value of the determinant can be thought of as a measure of how much multiplication by the matrix expands or contracts space. If the determinant is 0, then space is contracted completely along at least one dimension, causing it to lose all of its volume. If the determinant is 1, then the transformation preserves volume.

## 2.12 Example: Principal Components Analysis

One simple machine learning algorithm, **principal components analysis** or PCA can be derived using only knowledge of basic linear algebra.

Suppose we have a collection of  $m$  points  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  in  $\mathbb{R}^n$ . Suppose we would like to apply lossy compression to these points. Lossy compression means storing the points in a way that requires less memory but may lose some precision. We would like to lose as little precision as possible.

One way we can encode these points is to represent a lower-dimensional version of them. For each point  $\mathbf{x}^{(i)} \in \mathbb{R}^n$  we will find a corresponding code vector  $\mathbf{c}^{(i)} \in \mathbb{R}^l$ . If  $l$  is smaller than  $n$ , it will take less memory to store the code points than the original data. We will want to find some encoding function that produces the code for an input,  $f(\mathbf{x}) = \mathbf{c}$ , and a decoding function that produces the reconstructed input given its code,  $\mathbf{x} \approx g(f(\mathbf{x}))$ .

PCA is defined by our choice of the decoding function. Specifically, to make the decoder very simple, we choose to use matrix multiplication to map the code back into  $\mathbb{R}^n$ . Let  $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ , where  $\mathbf{D} \in \mathbb{R}^{n \times l}$  is the matrix defining the decoding.

Computing the optimal code for this decoder could be a difficult problem. To keep the encoding problem easy, PCA constrains the columns of  $\mathbf{D}$  to be orthogonal to each other. (Note that  $\mathbf{D}$  is still not technically “an orthogonal matrix” unless  $l = n$ )

With the problem as described so far, many solutions are possible, because we can increase the scale of  $\mathbf{D}_{:,i}$  if we decrease  $c_i$  proportionally for all points. To give the problem a unique solution, we constrain all of the columns of  $\mathbf{D}$  to have unit norm.

In order to turn this basic idea into an algorithm we can implement, the first thing we need to do is figure out how to generate the optimal code point  $\mathbf{c}^*$  for each input point  $\mathbf{x}$ . One way to do this is to minimize the distance between the input point  $\mathbf{x}$  and its reconstruction,  $g(\mathbf{c}^*)$ . We can measure this distance using a norm. In the principal components algorithm, we use the  $L^2$  norm:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2. \quad (2.54)$$

We can switch to the squared  $L^2$  norm instead of the  $L^2$  norm itself, because both are minimized by the same value of  $\mathbf{c}$ . Both are minimized by the same value of  $\mathbf{c}$  because the  $L^2$  norm is non-negative and the squaring operation is

monotonically increasing for non-negative arguments.

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2. \quad (2.55)$$

The function being minimized simplifies to

$$(\mathbf{x} - g(\mathbf{c}))^\top (\mathbf{x} - g(\mathbf{c})) \quad (2.56)$$

(by the definition of the  $L^2$  norm, equation 2.30)

$$= \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top g(\mathbf{c}) - g(\mathbf{c})^\top \mathbf{x} + g(\mathbf{c})^\top g(\mathbf{c}) \quad (2.57)$$

(by the distributive property)

$$= \mathbf{x}^\top \mathbf{x} - 2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c}) \quad (2.58)$$

(because the scalar  $g(\mathbf{c})^\top \mathbf{x}$  is equal to the transpose of itself).

We can now change the function being minimized again, to omit the first term, since this term does not depend on  $\mathbf{c}$ :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c}). \quad (2.59)$$

To make further progress, we must substitute in the definition of  $g(\mathbf{c})$ :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{D}^\top \mathbf{D}\mathbf{c} \quad (2.60)$$

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{I}_l \mathbf{c} \quad (2.61)$$

(by the orthogonality and unit norm constraints on  $\mathbf{D}$ )

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c} \quad (2.62)$$

We can solve this optimization problem using vector calculus (see section 4.3 if you do not know how to do this):

$$\nabla_{\mathbf{c}}(-2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c}) = \mathbf{0} \quad (2.63)$$

$$-2\mathbf{D}^\top \mathbf{x} + 2\mathbf{c} = \mathbf{0} \quad (2.64)$$

$$\mathbf{c} = \mathbf{D}^\top \mathbf{x}. \quad (2.65)$$



This makes the algorithm efficient: we can optimally encode  $\mathbf{x}$  just using a matrix-vector operation. To encode a vector, we apply the encoder function

$$f(\mathbf{x}) = \mathbf{D}^\top \mathbf{x}. \quad (2.66)$$

Using a further matrix multiplication, we can also define the PCA reconstruction operation:

$$r(\mathbf{x}) = g(f(\mathbf{x})) = \mathbf{D}\mathbf{D}^\top \mathbf{x}. \quad (2.67)$$

Next, we need to choose the encoding matrix  $\mathbf{D}$ . To do so, we revisit the idea of minimizing the  $L^2$  distance between inputs and reconstructions. Since we will use the same matrix  $\mathbf{D}$  to decode all of the points, we can no longer consider the points in isolation. Instead, we must minimize the Frobenius norm of the matrix of errors computed over all dimensions and all points:

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left( x_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \text{ subject to } \mathbf{D}^\top \mathbf{D} = \mathbf{I}_l \quad (2.68)$$

To derive the algorithm for finding  $\mathbf{D}^*$ , we will start by considering the case where  $l = 1$ . In this case,  $\mathbf{D}$  is just a single vector,  $\mathbf{d}$ . Substituting equation 2.67 into equation 2.68 and simplifying  $\mathbf{D}$  into  $\mathbf{d}$ , the problem reduces to

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^\top \mathbf{x}^{(i)}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1. \quad (2.69)$$

The above formulation is the most direct way of performing the substitution, but is not the most stylistically pleasing way to write the equation. It places the scalar value  $\mathbf{d}^\top \mathbf{x}^{(i)}$  on the right of the vector  $\mathbf{d}$ . It is more conventional to write scalar coefficients on the left of vector they operate on. We therefore usually write such a formula as

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}^\top \mathbf{x}^{(i)} \mathbf{d}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1, \quad (2.70)$$

or, exploiting the fact that a scalar is its own transpose, as

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d} \mathbf{d}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1. \quad (2.71)$$

The reader should aim to become familiar with such cosmetic rearrangements.

At this point, it can be helpful to rewrite the problem in terms of a single design matrix of examples, rather than as a sum over separate example vectors. This will allow us to use more compact notation. Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  be the matrix defined by stacking all of the vectors describing the points, such that  $\mathbf{X}_{i,:} = \mathbf{x}^{(i)\top}$ . We can now rewrite the problem as

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1. \quad (2.72)$$

Disregarding the constraint for the moment, we can simplify the Frobenius norm portion as follows:

$$\arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 \quad (2.73)$$

$$= \arg \min_{\mathbf{d}} \text{Tr} \left( \left( \mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right)^\top \left( \mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right) \right) \quad (2.74)$$

(by equation 2.49)

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top - \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} + \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.75)$$

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X}) - \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) - \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.76)$$

$$= \arg \min_{\mathbf{d}} -\text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) - \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.77)$$

(because terms not involving  $\mathbf{d}$  do not affect the  $\arg \min$ )

$$= \arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.78)$$

(because we can cycle the order of the matrices inside a trace, equation 2.52)

$$= \arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top \mathbf{d}\mathbf{d}^\top) \quad (2.79)$$

(using the same property again)

At this point, we re-introduce the constraint:

$$\arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top \mathbf{d}\mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.80)$$

$$= \arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.81)$$

(due to the constraint)

$$= \arg \min_{\mathbf{d}} -\text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.82)$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.83)$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.84)$$

This optimization problem may be solved using eigendecomposition. Specifically, the optimal  $\mathbf{d}$  is given by the eigenvector of  $\mathbf{X}^\top \mathbf{X}$  corresponding to the largest eigenvalue.

This derivation is specific to the case of  $l = 1$  and recovers only the first principal component. More generally, when we wish to recover a basis of principal components, the matrix  $\mathbf{D}$  is given by the  $l$  eigenvectors corresponding to the largest eigenvalues. This may be shown using proof by induction. We recommend writing this proof as an exercise.

Linear algebra is one of the fundamental mathematical disciplines that is necessary to understand deep learning. Another key area of mathematics that is ubiquitous in machine learning is probability theory, presented next.

## Chapter 3

# Probability and Information Theory

In this chapter, we describe probability theory and information theory.

Probability theory is a mathematical framework for representing uncertain statements. It provides a means of quantifying uncertainty and axioms for deriving new uncertain statements. In artificial intelligence applications, we use probability theory in two major ways. First, the laws of probability tell us how AI systems should reason, so we design our algorithms to compute or approximate various expressions derived using probability theory. Second, we can use probability and statistics to theoretically analyze the behavior of proposed AI systems.

Probability theory is a fundamental tool of many disciplines of science and engineering. We provide this chapter to ensure that readers whose background is primarily in software engineering with limited exposure to probability theory can understand the material in this book.

While probability theory allows us to make uncertain statements and reason in the presence of uncertainty, information theory allows us to quantify the amount of uncertainty in a probability distribution.

If you are already familiar with probability theory and information theory, you may wish to skip all of this chapter except for section 3.14, which describes the graphs we use to describe structured probabilistic models for machine learning. If you have absolutely no prior experience with these subjects, this chapter should be sufficient to successfully carry out deep learning research projects, but we do suggest that you consult an additional resource, such as [Jaynes \(2003\)](#).

## 3.1 Why Probability?

Many branches of computer science deal mostly with entities that are entirely deterministic and certain. A programmer can usually safely assume that a CPU will execute each machine instruction flawlessly. Errors in hardware do occur, but are rare enough that most software applications do not need to be designed to account for them. Given that many computer scientists and software engineers work in a relatively clean and certain environment, it can be surprising that machine learning makes heavy use of probability theory.

This is because machine learning must always deal with uncertain quantities, and sometimes may also need to deal with stochastic (non-deterministic) quantities. Uncertainty and stochasticity can arise from many sources. Researchers have made compelling arguments for quantifying uncertainty using probability since at least the 1980s. Many of the arguments presented here are summarized from or inspired by [Pearl \(1988\)](#).

Nearly all activities require some ability to reason in the presence of uncertainty. In fact, beyond mathematical statements that are true by definition, it is difficult to think of any proposition that is absolutely true or any event that is absolutely guaranteed to occur.

There are three possible sources of uncertainty:

1. Inherent stochasticity in the system being modeled. For example, most interpretations of quantum mechanics describe the dynamics of subatomic particles as being probabilistic. We can also create theoretical scenarios that we postulate to have random dynamics, such as a hypothetical card game where we assume that the cards are truly shuffled into a random order.
2. Incomplete observability. Even deterministic systems can appear stochastic when we cannot observe all of the variables that drive the behavior of the system. For example, in the Monty Hall problem, a game show contestant is asked to choose between three doors and wins a prize held behind the chosen door. Two doors lead to a goat while a third leads to a car. The outcome given the contestant's choice is deterministic, but from the contestant's point of view, the outcome is uncertain.
3. Incomplete modeling. When we use a model that must discard some of the information we have observed, the discarded information results in uncertainty in the model's predictions. For example, suppose we build a robot that can exactly observe the location of every object around it. If the

robot discretizes space when predicting the future location of these objects, then the discretization makes the robot immediately become uncertain about the precise position of objects: each object could be anywhere within the discrete cell that it was observed to occupy.

In many cases, it is more practical to use a simple but uncertain rule rather than a complex but certain one, even if the true rule is deterministic and our modeling system has the fidelity to accommodate a complex rule. For example, the simple rule “Most birds fly” is cheap to develop and is broadly useful, while a rule of the form, “Birds fly, except for very young birds that have not yet learned to fly, sick or injured birds that have lost the ability to fly, flightless species of birds including the cassowary, ostrich and kiwi. . .” is expensive to develop, maintain and communicate, and after all of this effort is still very brittle and prone to failure.

While it should be clear that we need a means of representing and reasoning about uncertainty, it is not immediately obvious that probability theory can provide all of the tools we want for artificial intelligence applications. Probability theory was originally developed to analyze the frequencies of events. It is easy to see how probability theory can be used to study events like drawing a certain hand of cards in a game of poker. These kinds of events are often repeatable. When we say that an outcome has a probability  $p$  of occurring, it means that if we repeated the experiment (e.g., draw a hand of cards) infinitely many times, then proportion  $p$  of the repetitions would result in that outcome. This kind of reasoning does not seem immediately applicable to propositions that are not repeatable. If a doctor analyzes a patient and says that the patient has a 40% chance of having the flu, this means something very different—we can not make infinitely many replicas of the patient, nor is there any reason to believe that different replicas of the patient would present with the same symptoms yet have varying underlying conditions. In the case of the doctor diagnosing the patient, we use probability to represent a **degree of belief**, with 1 indicating absolute certainty that the patient has the flu and 0 indicating absolute certainty that the patient does not have the flu. The former kind of probability, related directly to the rates at which events occur, is known as **frequentist probability**, while the latter, related to qualitative levels of certainty, is known as **Bayesian probability**.

If we list several properties that we expect common sense reasoning about uncertainty to have, then the only way to satisfy those properties is to treat Bayesian probabilities as behaving exactly the same as frequentist probabilities. For example, if we want to compute the probability that a player will win a poker game given that she has a certain set of cards, we use exactly the same formulas as when we compute the probability that a patient has a disease given that she

has certain symptoms. For more details about why a small set of common sense assumptions implies that the same axioms must control both kinds of probability, see [Ramsey \(1926\)](#).

Probability can be seen as the extension of logic to deal with uncertainty. Logic provides a set of formal rules for determining what propositions are implied to be true or false given the assumption that some other set of propositions is true or false. Probability theory provides a set of formal rules for determining the likelihood of a proposition being true given the likelihood of other propositions.

## 3.2 Random Variables

A **random variable** is a variable that can take on different values randomly. We typically denote the random variable itself with a lower case letter in plain typeface, and the values it can take on with lower case script letters. For example,  $x_1$  and  $x_2$  are both possible values that the random variable  $x$  can take on. For vector-valued variables, we would write the random variable as  $\mathbf{x}$  and one of its values as  $\mathbf{x}$ . On its own, a random variable is just a description of the states that are possible; it must be coupled with a probability distribution that specifies how likely each of these states are.

Random variables may be discrete or continuous. A discrete random variable is one that has a finite or countably infinite number of states. Note that these states are not necessarily the integers; they can also just be named states that are not considered to have any numerical value. A continuous random variable is associated with a real value.

## 3.3 Probability Distributions

A **probability distribution** is a description of how likely a random variable or set of random variables is to take on each of its possible states. The way we describe probability distributions depends on whether the variables are discrete or continuous.

### 3.3.1 Discrete Variables and Probability Mass Functions

A probability distribution over discrete variables may be described using a **probability mass function** (PMF). We typically denote probability mass functions with a capital  $P$ . Often we associate each random variable with a different probability

mass function and the reader must infer which probability mass function to use based on the identity of the random variable, rather than the name of the function;  $P(x)$  is usually not the same as  $P(y)$ .

The probability mass function maps from a state of a random variable to the probability of that random variable taking on that state. The probability that  $x = x$  is denoted as  $P(x)$ , with a probability of 1 indicating that  $x = x$  is certain and a probability of 0 indicating that  $x = x$  is impossible. Sometimes to disambiguate which PMF to use, we write the name of the random variable explicitly:  $P(x = x)$ . Sometimes we define a variable first, then use  $\sim$  notation to specify which distribution it follows later:  $x \sim P(x)$ .

Probability mass functions can act on many variables at the same time. Such a probability distribution over many variables is known as a **joint probability distribution**.  $P(x = x, y = y)$  denotes the probability that  $x = x$  and  $y = y$  simultaneously. We may also write  $P(x, y)$  for brevity.

To be a probability mass function on a random variable  $x$ , a function  $P$  must satisfy the following properties:

- The domain of  $P$  must be the set of all possible states of  $x$ .
- $\forall x \in x, 0 \leq P(x) \leq 1$ . An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring.
- $\sum_{x \in x} P(x) = 1$ . We refer to this property as being **normalized**. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring.

For example, consider a single discrete random variable  $x$  with  $k$  different states. We can place a **uniform distribution** on  $x$ —that is, make each of its states equally likely—by setting its probability mass function to

$$P(x = x_i) = \frac{1}{k} \tag{3.1}$$

for all  $i$ . We can see that this fits the requirements for a probability mass function. The value  $\frac{1}{k}$  is positive because  $k$  is a positive integer. We also see that

$$\sum_i P(x = x_i) = \sum_i \frac{1}{k} = \frac{k}{k} = 1, \tag{3.2}$$

so the distribution is properly normalized.



### 3.3.2 Continuous Variables and Probability Density Functions

When working with continuous random variables, we describe probability distributions using a **probability density function (PDF)** rather than a probability mass function. To be a probability density function, a function  $p$  must satisfy the following properties:

- The domain of  $p$  must be the set of all possible states of  $x$ .
- $\forall x \in \mathbf{x}, p(x) \geq 0$ . Note that we do not require  $p(x) \leq 1$ .
- $\int p(x)dx = 1$ .

A probability density function  $p(x)$  does not give the probability of a specific state directly, instead the probability of landing inside an infinitesimal region with volume  $\delta x$  is given by  $p(x)\delta x$ .

We can integrate the density function to find the actual probability mass of a set of points. Specifically, the probability that  $x$  lies in some set  $\mathbb{S}$  is given by the integral of  $p(x)$  over that set. In the univariate example, the probability that  $x$  lies in the interval  $[a, b]$  is given by  $\int_{[a,b]} p(x)dx$ .

For an example of a probability density function corresponding to a specific probability density over a continuous random variable, consider a uniform distribution on an interval of the real numbers. We can do this with a function  $u(x; a, b)$ , where  $a$  and  $b$  are the endpoints of the interval, with  $b > a$ . The “;” notation means “parametrized by”; we consider  $x$  to be the argument of the function, while  $a$  and  $b$  are parameters that define the function. To ensure that there is no probability mass outside the interval, we say  $u(x; a, b) = 0$  for all  $x \notin [a, b]$ . Within  $[a, b]$ ,  $u(x; a, b) = \frac{1}{b-a}$ . We can see that this is nonnegative everywhere. Additionally, it integrates to 1. We often denote that  $x$  follows the uniform distribution on  $[a, b]$  by writing  $x \sim U(a, b)$ .

## 3.4 Marginal Probability

Sometimes we know the probability distribution over a set of variables and we want to know the probability distribution over just a subset of them. The probability distribution over the subset is known as the **marginal probability** distribution.

For example, suppose we have discrete random variables  $x$  and  $y$ , and we know  $P(x, y)$ . We can find  $P(x)$  with the **sum rule**:

$$\forall x \in \mathbf{x}, P(x = x) = \sum_y P(x = x, y = y). \quad (3.3)$$

The name “marginal probability” comes from the process of computing marginal probabilities on paper. When the values of  $P(x, y)$  are written in a grid with different values of  $x$  in rows and different values of  $y$  in columns, it is natural to sum across a row of the grid, then write  $P(x)$  in the margin of the paper just to the right of the row.

For continuous variables, we need to use integration instead of summation:

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

## 3.5 Conditional Probability

In many cases, we are interested in the probability of some event, given that some other event has happened. This is called a **conditional probability**. We denote the conditional probability that  $y = y$  given  $x = x$  as  $P(y = y \mid x = x)$ . This conditional probability can be computed with the formula

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}. \quad (3.5)$$

The conditional probability is only defined when  $P(x = x) > 0$ . We cannot compute the conditional probability conditioned on an event that never happens.

It is important not to confuse conditional probability with computing what would happen if some action were undertaken. The conditional probability that a person is from Germany given that they speak German is quite high, but if a randomly selected person is taught to speak German, their country of origin does not change. Computing the consequences of an action is called making an **intervention query**. Intervention queries are the domain of **causal modeling**, which we do not explore in this book.

## 3.6 The Chain Rule of Conditional Probabilities

Any joint probability distribution over many random variables may be decomposed into conditional distributions over only one variable:

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)}). \quad (3.6)$$

This observation is known as the **chain rule** or **product rule** of probability. It follows immediately from the definition of conditional probability in equation 3.5.

For example, applying the definition twice, we get

$$\begin{aligned} P(a, b, c) &= P(a \mid b, c)P(b, c) \\ P(b, c) &= P(b \mid c)P(c) \\ P(a, b, c) &= P(a \mid b, c)P(b \mid c)P(c). \end{aligned}$$

### 3.7 Independence and Conditional Independence

Two random variables  $x$  and  $y$  are **independent** if their probability distribution can be expressed as a product of two factors, one involving only  $x$  and one involving only  $y$ :

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, p(x = x, y = y) = p(x = x)p(y = y). \quad (3.7)$$

Two random variables  $x$  and  $y$  are **conditionally independent** given a random variable  $z$  if the conditional probability distribution over  $x$  and  $y$  factorizes in this way for every value of  $z$ :

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, p(x = x, y = y \mid z = z) = p(x = x \mid z = z)p(y = y \mid z = z). \quad (3.8)$$

We can denote independence and conditional independence with compact notation:  $x \perp y$  means that  $x$  and  $y$  are independent, while  $x \perp y \mid z$  means that  $x$  and  $y$  are conditionally independent given  $z$ .

### 3.8 Expectation, Variance and Covariance

The **expectation** or **expected value** of some function  $f(x)$  with respect to a probability distribution  $P(x)$  is the average or mean value that  $f$  takes on when  $x$  is drawn from  $P$ . For discrete variables this can be computed with a summation:

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x), \quad (3.9)$$

while for continuous variables, it is computed with an integral:

$$\mathbb{E}_{x \sim p}[f(x)] = \int p(x)f(x)dx. \quad (3.10)$$

When the identity of the distribution is clear from the context, we may simply write the name of the random variable that the expectation is over, as in  $\mathbb{E}_x[f(x)]$ . If it is clear which random variable the expectation is over, we may omit the subscript entirely, as in  $\mathbb{E}[f(x)]$ . By default, we can assume that  $\mathbb{E}[\cdot]$  averages over the values of all the random variables inside the brackets. Likewise, when there is no ambiguity, we may omit the square brackets.

Expectations are linear, for example,

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)], \quad (3.11)$$

when  $\alpha$  and  $\beta$  are not dependent on  $x$ .

The **variance** gives a measure of how much the values of a function of a random variable  $x$  vary as we sample different values of  $x$  from its probability distribution:

$$\text{Var}(f(x)) = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right]. \quad (3.12)$$

When the variance is low, the values of  $f(x)$  cluster near their expected value. The square root of the variance is known as the **standard deviation**.

The **covariance** gives some sense of how much two values are linearly related to each other, as well as the scale of these variables:

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)]) (g(y) - \mathbb{E}[g(y)])]. \quad (3.13)$$

High absolute values of the covariance mean that the values change very much and are both far from their respective means at the same time. If the sign of the covariance is positive, then both variables tend to take on relatively high values simultaneously. If the sign of the covariance is negative, then one variable tends to take on a relatively high value at the times that the other takes on a relatively low value and vice versa. Other measures such as **correlation** normalize the contribution of each variable in order to measure only how much the variables are related, rather than also being affected by the scale of the separate variables.

The notions of covariance and dependence are related, but are in fact distinct concepts. They are related because two variables that are independent have zero covariance, and two variables that have non-zero covariance are dependent. However, independence is a distinct property from covariance. For two variables to have zero covariance, there must be no linear dependence between them. Independence is a stronger requirement than zero covariance, because independence also excludes nonlinear relationships. It is possible for two variables to be dependent but have zero covariance. For example, suppose we first sample a real number  $x$  from a uniform distribution over the interval  $[-1, 1]$ . We next sample a random variable

$s$ . With probability  $\frac{1}{2}$ , we choose the value of  $s$  to be 1. Otherwise, we choose the value of  $s$  to be  $-1$ . We can then generate a random variable  $y$  by assigning  $y = sx$ . Clearly,  $x$  and  $y$  are not independent, because  $x$  completely determines the magnitude of  $y$ . However,  $\text{Cov}(x, y) = 0$ .

The **covariance matrix** of a random vector  $\mathbf{x} \in \mathbb{R}^n$  is an  $n \times n$  matrix, such that

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j). \quad (3.14)$$

The diagonal elements of the covariance give the variance:

$$\text{Cov}(x_i, x_i) = \text{Var}(x_i). \quad (3.15)$$

## 3.9 Common Probability Distributions

Several simple probability distributions are useful in many contexts in machine learning.

### 3.9.1 Bernoulli Distribution

The **Bernoulli** distribution is a distribution over a single binary random variable. It is controlled by a single parameter  $\phi \in [0, 1]$ , which gives the probability of the random variable being equal to 1. It has the following properties:

$$P(x = 1) = \phi \quad (3.16)$$

$$P(x = 0) = 1 - \phi \quad (3.17)$$

$$P(x = x) = \phi^x (1 - \phi)^{1-x} \quad (3.18)$$

$$\mathbb{E}_x[x] = \phi \quad (3.19)$$

$$\text{Var}_x(x) = \phi(1 - \phi) \quad (3.20)$$

### 3.9.2 Multinoulli Distribution

The **multinoulli** or **categorical** distribution is a distribution over a single discrete variable with  $k$  different states, where  $k$  is finite.<sup>1</sup> The multinoulli distribution is

---

<sup>1</sup> “Multinoulli” is a term that was recently coined by Gustavo Lacerdo and popularized by [Murphy \(2012\)](#). The multinoulli distribution is a special case of the **multinomial** distribution. A multinomial distribution is the distribution over vectors in  $\{0, \dots, n\}^k$  representing how many times each of the  $k$  categories is visited when  $n$  samples are drawn from a multinoulli distribution. Many texts use the term “multinomial” to refer to multinoulli distributions without clarifying that they refer only to the  $n = 1$  case.

parametrized by a vector  $\mathbf{p} \in [0, 1]^{k-1}$ , where  $p_i$  gives the probability of the  $i$ -th state. The final,  $k$ -th state's probability is given by  $1 - \mathbf{1}^\top \mathbf{p}$ . Note that we must constrain  $\mathbf{1}^\top \mathbf{p} \leq 1$ . Multinoulli distributions are often used to refer to distributions over categories of objects, so we do not usually assume that state 1 has numerical value 1, etc. For this reason, we do not usually need to compute the expectation or variance of multinoulli-distributed random variables.

The Bernoulli and multinoulli distributions are sufficient to describe any distribution over their domain. They are able to describe any distribution over their domain not so much because they are particularly powerful but rather because their domain is simple; they model discrete variables for which it is feasible to enumerate all of the states. When dealing with continuous variables, there are uncountably many states, so any distribution described by a small number of parameters must impose strict limits on the distribution.

### 3.9.3 Gaussian Distribution

The most commonly used distribution over real numbers is the **normal distribution**, also known as the **Gaussian distribution**:

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

See figure 3.1 for a plot of the density function.

The two parameters  $\mu \in \mathbb{R}$  and  $\sigma \in (0, \infty)$  control the normal distribution. The parameter  $\mu$  gives the coordinate of the central peak. This is also the mean of the distribution:  $\mathbb{E}[x] = \mu$ . The standard deviation of the distribution is given by  $\sigma$ , and the variance by  $\sigma^2$ .

When we evaluate the PDF, we need to square and invert  $\sigma$ . When we need to frequently evaluate the PDF with different parameter values, a more efficient way of parametrizing the distribution is to use a parameter  $\beta \in (0, \infty)$  to control the **precision** or inverse variance of the distribution:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

Normal distributions are a sensible choice for many applications. In the absence of prior knowledge about what form a distribution over the real numbers should take, the normal distribution is a good default choice for two major reasons.

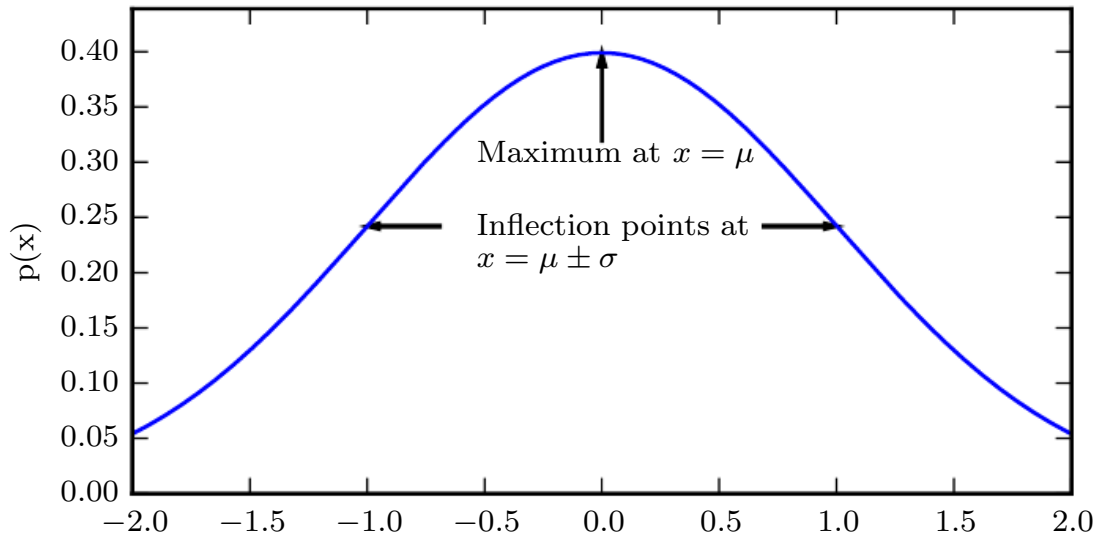


Figure 3.1: **The normal distribution:** The normal distribution  $\mathcal{N}(x; \mu, \sigma^2)$  exhibits a classic “bell curve” shape, with the  $x$  coordinate of its central peak given by  $\mu$ , and the width of its peak controlled by  $\sigma$ . In this example, we depict the **standard normal distribution**, with  $\mu = 0$  and  $\sigma = 1$ .

First, many distributions we wish to model are truly close to being normal distributions. The **central limit theorem** shows that the sum of many independent random variables is approximately normally distributed. This means that in practice, many complicated systems can be modeled successfully as normally distributed noise, even if the system can be decomposed into parts with more structured behavior.

Second, out of all possible probability distributions with the same variance, the normal distribution encodes the maximum amount of uncertainty over the real numbers. We can thus think of the normal distribution as being the one that inserts the least amount of prior knowledge into a model. Fully developing and justifying this idea requires more mathematical tools, and is postponed to section 19.4.2.

The normal distribution generalizes to  $\mathbb{R}^n$ , in which case it is known as the **multivariate normal distribution**. It may be parametrized with a positive definite symmetric matrix  $\Sigma$ :

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right). \quad (3.23)$$

The parameter  $\boldsymbol{\mu}$  still gives the mean of the distribution, though now it is vector-valued. The parameter  $\boldsymbol{\Sigma}$  gives the covariance matrix of the distribution. As in the univariate case, when we wish to evaluate the PDF several times for many different values of the parameters, the covariance is not a computationally efficient way to parametrize the distribution, since we need to invert  $\boldsymbol{\Sigma}$  to evaluate the PDF. We can instead use a **precision matrix**  $\boldsymbol{\beta}$ :

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\beta}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.24)$$

We often fix the covariance matrix to be a diagonal matrix. An even simpler version is the **isotropic** Gaussian distribution, whose covariance matrix is a scalar times the identity matrix.

### 3.9.4 Exponential and Laplace Distributions

In the context of deep learning, we often want to have a probability distribution with a sharp point at  $x = 0$ . To accomplish this, we can use the **exponential distribution**:

$$p(x; \lambda) = \lambda \mathbf{1}_{x \geq 0} \exp(-\lambda x). \quad (3.25)$$

The exponential distribution uses the indicator function  $\mathbf{1}_{x \geq 0}$  to assign probability zero to all negative values of  $x$ .

A closely related probability distribution that allows us to place a sharp peak of probability mass at an arbitrary point  $\mu$  is the **Laplace distribution**

$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right). \quad (3.26)$$

### 3.9.5 The Dirac Distribution and Empirical Distribution

In some cases, we wish to specify that all of the mass in a probability distribution clusters around a single point. This can be accomplished by defining a PDF using the Dirac delta function,  $\delta(x)$ :

$$p(x) = \delta(x - \mu). \quad (3.27)$$

The Dirac delta function is defined such that it is zero-valued everywhere except 0, yet integrates to 1. The Dirac delta function is not an ordinary function that associates each value  $x$  with a real-valued output, instead it is a different kind of



mathematical object called a **generalized function** that is defined in terms of its properties when integrated. We can think of the Dirac delta function as being the limit point of a series of functions that put less and less mass on all points other than zero.

By defining  $p(x)$  to be  $\delta$  shifted by  $-\mu$  we obtain an infinitely narrow and infinitely high peak of probability mass where  $x = \mu$ .

A common use of the Dirac delta distribution is as a component of an **empirical distribution**,

$$\hat{p}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)}) \quad (3.28)$$

which puts probability mass  $\frac{1}{m}$  on each of the  $m$  points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$  forming a given dataset or collection of samples. The Dirac delta distribution is only necessary to define the empirical distribution over continuous variables. For discrete variables, the situation is simpler: an empirical distribution can be conceptualized as a multinoulli distribution, with a probability associated to each possible input value that is simply equal to the **empirical frequency** of that value in the training set.

We can view the empirical distribution formed from a dataset of training examples as specifying the distribution that we sample from when we train a model on this dataset. Another important perspective on the empirical distribution is that it is the probability density that maximizes the likelihood of the training data (see section 5.5).

### 3.9.6 Mixtures of Distributions

It is also common to define probability distributions by combining other simpler probability distributions. One common way of combining distributions is to construct a **mixture distribution**. A mixture distribution is made up of several component distributions. On each trial, the choice of which component distribution generates the sample is determined by sampling a component identity from a multinoulli distribution:

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} \mid c = i) \quad (3.29)$$

where  $P(c)$  is the multinoulli distribution over component identities.

We have already seen one example of a mixture distribution: the empirical distribution over real-valued variables is a mixture distribution with one Dirac component for each training example.

The mixture model is one simple strategy for combining probability distributions to create a richer distribution. In chapter 16, we explore the art of building complex probability distributions from simple ones in more detail.

The mixture model allows us to briefly glimpse a concept that will be of paramount importance later—the **latent variable**. A latent variable is a random variable that we cannot observe directly. The component identity variable  $c$  of the mixture model provides an example. Latent variables may be related to  $\mathbf{x}$  through the joint distribution, in this case,  $P(\mathbf{x}, c) = P(\mathbf{x} | c)P(c)$ . The distribution  $P(c)$  over the latent variable and the distribution  $P(\mathbf{x} | c)$  relating the latent variables to the visible variables determines the shape of the distribution  $P(\mathbf{x})$  even though it is possible to describe  $P(\mathbf{x})$  without reference to the latent variable. Latent variables are discussed further in section 16.5.

A very powerful and common type of mixture model is the **Gaussian mixture** model, in which the components  $p(\mathbf{x} | c = i)$  are Gaussians. Each component has a separately parametrized mean  $\boldsymbol{\mu}^{(i)}$  and covariance  $\boldsymbol{\Sigma}^{(i)}$ . Some mixtures can have more constraints. For example, the covariances could be shared across components via the constraint  $\boldsymbol{\Sigma}^{(i)} = \boldsymbol{\Sigma}, \forall i$ . As with a single Gaussian distribution, the mixture of Gaussians might constrain the covariance matrix for each component to be diagonal or isotropic.

In addition to the means and covariances, the parameters of a Gaussian mixture specify the **prior probability**  $\alpha_i = P(c = i)$  given to each component  $i$ . The word “prior” indicates that it expresses the model’s beliefs about  $c$  *before* it has observed  $\mathbf{x}$ . By comparison,  $P(c | \mathbf{x})$  is a **posterior probability**, because it is computed *after* observation of  $\mathbf{x}$ . A Gaussian mixture model is a **universal approximator** of densities, in the sense that any smooth density can be approximated with any specific, non-zero amount of error by a Gaussian mixture model with enough components.

Figure 3.2 shows samples from a Gaussian mixture model.

## 3.10 Useful Properties of Common Functions

Certain functions arise often while working with probability distributions, especially the probability distributions used in deep learning models.

One of these functions is the **logistic sigmoid**:

$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \quad (3.30)$$

The logistic sigmoid is commonly used to produce the  $\phi$  parameter of a Bernoulli

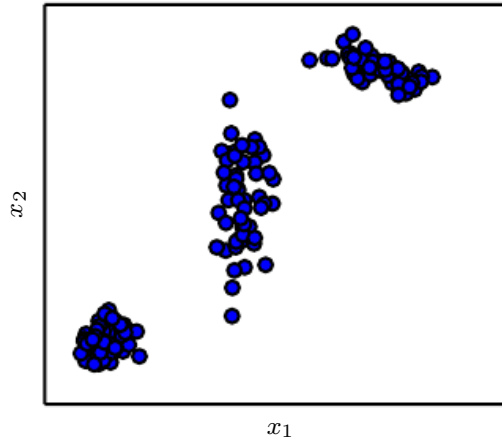


Figure 3.2: Samples from a Gaussian mixture model. In this example, there are three components. From left to right, the first component has an isotropic covariance matrix, meaning it has the same amount of variance in each direction. The second has a diagonal covariance matrix, meaning it can control the variance separately along each axis-aligned direction. This example has more variance along the  $x_2$  axis than along the  $x_1$  axis. The third component has a full-rank covariance matrix, allowing it to control the variance separately along an arbitrary basis of directions.

distribution because its range is  $(0,1)$ , which lies within the valid range of values for the  $\phi$  parameter. See figure 3.3 for a graph of the sigmoid function. The sigmoid function **saturates** when its argument is very positive or very negative, meaning that the function becomes very flat and insensitive to small changes in its input.

Another commonly encountered function is the **softplus** function (Dugas *et al.*, 2001):

$$\zeta(x) = \log(1 + \exp(x)). \quad (3.31)$$

The softplus function can be useful for producing the  $\beta$  or  $\sigma$  parameter of a normal distribution because its range is  $(0, \infty)$ . It also arises commonly when manipulating expressions involving sigmoids. The name of the softplus function comes from the fact that it is a smoothed or “softened” version of

$$x^+ = \max(0, x). \quad (3.32)$$

See figure 3.4 for a graph of the softplus function.

The following properties are all useful enough that you may wish to memorize them:

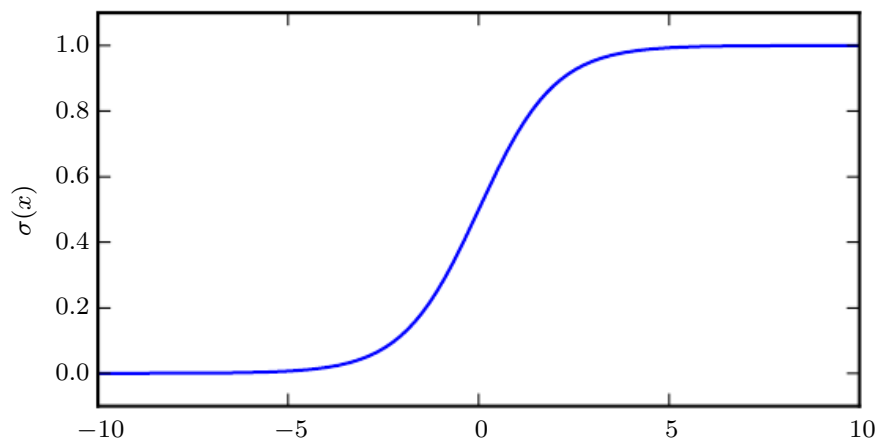


Figure 3.3: The logistic sigmoid function.

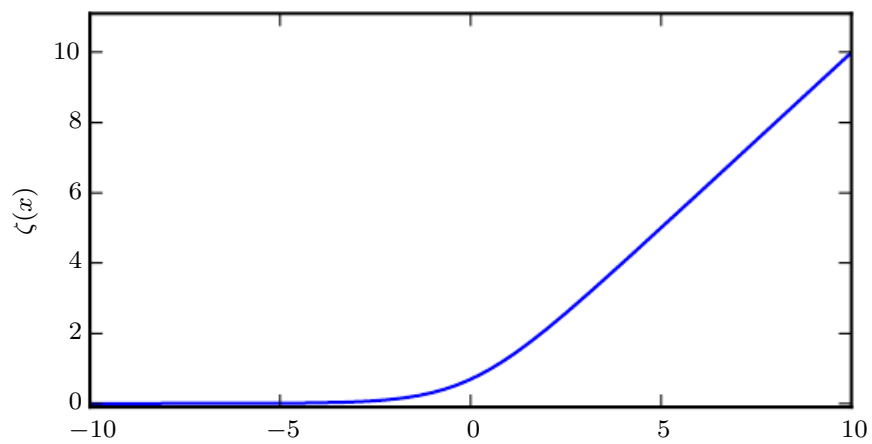


Figure 3.4: The softplus function.

$$\sigma(x) = \frac{\exp(x)}{\exp(x) + \exp(0)} \quad (3.33)$$

$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x)) \quad (3.34)$$

$$1 - \sigma(x) = \sigma(-x) \quad (3.35)$$

$$\log \sigma(x) = -\zeta(-x) \quad (3.36)$$

$$\frac{d}{dx}\zeta(x) = \sigma(x) \quad (3.37)$$

$$\forall x \in (0, 1), \sigma^{-1}(x) = \log\left(\frac{x}{1-x}\right) \quad (3.38)$$

$$\forall x > 0, \zeta^{-1}(x) = \log(\exp(x) - 1) \quad (3.39)$$

$$\zeta(x) = \int_{-\infty}^x \sigma(y)dy \quad (3.40)$$

$$\zeta(x) - \zeta(-x) = x \quad (3.41)$$

The function  $\sigma^{-1}(x)$  is called the **logit** in statistics, but this term is more rarely used in machine learning.

Equation 3.41 provides extra justification for the name “softplus.” The softplus function is intended as a smoothed version of the **positive part** function,  $x^+ = \max\{0, x\}$ . The positive part function is the counterpart of the **negative part** function,  $x^- = \max\{0, -x\}$ . To obtain a smooth function that is analogous to the negative part, one can use  $\zeta(-x)$ . Just as  $x$  can be recovered from its positive part and negative part via the identity  $x^+ - x^- = x$ , it is also possible to recover  $x$  using the same relationship between  $\zeta(x)$  and  $\zeta(-x)$ , as shown in equation 3.41.

### 3.11 Bayes’ Rule

We often find ourselves in a situation where we know  $P(y | x)$  and need to know  $P(x | y)$ . Fortunately, if we also know  $P(x)$ , we can compute the desired quantity using **Bayes’ rule**:

$$P(x | y) = \frac{P(x)P(y | x)}{P(y)}. \quad (3.42)$$

Note that while  $P(y)$  appears in the formula, it is usually feasible to compute  $P(y) = \sum_x P(y | x)P(x)$ , so we do not need to begin with knowledge of  $P(y)$ .

Bayes' rule is straightforward to derive from the definition of conditional probability, but it is useful to know the name of this formula since many texts refer to it by name. It is named after the Reverend Thomas Bayes, who first discovered a special case of the formula. The general version presented here was independently discovered by Pierre-Simon Laplace.

## 3.12 Technical Details of Continuous Variables

A proper formal understanding of continuous random variables and probability density functions requires developing probability theory in terms of a branch of mathematics known as **measure theory**. Measure theory is beyond the scope of this textbook, but we can briefly sketch some of the issues that measure theory is employed to resolve.

In section 3.3.2, we saw that the probability of a continuous vector-valued  $\mathbf{x}$  lying in some set  $\mathbb{S}$  is given by the integral of  $p(\mathbf{x})$  over the set  $\mathbb{S}$ . Some choices of set  $\mathbb{S}$  can produce paradoxes. For example, it is possible to construct two sets  $\mathbb{S}_1$  and  $\mathbb{S}_2$  such that  $p(\mathbf{x} \in \mathbb{S}_1) + p(\mathbf{x} \in \mathbb{S}_2) > 1$  but  $\mathbb{S}_1 \cap \mathbb{S}_2 = \emptyset$ . These sets are generally constructed making very heavy use of the infinite precision of real numbers, for example by making fractal-shaped sets or sets that are defined by transforming the set of rational numbers.<sup>2</sup> One of the key contributions of measure theory is to provide a characterization of the set of sets that we can compute the probability of without encountering paradoxes. In this book, we only integrate over sets with relatively simple descriptions, so this aspect of measure theory never becomes a relevant concern.

For our purposes, measure theory is more useful for describing theorems that apply to most points in  $\mathbb{R}^n$  but do not apply to some corner cases. Measure theory provides a rigorous way of describing that a set of points is negligibly small. Such a set is said to have **measure zero**. We do not formally define this concept in this textbook. For our purposes, it is sufficient to understand the intuition that a set of measure zero occupies no volume in the space we are measuring. For example, within  $\mathbb{R}^2$ , a line has measure zero, while a filled polygon has positive measure. Likewise, an individual point has measure zero. Any union of countably many sets that each have measure zero also has measure zero (so the set of all the rational numbers has measure zero, for instance).

Another useful term from measure theory is **almost everywhere**. A property that holds almost everywhere holds throughout all of space except for on a set of

---

<sup>2</sup>The Banach-Tarski theorem provides a fun example of such sets.

measure zero. Because the exceptions occupy a negligible amount of space, they can be safely ignored for many applications. Some important results in probability theory hold for all discrete values but only hold “almost everywhere” for continuous values.

Another technical detail of continuous variables relates to handling continuous random variables that are deterministic functions of one another. Suppose we have two random variables,  $\mathbf{x}$  and  $\mathbf{y}$ , such that  $\mathbf{y} = g(\mathbf{x})$ , where  $g$  is an invertible, continuous, differentiable transformation. One might expect that  $p_y(\mathbf{y}) = p_x(g^{-1}(\mathbf{y}))$ . This is actually not the case.

As a simple example, suppose we have scalar random variables  $x$  and  $y$ . Suppose  $y = \frac{x}{2}$  and  $x \sim U(0, 1)$ . If we use the rule  $p_y(y) = p_x(2y)$  then  $p_y$  will be 0 everywhere except the interval  $[0, \frac{1}{2}]$ , and it will be 1 on this interval. This means

$$\int p_y(y)dy = \frac{1}{2}, \quad (3.43)$$

which violates the definition of a probability distribution. This is a common mistake. The problem with this approach is that it fails to account for the distortion of space introduced by the function  $g$ . Recall that the probability of  $\mathbf{x}$  lying in an infinitesimally small region with volume  $\delta\mathbf{x}$  is given by  $p(\mathbf{x})\delta\mathbf{x}$ . Since  $g$  can expand or contract space, the infinitesimal volume surrounding  $\mathbf{x}$  in  $\mathbf{x}$  space may have different volume in  $\mathbf{y}$  space.

To see how to correct the problem, we return to the scalar case. We need to preserve the property

$$|p_y(g(x))dy| = |p_x(x)dx|. \quad (3.44)$$

Solving from this, we obtain

$$p_y(y) = p_x(g^{-1}(y)) \left| \frac{\partial x}{\partial y} \right| \quad (3.45)$$

or equivalently

$$p_x(x) = p_y(g(x)) \left| \frac{\partial g(x)}{\partial x} \right|. \quad (3.46)$$

In higher dimensions, the derivative generalizes to the determinant of the **Jacobian matrix**—the matrix with  $J_{i,j} = \frac{\partial x_i}{\partial y_j}$ . Thus, for real-valued vectors  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|. \quad (3.47)$$

### 3.13 Information Theory

Information theory is a branch of applied mathematics that revolves around quantifying how much information is present in a signal. It was originally invented to study sending messages from discrete alphabets over a noisy channel, such as communication via radio transmission. In this context, information theory tells how to design optimal codes and calculate the expected length of messages sampled from specific probability distributions using various encoding schemes. In the context of machine learning, we can also apply information theory to continuous variables where some of these message length interpretations do not apply. This field is fundamental to many areas of electrical engineering and computer science. In this textbook, we mostly use a few key ideas from information theory to characterize probability distributions or quantify similarity between probability distributions. For more detail on information theory, see [Cover and Thomas \(2006\)](#) or [MacKay \(2003\)](#).

The basic intuition behind information theory is that learning that an unlikely event has occurred is more informative than learning that a likely event has occurred. A message saying “the sun rose this morning” is so uninformative as to be unnecessary to send, but a message saying “there was a solar eclipse this morning” is very informative.

We would like to quantify information in a way that formalizes this intuition. Specifically,

- Likely events should have low information content, and in the extreme case, events that are guaranteed to happen should have no information content whatsoever.
- Less likely events should have higher information content.
- Independent events should have additive information. For example, finding out that a tossed coin has come up as heads twice should convey twice as much information as finding out that a tossed coin has come up as heads once.

In order to satisfy all three of these properties, we define the **self-information** of an event  $x = x$  to be

$$I(x) = -\log P(x). \quad (3.48)$$

In this book, we always use  $\log$  to mean the natural logarithm, with base  $e$ . Our definition of  $I(x)$  is therefore written in units of **nats**. One nat is the amount of



information gained by observing an event of probability  $\frac{1}{e}$ . Other texts use base-2 logarithms and units called **bits** or **shannons**; information measured in bits is just a rescaling of information measured in nats.

When  $x$  is continuous, we use the same definition of information by analogy, but some of the properties from the discrete case are lost. For example, an event with unit density still has zero information, despite not being an event that is guaranteed to occur.

Self-information deals only with a single outcome. We can quantify the amount of uncertainty in an entire probability distribution using the **Shannon entropy**:

$$H(x) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)]. \quad (3.49)$$

also denoted  $H(P)$ . In other words, the Shannon entropy of a distribution is the expected amount of information in an event drawn from that distribution. It gives a lower bound on the number of bits (if the logarithm is base 2, otherwise the units are different) needed on average to encode symbols drawn from a distribution  $P$ . Distributions that are nearly deterministic (where the outcome is nearly certain) have low entropy; distributions that are closer to uniform have high entropy. See figure 3.5 for a demonstration. When  $x$  is continuous, the Shannon entropy is known as the **differential entropy**.

If we have two separate probability distributions  $P(x)$  and  $Q(x)$  over the same random variable  $x$ , we can measure how different these two distributions are using the **Kullback-Leibler (KL) divergence**:

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

In the case of discrete variables, it is the extra amount of information (measured in bits if we use the base 2 logarithm, but in machine learning we usually use nats and the natural logarithm) needed to send a message containing symbols drawn from probability distribution  $P$ , when we use a code that was designed to minimize the length of messages drawn from probability distribution  $Q$ .

The KL divergence has many useful properties, most notably that it is non-negative. The KL divergence is 0 if and only if  $P$  and  $Q$  are the same distribution in the case of discrete variables, or equal “almost everywhere” in the case of continuous variables. Because the KL divergence is non-negative and measures the difference between two distributions, it is often conceptualized as measuring some sort of distance between these distributions. However, it is not a true distance measure because it is not symmetric:  $D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P)$  for some  $P$  and  $Q$ . This

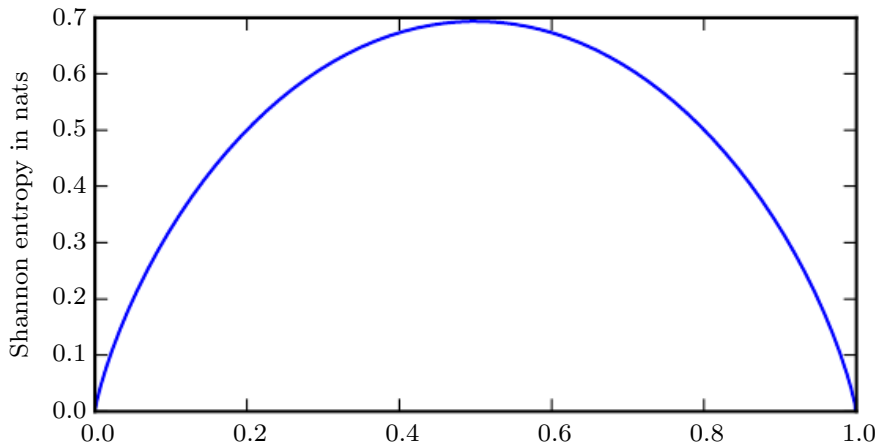


Figure 3.5: This plot shows how distributions that are closer to deterministic have low Shannon entropy while distributions that are close to uniform have high Shannon entropy. On the horizontal axis, we plot  $p$ , the probability of a binary random variable being equal to 1. The entropy is given by  $(p-1)\log(1-p) - p\log p$ . When  $p$  is near 0, the distribution is nearly deterministic, because the random variable is nearly always 0. When  $p$  is near 1, the distribution is nearly deterministic, because the random variable is nearly always 1. When  $p = 0.5$ , the entropy is maximal, because the distribution is uniform over the two outcomes.

asymmetry means that there are important consequences to the choice of whether to use  $D_{\text{KL}}(P\|Q)$  or  $D_{\text{KL}}(Q\|P)$ . See figure 3.6 for more detail.

A quantity that is closely related to the KL divergence is the **cross-entropy**  $H(P, Q) = H(P) + D_{\text{KL}}(P\|Q)$ , which is similar to the KL divergence but lacking the term on the left:

$$H(P, Q) = -\mathbb{E}_{x \sim P} \log Q(x). \quad (3.51)$$

Minimizing the cross-entropy with respect to  $Q$  is equivalent to minimizing the KL divergence, because  $Q$  does not participate in the omitted term.

When computing many of these quantities, it is common to encounter expressions of the form  $0 \log 0$ . By convention, in the context of information theory, we treat these expressions as  $\lim_{x \rightarrow 0} x \log x = 0$ .

### 3.14 Structured Probabilistic Models

Machine learning algorithms often involve probability distributions over a very large number of random variables. Often, these probability distributions involve direct interactions between relatively few variables. Using a single function to

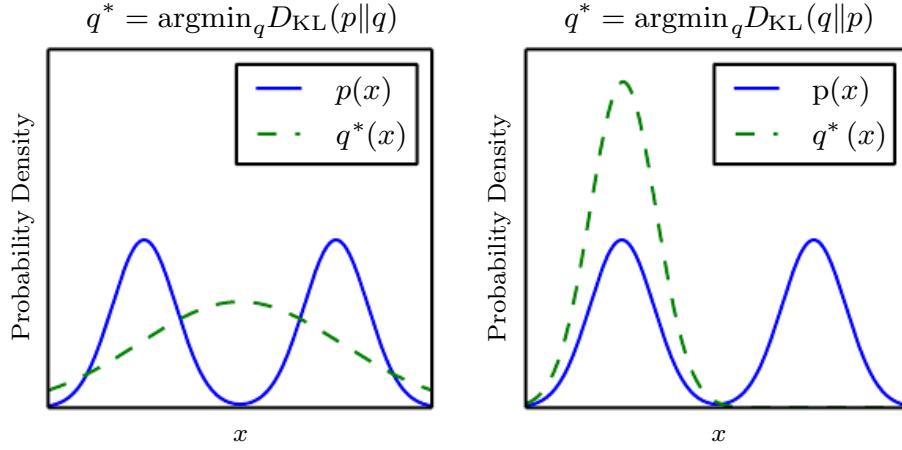


Figure 3.6: The KL divergence is asymmetric. Suppose we have a distribution  $p(x)$  and wish to approximate it with another distribution  $q(x)$ . We have the choice of minimizing either  $D_{\text{KL}}(p||q)$  or  $D_{\text{KL}}(q||p)$ . We illustrate the effect of this choice using a mixture of two Gaussians for  $p$ , and a single Gaussian for  $q$ . The choice of which direction of the KL divergence to use is problem-dependent. Some applications require an approximation that usually places high probability anywhere that the true distribution places high probability, while other applications require an approximation that rarely places high probability anywhere that the true distribution places low probability. The choice of the direction of the KL divergence reflects which of these considerations takes priority for each application. *(Left)* The effect of minimizing  $D_{\text{KL}}(p||q)$ . In this case, we select a  $q$  that has high probability where  $p$  has high probability. When  $p$  has multiple modes,  $q$  chooses to blur the modes together, in order to put high probability mass on all of them. *(Right)* The effect of minimizing  $D_{\text{KL}}(q||p)$ . In this case, we select a  $q$  that has low probability where  $p$  has low probability. When  $p$  has multiple modes that are sufficiently widely separated, as in this figure, the KL divergence is minimized by choosing a single mode, in order to avoid putting probability mass in the low-probability areas between modes of  $p$ . Here, we illustrate the outcome when  $q$  is chosen to emphasize the left mode. We could also have achieved an equal value of the KL divergence by choosing the right mode. If the modes are not separated by a sufficiently strong low probability region, then this direction of the KL divergence can still choose to blur the modes.

describe the entire joint probability distribution can be very inefficient (both computationally and statistically).

Instead of using a single function to represent a probability distribution, we can split a probability distribution into many factors that we multiply together. For example, suppose we have three random variables:  $a$ ,  $b$  and  $c$ . Suppose that  $a$  influences the value of  $b$  and  $b$  influences the value of  $c$ , but that  $a$  and  $c$  are independent given  $b$ . We can represent the probability distribution over all three variables as a product of probability distributions over two variables:

$$p(a, b, c) = p(a)p(b | a)p(c | b). \quad (3.52)$$

These factorizations can greatly reduce the number of parameters needed to describe the distribution. Each factor uses a number of parameters that is exponential in the number of variables in the factor. This means that we can greatly reduce the cost of representing a distribution if we are able to find a factorization into distributions over fewer variables.

We can describe these kinds of factorizations using graphs. Here we use the word “graph” in the sense of graph theory: a set of vertices that may be connected to each other with edges. When we represent the factorization of a probability distribution with a graph, we call it a **structured probabilistic model** or **graphical model**.

There are two main kinds of structured probabilistic models: directed and undirected. Both kinds of graphical models use a graph  $\mathcal{G}$  in which each node in the graph corresponds to a random variable, and an edge connecting two random variables means that the probability distribution is able to represent direct interactions between those two random variables.

**Directed** models use graphs with directed edges, and they represent factorizations into conditional probability distributions, as in the example above. Specifically, a directed model contains one factor for every random variable  $x_i$  in the distribution, and that factor consists of the conditional distribution over  $x_i$  given the parents of  $x_i$ , denoted  $Pa_{\mathcal{G}}(x_i)$ :

$$p(\mathbf{x}) = \prod_i p(x_i | Pa_{\mathcal{G}}(x_i)). \quad (3.53)$$

See figure 3.7 for an example of a directed graph and the factorization of probability distributions it represents.

**Undirected** models use graphs with undirected edges, and they represent factorizations into a set of functions; unlike in the directed case, these functions

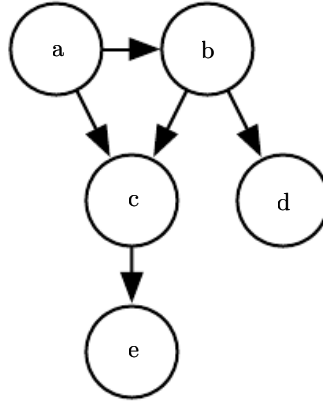


Figure 3.7: A directed graphical model over random variables  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . This graph corresponds to probability distributions that can be factored as

$$p(a, b, c, d, e) = p(a)p(b | a)p(c | a, b)p(d | b)p(e | c). \quad (3.54)$$

This graph allows us to quickly see some properties of the distribution. For example,  $a$  and  $c$  interact directly, but  $a$  and  $e$  interact only indirectly via  $c$ .

are usually not probability distributions of any kind. Any set of nodes that are all connected to each other in  $\mathcal{G}$  is called a clique. Each clique  $\mathcal{C}^{(i)}$  in an undirected model is associated with a factor  $\phi^{(i)}(\mathcal{C}^{(i)})$ . These factors are just functions, not probability distributions. The output of each factor must be non-negative, but there is no constraint that the factor must sum or integrate to 1 like a probability distribution.

The probability of a configuration of random variables is **proportional** to the product of all of these factors—assignments that result in larger factor values are more likely. Of course, there is no guarantee that this product will sum to 1. We therefore divide by a normalizing constant  $Z$ , defined to be the sum or integral over all states of the product of the  $\phi$  functions, in order to obtain a normalized probability distribution:

$$p(\mathbf{x}) = \frac{1}{Z} \prod_i \phi^{(i)}(\mathcal{C}^{(i)}). \quad (3.55)$$

See figure 3.8 for an example of an undirected graph and the factorization of probability distributions it represents.

Keep in mind that these graphical representations of factorizations are a language for describing probability distributions. They are not mutually exclusive families of probability distributions. Being directed or undirected is not a property of a probability distribution; it is a property of a particular **description** of a

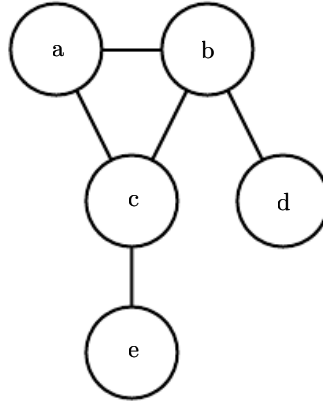


Figure 3.8: An undirected graphical model over random variables  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . This graph corresponds to probability distributions that can be factored as

$$p(a, b, c, d, e) = \frac{1}{Z} \phi^{(1)}(a, b, c) \phi^{(2)}(b, d) \phi^{(3)}(c, e). \quad (3.56)$$

This graph allows us to quickly see some properties of the distribution. For example,  $a$  and  $c$  interact directly, but  $a$  and  $e$  interact only indirectly via  $c$ .

probability distribution, but any probability distribution may be described in both ways.

Throughout parts [I](#) and [II](#) of this book, we will use structured probabilistic models merely as a language to describe which direct probabilistic relationships different machine learning algorithms choose to represent. No further understanding of structured probabilistic models is needed until the discussion of research topics, in part [III](#), where we will explore structured probabilistic models in much greater detail.

This chapter has reviewed the basic concepts of probability theory that are most relevant to deep learning. One more set of fundamental mathematical tools remains: numerical methods.

## Chapter 4

# Numerical Computation

Machine learning algorithms usually require a high amount of numerical computation. This typically refers to algorithms that solve mathematical problems by methods that update estimates of the solution via an iterative process, rather than analytically deriving a formula providing a symbolic expression for the correct solution. Common operations include optimization (finding the value of an argument that minimizes or maximizes a function) and solving systems of linear equations. Even just evaluating a mathematical function on a digital computer can be difficult when the function involves real numbers, which cannot be represented precisely using a finite amount of memory.

### 4.1 Overflow and Underflow

The fundamental difficulty in performing continuous math on a digital computer is that we need to represent infinitely many real numbers with a finite number of bit patterns. This means that for almost all real numbers, we incur some approximation error when we represent the number in the computer. In many cases, this is just rounding error. Rounding error is problematic, especially when it compounds across many operations, and can cause algorithms that work in theory to fail in practice if they are not designed to minimize the accumulation of rounding error.

One form of rounding error that is particularly devastating is **underflow**. Underflow occurs when numbers near zero are rounded to zero. Many functions behave qualitatively differently when their argument is zero rather than a small positive number. For example, we usually want to avoid division by zero (some

software environments will raise exceptions when this occurs, others will return a result with a placeholder not-a-number value) or taking the logarithm of zero (this is usually treated as  $-\infty$ , which then becomes not-a-number if it is used for many further arithmetic operations).

Another highly damaging form of numerical error is **overflow**. Overflow occurs when numbers with large magnitude are approximated as  $\infty$  or  $-\infty$ . Further arithmetic will usually change these infinite values into not-a-number values.

One example of a function that must be stabilized against underflow and overflow is the softmax function. The softmax function is often used to predict the probabilities associated with a multinoulli distribution. The softmax function is defined to be

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}. \quad (4.1)$$

Consider what happens when all of the  $x_i$  are equal to some constant  $c$ . Analytically, we can see that all of the outputs should be equal to  $\frac{1}{n}$ . Numerically, this may not occur when  $c$  has large magnitude. If  $c$  is very negative, then  $\exp(c)$  will underflow. This means the denominator of the softmax will become 0, so the final result is undefined. When  $c$  is very large and positive,  $\exp(c)$  will overflow, again resulting in the expression as a whole being undefined. Both of these difficulties can be resolved by instead evaluating  $\text{softmax}(\mathbf{z})$  where  $\mathbf{z} = \mathbf{x} - \max_i x_i$ . Simple algebra shows that the value of the softmax function is not changed analytically by adding or subtracting a scalar from the input vector. Subtracting  $\max_i x_i$  results in the largest argument to  $\exp$  being 0, which rules out the possibility of overflow. Likewise, at least one term in the denominator has a value of 1, which rules out the possibility of underflow in the denominator leading to a division by zero.

There is still one small problem. Underflow in the numerator can still cause the expression as a whole to evaluate to zero. This means that if we implement  $\log \text{softmax}(\mathbf{x})$  by first running the softmax subroutine then passing the result to the log function, we could erroneously obtain  $-\infty$ . Instead, we must implement a separate function that calculates  $\log \text{softmax}$  in a numerically stable way. The  $\log \text{softmax}$  function can be stabilized using the same trick as we used to stabilize the softmax function.

For the most part, we do not explicitly detail all of the numerical considerations involved in implementing the various algorithms described in this book. Developers of low-level libraries should keep numerical issues in mind when implementing deep learning algorithms. Most readers of this book can simply rely on low-level libraries that provide stable implementations. In some cases, it is possible to implement a new algorithm and have the new implementation automatically



stabilized. Theano (Bergstra *et al.*, 2010; Bastien *et al.*, 2012) is an example of a software package that automatically detects and stabilizes many common numerically unstable expressions that arise in the context of deep learning.

## 4.2 Poor Conditioning

Conditioning refers to how rapidly a function changes with respect to small changes in its inputs. Functions that change rapidly when their inputs are perturbed slightly can be problematic for scientific computation because rounding errors in the inputs can result in large changes in the output.

Consider the function  $f(\mathbf{x}) = \mathbf{A}^{-1}\mathbf{x}$ . When  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has an eigenvalue decomposition, its **condition number** is

$$\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|. \quad (4.2)$$

This is the ratio of the magnitude of the largest and smallest eigenvalue. When this number is large, matrix inversion is particularly sensitive to error in the input.

This sensitivity is an intrinsic property of the matrix itself, not the result of rounding error during matrix inversion. Poorly conditioned matrices amplify pre-existing errors when we multiply by the true matrix inverse. In practice, the error will be compounded further by numerical errors in the inversion process itself.

## 4.3 Gradient-Based Optimization

Most deep learning algorithms involve optimization of some sort. Optimization refers to the task of either minimizing or maximizing some function  $f(\mathbf{x})$  by altering  $\mathbf{x}$ . We usually phrase most optimization problems in terms of minimizing  $f(\mathbf{x})$ . Maximization may be accomplished via a minimization algorithm by minimizing  $-f(\mathbf{x})$ .

The function we want to minimize or maximize is called the **objective function** or **criterion**. When we are minimizing it, we may also call it the **cost function**, **loss function**, or **error function**. In this book, we use these terms interchangeably, though some machine learning publications assign special meaning to some of these terms.

We often denote the value that minimizes or maximizes a function with a superscript  $*$ . For example, we might say  $\mathbf{x}^* = \arg \min f(\mathbf{x})$ .

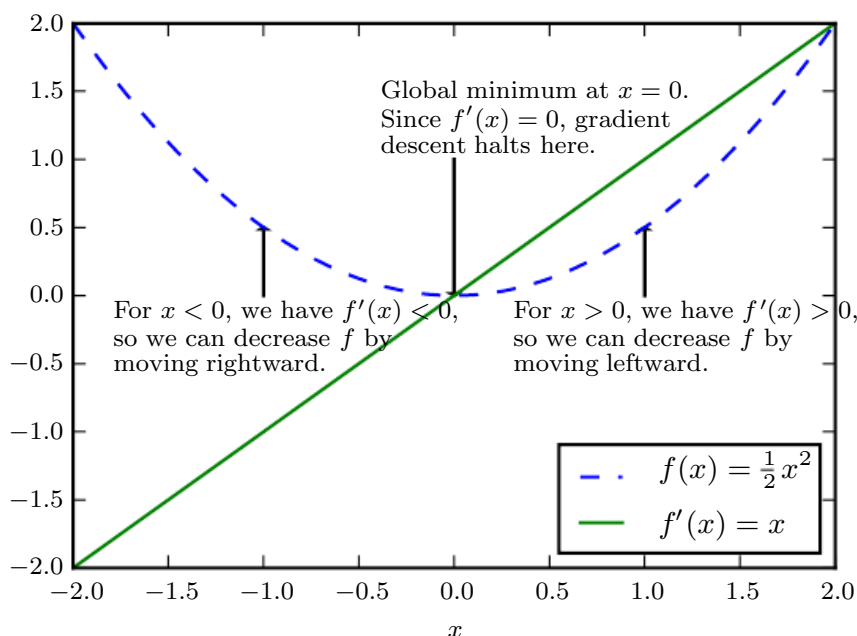


Figure 4.1: An illustration of how the gradient descent algorithm uses the derivatives of a function can be used to follow the function downhill to a minimum.

We assume the reader is already familiar with calculus, but provide a brief review of how calculus concepts relate to optimization here.

Suppose we have a function  $y = f(x)$ , where both  $x$  and  $y$  are real numbers. The **derivative** of this function is denoted as  $f'(x)$  or as  $\frac{dy}{dx}$ . The derivative  $f'(x)$  gives the slope of  $f(x)$  at the point  $x$ . In other words, it specifies how to scale a small change in the input in order to obtain the corresponding change in the output:  $f(x + \epsilon) \approx f(x) + \epsilon f'(x)$ .

The derivative is therefore useful for minimizing a function because it tells us how to change  $x$  in order to make a small improvement in  $y$ . For example, we know that  $f(x - \epsilon \text{sign}(f'(x)))$  is less than  $f(x)$  for small enough  $\epsilon$ . We can thus reduce  $f(x)$  by moving  $x$  in small steps with opposite sign of the derivative. This technique is called **gradient descent** (Cauchy, 1847). See figure 4.1 for an example of this technique.

When  $f'(x) = 0$ , the derivative provides no information about which direction to move. Points where  $f'(x) = 0$  are known as **critical points** or **stationary points**. A **local minimum** is a point where  $f(x)$  is lower than at all neighboring points, so it is no longer possible to decrease  $f(x)$  by making infinitesimal steps. A **local maximum** is a point where  $f(x)$  is higher than at all neighboring points,

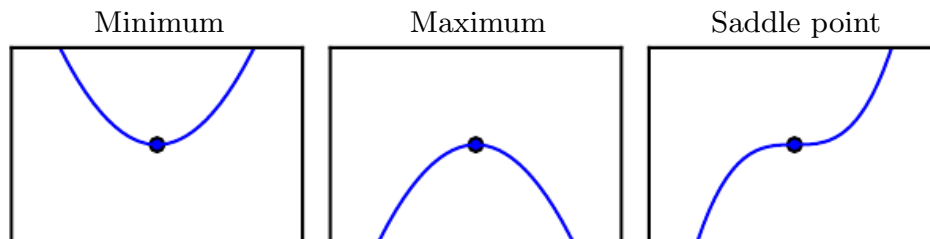


Figure 4.2: Examples of each of the three types of critical points in 1-D. A critical point is a point with zero slope. Such a point can either be a local minimum, which is lower than the neighboring points, a local maximum, which is higher than the neighboring points, or a saddle point, which has neighbors that are both higher and lower than the point itself.

so it is not possible to increase  $f(x)$  by making infinitesimal steps. Some critical points are neither maxima nor minima. These are known as **saddle points**. See figure 4.2 for examples of each type of critical point.

A point that obtains the absolute lowest value of  $f(x)$  is a **global minimum**. It is possible for there to be only one global minimum or multiple global minima of the function. It is also possible for there to be local minima that are not globally optimal. In the context of deep learning, we optimize functions that may have many local minima that are not optimal, and many saddle points surrounded by very flat regions. All of this makes optimization very difficult, especially when the input to the function is multidimensional. We therefore usually settle for finding a value of  $f$  that is very low, but not necessarily minimal in any formal sense. See figure 4.3 for an example.

We often minimize functions that have multiple inputs:  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . For the concept of “minimization” to make sense, there must still be only one (scalar) output.

For functions with multiple inputs, we must make use of the concept of **partial derivatives**. The partial derivative  $\frac{\partial}{\partial x_i} f(\mathbf{x})$  measures how  $f$  changes as only the variable  $x_i$  increases at point  $\mathbf{x}$ . The **gradient** generalizes the notion of derivative to the case where the derivative is with respect to a vector: the gradient of  $f$  is the vector containing all of the partial derivatives, denoted  $\nabla_{\mathbf{x}} f(\mathbf{x})$ . Element  $i$  of the gradient is the partial derivative of  $f$  with respect to  $x_i$ . In multiple dimensions,

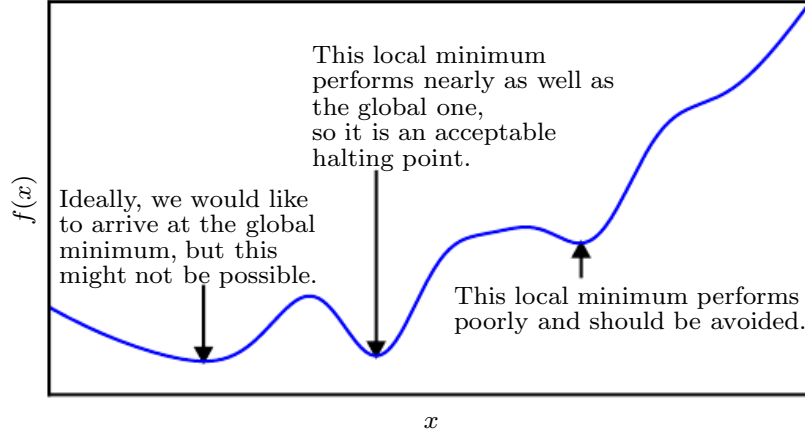


Figure 4.3: Optimization algorithms may fail to find a global minimum when there are multiple local minima or plateaus present. In the context of deep learning, we generally accept such solutions even though they are not truly minimal, so long as they correspond to significantly low values of the cost function.

critical points are points where every element of the gradient is equal to zero.

The **directional derivative** in direction  $\mathbf{u}$  (a unit vector) is the slope of the function  $f$  in direction  $\mathbf{u}$ . In other words, the directional derivative is the derivative of the function  $f(\mathbf{x} + \alpha\mathbf{u})$  with respect to  $\alpha$ , evaluated at  $\alpha = 0$ . Using the chain rule, we can see that  $\frac{\partial}{\partial \alpha} f(\mathbf{x} + \alpha\mathbf{u})$  evaluates to  $\mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$  when  $\alpha = 0$ .

To minimize  $f$ , we would like to find the direction in which  $f$  decreases the fastest. We can do this using the directional derivative:

$$\min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (4.3)$$

$$= \min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \|\mathbf{u}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos \theta \quad (4.4)$$

where  $\theta$  is the angle between  $\mathbf{u}$  and the gradient. Substituting in  $\|\mathbf{u}\|_2 = 1$  and ignoring factors that do not depend on  $\mathbf{u}$ , this simplifies to  $\min_{\mathbf{u}} \cos \theta$ . This is minimized when  $\mathbf{u}$  points in the opposite direction as the gradient. In other words, the gradient points directly uphill, and the negative gradient points directly downhill. We can decrease  $f$  by moving in the direction of the negative gradient. This is known as the **method of steepest descent** or **gradient descent**.

Steepest descent proposes a new point

$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (4.5)$$

where  $\epsilon$  is the **learning rate**, a positive scalar determining the size of the step. We can choose  $\epsilon$  in several different ways. A popular approach is to set  $\epsilon$  to a small constant. Sometimes, we can solve for the step size that makes the directional derivative vanish. Another approach is to evaluate  $f(\mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}))$  for several values of  $\epsilon$  and choose the one that results in the smallest objective function value. This last strategy is called a **line search**.

Steepest descent converges when every element of the gradient is zero (or, in practice, very close to zero). In some cases, we may be able to avoid running this iterative algorithm, and just jump directly to the critical point by solving the equation  $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$  for  $\mathbf{x}$ .

Although gradient descent is limited to optimization in continuous spaces, the general concept of repeatedly making a small move (that is approximately the best small move) towards better configurations can be generalized to discrete spaces. Ascending an objective function of discrete parameters is called **hill climbing** (Russel and Norvig, 2003).

### 4.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

Sometimes we need to find all of the partial derivatives of a function whose input and output are both vectors. The matrix containing all such partial derivatives is known as a **Jacobian matrix**. Specifically, if we have a function  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ , then the Jacobian matrix  $\mathbf{J} \in \mathbb{R}^{n \times m}$  of  $\mathbf{f}$  is defined such that  $J_{i,j} = \frac{\partial}{\partial x_j} f(\mathbf{x})_i$ .

We are also sometimes interested in a derivative of a derivative. This is known as a **second derivative**. For example, for a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the derivative with respect to  $x_i$  of the derivative of  $f$  with respect to  $x_j$  is denoted as  $\frac{\partial^2}{\partial x_i \partial x_j} f$ . In a single dimension, we can denote  $\frac{d^2}{dx^2} f$  by  $f''(x)$ . The second derivative tells us how the first derivative will change as we vary the input. This is important because it tells us whether a gradient step will cause as much of an improvement as we would expect based on the gradient alone. We can think of the second derivative as measuring **curvature**. Suppose we have a quadratic function (many functions that arise in practice are not quadratic but can be approximated well as quadratic, at least locally). If such a function has a second derivative of zero, then there is no curvature. It is a perfectly flat line, and its value can be predicted using only the gradient. If the gradient is 1, then we can make a step of size  $\epsilon$  along the negative gradient, and the cost function will decrease by  $\epsilon$ . If the second derivative is negative, the function curves downward, so the cost function will actually decrease by more than  $\epsilon$ . Finally, if the second derivative is positive, the function curves upward, so the cost function can decrease by less than  $\epsilon$ . See

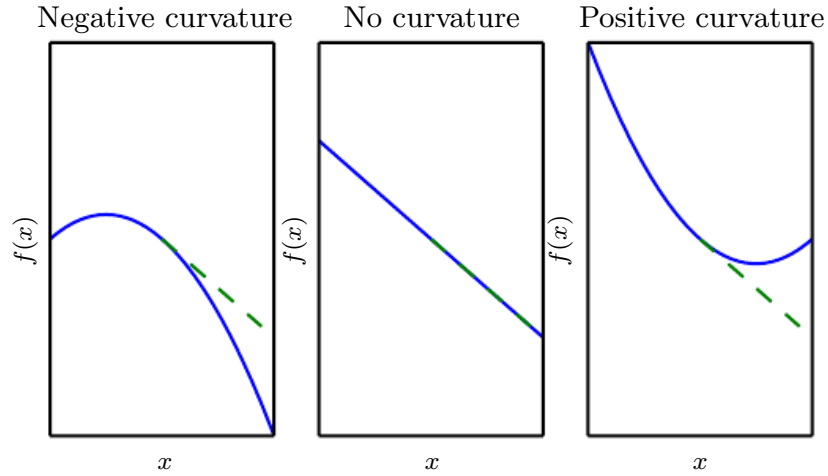


Figure 4.4: The second derivative determines the curvature of a function. Here we show quadratic functions with various curvature. The dashed line indicates the value of the cost function we would expect based on the gradient information alone as we make a gradient step downhill. In the case of negative curvature, the cost function actually decreases faster than the gradient predicts. In the case of no curvature, the gradient predicts the decrease correctly. In the case of positive curvature, the function decreases slower than expected and eventually begins to increase, so steps that are too large can actually increase the function inadvertently.

figure 4.4 to see how different forms of curvature affect the relationship between the value of the cost function predicted by the gradient and the true value.

When our function has multiple input dimensions, there are many second derivatives. These derivatives can be collected together into a matrix called the **Hessian matrix**. The Hessian matrix  $\mathbf{H}(f)(\mathbf{x})$  is defined such that

$$\mathbf{H}(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}). \quad (4.6)$$

Equivalently, the Hessian is the Jacobian of the gradient.

Anywhere that the second partial derivatives are continuous, the differential operators are commutative, i.e. their order can be swapped:

$$\frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}) = \frac{\partial^2}{\partial x_j \partial x_i} f(\mathbf{x}). \quad (4.7)$$

This implies that  $H_{i,j} = H_{j,i}$ , so the Hessian matrix is symmetric at such points. Most of the functions we encounter in the context of deep learning have a symmetric Hessian almost everywhere. Because the Hessian matrix is real and symmetric, we can decompose it into a set of real eigenvalues and an orthogonal basis of

eigenvectors. The second derivative in a specific direction represented by a unit vector  $\mathbf{d}$  is given by  $\mathbf{d}^\top \mathbf{H} \mathbf{d}$ . When  $\mathbf{d}$  is an eigenvector of  $\mathbf{H}$ , the second derivative in that direction is given by the corresponding eigenvalue. For other directions of  $\mathbf{d}$ , the directional second derivative is a weighted average of all of the eigenvalues, with weights between 0 and 1, and eigenvectors that have smaller angle with  $\mathbf{d}$  receiving more weight. The maximum eigenvalue determines the maximum second derivative and the minimum eigenvalue determines the minimum second derivative.

The (directional) second derivative tells us how well we can expect a gradient descent step to perform. We can make a second-order Taylor series approximation to the function  $f(\mathbf{x})$  around the current point  $\mathbf{x}^{(0)}$ :

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(\mathbf{x} - \mathbf{x}^{(0)}). \quad (4.8)$$

where  $\mathbf{g}$  is the gradient and  $\mathbf{H}$  is the Hessian at  $\mathbf{x}^{(0)}$ . If we use a learning rate of  $\epsilon$ , then the new point  $\mathbf{x}$  will be given by  $\mathbf{x}^{(0)} - \epsilon \mathbf{g}$ . Substituting this into our approximation, we obtain

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}. \quad (4.9)$$

There are three terms here: the original value of the function, the expected improvement due to the slope of the function, and the correction we must apply to account for the curvature of the function. When this last term is too large, the gradient descent step can actually move uphill. When  $\mathbf{g}^\top \mathbf{H} \mathbf{g}$  is zero or negative, the Taylor series approximation predicts that increasing  $\epsilon$  forever will decrease  $f$  forever. In practice, the Taylor series is unlikely to remain accurate for large  $\epsilon$ , so one must resort to more heuristic choices of  $\epsilon$  in this case. When  $\mathbf{g}^\top \mathbf{H} \mathbf{g}$  is positive, solving for the optimal step size that decreases the Taylor series approximation of the function the most yields

$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}. \quad (4.10)$$

In the worst case, when  $\mathbf{g}$  aligns with the eigenvector of  $\mathbf{H}$  corresponding to the maximal eigenvalue  $\lambda_{\max}$ , then this optimal step size is given by  $\frac{1}{\lambda_{\max}}$ . To the extent that the function we minimize can be approximated well by a quadratic function, the eigenvalues of the Hessian thus determine the scale of the learning rate.

The second derivative can be used to determine whether a critical point is a local maximum, a local minimum, or saddle point. Recall that on a critical point,  $f'(x) = 0$ . When the second derivative  $f''(x) > 0$ , the first derivative  $f'(x)$  increases as we move to the right and decreases as we move to the left. This means

$f'(x - \epsilon) < 0$  and  $f'(x + \epsilon) > 0$  for small enough  $\epsilon$ . In other words, as we move right, the slope begins to point uphill to the right, and as we move left, the slope begins to point uphill to the left. Thus, when  $f'(x) = 0$  and  $f''(x) > 0$ , we can conclude that  $x$  is a local minimum. Similarly, when  $f'(x) = 0$  and  $f''(x) < 0$ , we can conclude that  $x$  is a local maximum. This is known as the **second derivative test**. Unfortunately, when  $f''(x) = 0$ , the test is inconclusive. In this case  $x$  may be a saddle point, or a part of a flat region.

In multiple dimensions, we need to examine all of the second derivatives of the function. Using the eigendecomposition of the Hessian matrix, we can generalize the second derivative test to multiple dimensions. At a critical point, where  $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ , we can examine the eigenvalues of the Hessian to determine whether the critical point is a local maximum, local minimum, or saddle point. When the Hessian is positive definite (all its eigenvalues are positive), the point is a local minimum. This can be seen by observing that the directional second derivative in any direction must be positive, and making reference to the univariate second derivative test. Likewise, when the Hessian is negative definite (all its eigenvalues are negative), the point is a local maximum. In multiple dimensions, it is actually possible to find positive evidence of saddle points in some cases. When at least one eigenvalue is positive and at least one eigenvalue is negative, we know that  $\mathbf{x}$  is a local maximum on one cross section of  $f$  but a local minimum on another cross section. See figure 4.5 for an example. Finally, the multidimensional second derivative test can be inconclusive, just like the univariate version. The test is inconclusive whenever all of the non-zero eigenvalues have the same sign, but at least one eigenvalue is zero. This is because the univariate second derivative test is inconclusive in the cross section corresponding to the zero eigenvalue.

In multiple dimensions, there is a different second derivative for each direction at a single point. The condition number of the Hessian at this point measures how much the second derivatives differ from each other. When the Hessian has a poor condition number, gradient descent performs poorly. This is because in one direction, the derivative increases rapidly, while in another direction, it increases slowly. Gradient descent is unaware of this change in the derivative so it does not know that it needs to explore preferentially in the direction where the derivative remains negative for longer. It also makes it difficult to choose a good step size. The step size must be small enough to avoid overshooting the minimum and going uphill in directions with strong positive curvature. This usually means that the step size is too small to make significant progress in other directions with less curvature. See figure 4.6 for an example.

This issue can be resolved by using information from the Hessian matrix to guide



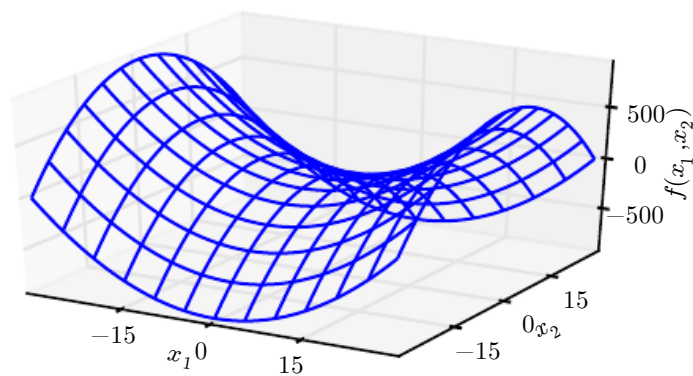


Figure 4.5: A saddle point containing both positive and negative curvature. The function in this example is  $f(\mathbf{x}) = x_1^2 - x_2^2$ . Along the axis corresponding to  $x_1$ , the function curves upward. This axis is an eigenvector of the Hessian and has a positive eigenvalue. Along the axis corresponding to  $x_2$ , the function curves downward. This direction is an eigenvector of the Hessian with negative eigenvalue. The name “saddle point” derives from the saddle-like shape of this function. This is the quintessential example of a function with a saddle point. In more than one dimension, it is not necessary to have an eigenvalue of 0 in order to get a saddle point: it is only necessary to have both positive and negative eigenvalues. We can think of a saddle point with both signs of eigenvalues as being a local maximum within one cross section and a local minimum within another cross section.

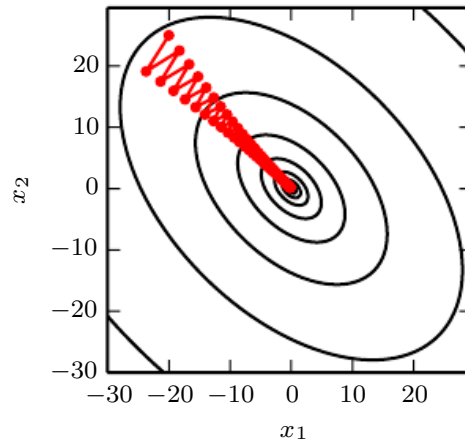


Figure 4.6: Gradient descent fails to exploit the curvature information contained in the Hessian matrix. Here we use gradient descent to minimize a quadratic function  $f(\mathbf{x})$  whose Hessian matrix has condition number 5. This means that the direction of most curvature has five times more curvature than the direction of least curvature. In this case, the most curvature is in the direction  $[1, 1]^\top$  and the least curvature is in the direction  $[1, -1]^\top$ . The red lines indicate the path followed by gradient descent. This very elongated quadratic function resembles a long canyon. Gradient descent wastes time repeatedly descending canyon walls, because they are the steepest feature. Because the step size is somewhat too large, it has a tendency to overshoot the bottom of the function and thus needs to descend the opposite canyon wall on the next iteration. The large positive eigenvalue of the Hessian corresponding to the eigenvector pointed in this direction indicates that this directional derivative is rapidly increasing, so an optimization algorithm based on the Hessian could predict that the steepest direction is not actually a promising search direction in this context.

the search. The simplest method for doing so is known as **Newton's method**. Newton's method is based on using a second-order Taylor series expansion to approximate  $f(\mathbf{x})$  near some point  $\mathbf{x}^{(0)}$ :

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(f)(\mathbf{x}^{(0)}) (\mathbf{x} - \mathbf{x}^{(0)}). \quad (4.11)$$

If we then solve for the critical point of this function, we obtain:

$$\mathbf{x}^* = \mathbf{x}^{(0)} - \mathbf{H}(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}). \quad (4.12)$$

When  $f$  is a positive definite quadratic function, Newton's method consists of applying equation 4.12 once to jump to the minimum of the function directly. When  $f$  is not truly quadratic but can be locally approximated as a positive definite quadratic, Newton's method consists of applying equation 4.12 multiple times. Iteratively updating the approximation and jumping to the minimum of the approximation can reach the critical point much faster than gradient descent would. This is a useful property near a local minimum, but it can be a harmful property near a saddle point. As discussed in section 8.2.3, Newton's method is only appropriate when the nearby critical point is a minimum (all the eigenvalues of the Hessian are positive), whereas gradient descent is not attracted to saddle points unless the gradient points toward them.

Optimization algorithms that use only the gradient, such as gradient descent, are called **first-order optimization algorithms**. Optimization algorithms that also use the Hessian matrix, such as Newton's method, are called **second-order optimization algorithms** (Nocedal and Wright, 2006).

The optimization algorithms employed in most contexts in this book are applicable to a wide variety of functions, but come with almost no guarantees. Deep learning algorithms tend to lack guarantees because the family of functions used in deep learning is quite complicated. In many other fields, the dominant approach to optimization is to design optimization algorithms for a limited family of functions.

In the context of deep learning, we sometimes gain some guarantees by restricting ourselves to functions that are either **Lipschitz continuous** or have Lipschitz continuous derivatives. A Lipschitz continuous function is a function  $f$  whose rate of change is bounded by a **Lipschitz constant**  $\mathcal{L}$ :

$$\forall \mathbf{x}, \forall \mathbf{y}, |f(\mathbf{x}) - f(\mathbf{y})| \leq \mathcal{L} \|\mathbf{x} - \mathbf{y}\|_2. \quad (4.13)$$

This property is useful because it allows us to quantify our assumption that a small change in the input made by an algorithm such as gradient descent will have

a small change in the output. Lipschitz continuity is also a fairly weak constraint, and many optimization problems in deep learning can be made Lipschitz continuous with relatively minor modifications.

Perhaps the most successful field of specialized optimization is **convex optimization**. Convex optimization algorithms are able to provide many more guarantees by making stronger restrictions. Convex optimization algorithms are applicable only to convex functions—functions for which the Hessian is positive semidefinite everywhere. Such functions are well-behaved because they lack saddle points and all of their local minima are necessarily global minima. However, most problems in deep learning are difficult to express in terms of convex optimization. Convex optimization is used only as a subroutine of some deep learning algorithms. Ideas from the analysis of convex optimization algorithms can be useful for proving the convergence of deep learning algorithms. However, in general, the importance of convex optimization is greatly diminished in the context of deep learning. For more information about convex optimization, see [Boyd and Vandenberghe \(2004\)](#) or [Rockafellar \(1997\)](#).

## 4.4 Constrained Optimization

Sometimes we wish not only to maximize or minimize a function  $f(\mathbf{x})$  over all possible values of  $\mathbf{x}$ . Instead we may wish to find the maximal or minimal value of  $f(\mathbf{x})$  for values of  $\mathbf{x}$  in some set  $\mathbb{S}$ . This is known as **constrained optimization**. Points  $\mathbf{x}$  that lie within the set  $\mathbb{S}$  are called **feasible** points in constrained optimization terminology.

We often wish to find a solution that is small in some sense. A common approach in such situations is to impose a norm constraint, such as  $\|\mathbf{x}\| \leq 1$ .

One simple approach to constrained optimization is simply to modify gradient descent taking the constraint into account. If we use a small constant step size  $\epsilon$ , we can make gradient descent steps, then project the result back into  $\mathbb{S}$ . If we use a line search, we can search only over step sizes  $\epsilon$  that yield new  $\mathbf{x}$  points that are feasible, or we can project each point on the line back into the constraint region. When possible, this method can be made more efficient by projecting the gradient into the tangent space of the feasible region before taking the step or beginning the line search ([Rosen, 1960](#)).

A more sophisticated approach is to design a different, unconstrained optimization problem whose solution can be converted into a solution to the original, constrained optimization problem. For example, if we want to minimize  $f(\mathbf{x})$  for

$\mathbf{x} \in \mathbb{R}^2$  with  $\mathbf{x}$  constrained to have exactly unit  $L^2$  norm, we can instead minimize  $g(\theta) = f([\cos \theta, \sin \theta]^\top)$  with respect to  $\theta$ , then return  $[\cos \theta, \sin \theta]$  as the solution to the original problem. This approach requires creativity; the transformation between optimization problems must be designed specifically for each case we encounter.

The **Karush–Kuhn–Tucker** (KKT) approach<sup>1</sup> provides a very general solution to constrained optimization. With the KKT approach, we introduce a new function called the **generalized Lagrangian** or **generalized Lagrange function**.

To define the Lagrangian, we first need to describe  $\mathbb{S}$  in terms of equations and inequalities. We want a description of  $\mathbb{S}$  in terms of  $m$  functions  $g^{(i)}$  and  $n$  functions  $h^{(j)}$  so that  $\mathbb{S} = \{\mathbf{x} \mid \forall i, g^{(i)}(\mathbf{x}) = 0 \text{ and } \forall j, h^{(j)}(\mathbf{x}) \leq 0\}$ . The equations involving  $g^{(i)}$  are called the **equality constraints** and the inequalities involving  $h^{(j)}$  are called **inequality constraints**.

We introduce new variables  $\lambda_i$  and  $\alpha_j$  for each constraint, these are called the KKT multipliers. The generalized Lagrangian is then defined as

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_i \lambda_i g^{(i)}(\mathbf{x}) + \sum_j \alpha_j h^{(j)}(\mathbf{x}). \quad (4.14)$$

We can now solve a constrained minimization problem using unconstrained optimization of the generalized Lagrangian. Observe that, so long as at least one feasible point exists and  $f(\mathbf{x})$  is not permitted to have value  $\infty$ , then

$$\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}). \quad (4.15)$$

has the same optimal objective function value and set of optimal points  $\mathbf{x}$  as

$$\min_{\mathbf{x} \in \mathbb{S}} f(\mathbf{x}). \quad (4.16)$$

This follows because any time the constraints are satisfied,

$$\max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}), \quad (4.17)$$

while any time a constraint is violated,

$$\max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = \infty. \quad (4.18)$$

---

<sup>1</sup>The KKT approach generalizes the method of **Lagrange multipliers** which allows equality constraints but not inequality constraints.

These properties guarantee that no infeasible point can be optimal, and that the optimum within the feasible points is unchanged.

To perform constrained maximization, we can construct the generalized Lagrange function of  $-f(\mathbf{x})$ , which leads to this optimization problem:

$$\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} -f(\mathbf{x}) + \sum_i \lambda_i g^{(i)}(\mathbf{x}) + \sum_j \alpha_j h^{(j)}(\mathbf{x}). \quad (4.19)$$

We may also convert this to a problem with maximization in the outer loop:

$$\max_{\mathbf{x}} \min_{\boldsymbol{\lambda}} \min_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} f(\mathbf{x}) + \sum_i \lambda_i g^{(i)}(\mathbf{x}) - \sum_j \alpha_j h^{(j)}(\mathbf{x}). \quad (4.20)$$

The sign of the term for the equality constraints does not matter; we may define it with addition or subtraction as we wish, because the optimization is free to choose any sign for each  $\lambda_i$ .

The inequality constraints are particularly interesting. We say that a constraint  $h^{(i)}(\mathbf{x})$  is **active** if  $h^{(i)}(\mathbf{x}^*) = 0$ . If a constraint is not active, then the solution to the problem found using that constraint would remain at least a local solution if that constraint were removed. It is possible that an inactive constraint excludes other solutions. For example, a convex problem with an entire region of globally optimal points (a wide, flat, region of equal cost) could have a subset of this region eliminated by constraints, or a non-convex problem could have better local stationary points excluded by a constraint that is inactive at convergence. However, the point found at convergence remains a stationary point whether or not the inactive constraints are included. Because an inactive  $h^{(i)}$  has negative value, then the solution to  $\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha})$  will have  $\alpha_i = 0$ . We can thus observe that at the solution,  $\boldsymbol{\alpha} \odot \mathbf{h}(\mathbf{x}) = \mathbf{0}$ . In other words, for all  $i$ , we know that at least one of the constraints  $\alpha_i \geq 0$  and  $h^{(i)}(\mathbf{x}) \leq 0$  must be active at the solution. To gain some intuition for this idea, we can say that either the solution is on the boundary imposed by the inequality and we must use its KKT multiplier to influence the solution to  $\mathbf{x}$ , or the inequality has no influence on the solution and we represent this by zeroing out its KKT multiplier.

A simple set of properties describe the optimal points of constrained optimization problems. These properties are called the Karush-Kuhn-Tucker (KKT) conditions ([Karush, 1939](#); [Kuhn and Tucker, 1951](#)). They are necessary conditions, but not always sufficient conditions, for a point to be optimal. The conditions are:

- The gradient of the generalized Lagrangian is zero.
- All constraints on both  $\mathbf{x}$  and the KKT multipliers are satisfied.

- The inequality constraints exhibit “complementary slackness”:  $\alpha \odot \mathbf{h}(\mathbf{x}) = \mathbf{0}$ .

For more information about the KKT approach, see [Nocedal and Wright \(2006\)](#).

## 4.5 Example: Linear Least Squares

Suppose we want to find the value of  $\mathbf{x}$  that minimizes

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2. \quad (4.21)$$

There are specialized linear algebra algorithms that can solve this problem efficiently. However, we can also explore how to solve it using gradient-based optimization as a simple example of how these techniques work.

First, we need to obtain the gradient:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}) = \mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top \mathbf{b}. \quad (4.22)$$

We can then follow this gradient downhill, taking small steps. See algorithm 4.1 for details.

---

**Algorithm 4.1** An algorithm to minimize  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$  with respect to  $\mathbf{x}$  using gradient descent, starting from an arbitrary value of  $\mathbf{x}$ .

---

Set the step size ( $\epsilon$ ) and tolerance ( $\delta$ ) to small, positive numbers.

**while**  $\|\mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top \mathbf{b}\|_2 > \delta$  **do**  
     $\mathbf{x} \leftarrow \mathbf{x} - \epsilon (\mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top \mathbf{b})$   
**end while**

---

One can also solve this problem using Newton’s method. In this case, because the true function is quadratic, the quadratic approximation employed by Newton’s method is exact, and the algorithm converges to the global minimum in a single step.

Now suppose we wish to minimize the same function, but subject to the constraint  $\mathbf{x}^\top \mathbf{x} \leq 1$ . To do so, we introduce the Lagrangian

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda (\mathbf{x}^\top \mathbf{x} - 1). \quad (4.23)$$

We can now solve the problem

$$\min_{\mathbf{x}} \max_{\lambda, \lambda \geq 0} L(\mathbf{x}, \lambda). \quad (4.24)$$

The smallest-norm solution to the unconstrained least squares problem may be found using the Moore-Penrose pseudoinverse:  $\mathbf{x} = \mathbf{A}^+ \mathbf{b}$ . If this point is feasible, then it is the solution to the constrained problem. Otherwise, we must find a solution where the constraint is active. By differentiating the Lagrangian with respect to  $\mathbf{x}$ , we obtain the equation

$$\mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{A}^\top \mathbf{b} + 2\lambda \mathbf{x} = 0. \quad (4.25)$$

This tells us that the solution will take the form

$$\mathbf{x} = (\mathbf{A}^\top \mathbf{A} + 2\lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{b}. \quad (4.26)$$

The magnitude of  $\lambda$  must be chosen such that the result obeys the constraint. We can find this value by performing gradient ascent on  $\lambda$ . To do so, observe

$$\frac{\partial}{\partial \lambda} L(\mathbf{x}, \lambda) = \mathbf{x}^\top \mathbf{x} - 1. \quad (4.27)$$

When the norm of  $\mathbf{x}$  exceeds 1, this derivative is positive, so to follow the derivative uphill and increase the Lagrangian with respect to  $\lambda$ , we increase  $\lambda$ . Because the coefficient on the  $\mathbf{x}^\top \mathbf{x}$  penalty has increased, solving the linear equation for  $\mathbf{x}$  will now yield a solution with smaller norm. The process of solving the linear equation and adjusting  $\lambda$  continues until  $\mathbf{x}$  has the correct norm and the derivative on  $\lambda$  is 0.

This concludes the mathematical preliminaries that we use to develop machine learning algorithms. We are now ready to build and analyze some full-fledged learning systems.



## Chapter 5

# Machine Learning Basics

Deep learning is a specific kind of machine learning. In order to understand deep learning well, one must have a solid understanding of the basic principles of machine learning. This chapter provides a brief course in the most important general principles that will be applied throughout the rest of the book. Novice readers or those who want a wider perspective are encouraged to consider machine learning textbooks with a more comprehensive coverage of the fundamentals, such as [Murphy \(2012\)](#) or [Bishop \(2006\)](#). If you are already familiar with machine learning basics, feel free to skip ahead to section [5.11](#). That section covers some perspectives on traditional machine learning techniques that have strongly influenced the development of deep learning algorithms.

We begin with a definition of what a learning algorithm is, and present an example: the linear regression algorithm. We then proceed to describe how the challenge of fitting the training data differs from the challenge of finding patterns that generalize to new data. Most machine learning algorithms have settings called hyperparameters that must be determined external to the learning algorithm itself; we discuss how to set these using additional data. Machine learning is essentially a form of applied statistics with increased emphasis on the use of computers to statistically estimate complicated functions and a decreased emphasis on proving confidence intervals around these functions; we therefore present the two central approaches to statistics: frequentist estimators and Bayesian inference. Most machine learning algorithms can be divided into the categories of supervised learning and unsupervised learning; we describe these categories and give some examples of simple learning algorithms from each category. Most deep learning algorithms are based on an optimization algorithm called stochastic gradient descent. We describe how to combine various algorithm components such as

an optimization algorithm, a cost function, a model, and a dataset to build a machine learning algorithm. Finally, in section 5.11, we describe some of the factors that have limited the ability of traditional machine learning to generalize. These challenges have motivated the development of deep learning algorithms that overcome these obstacles.

## 5.1 Learning Algorithms

A machine learning algorithm is an algorithm that is able to learn from data. But what do we mean by learning? Mitchell (1997) provides the definition “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .” One can imagine a very wide variety of experiences  $E$ , tasks  $T$ , and performance measures  $P$ , and we do not make any attempt in this book to provide a formal definition of what may be used for each of these entities. Instead, the following sections provide intuitive descriptions and examples of the different kinds of tasks, performance measures and experiences that can be used to construct machine learning algorithms.

### 5.1.1 The Task, $T$

Machine learning allows us to tackle tasks that are too difficult to solve with fixed programs written and designed by human beings. From a scientific and philosophical point of view, machine learning is interesting because developing our understanding of machine learning entails developing our understanding of the principles that underlie intelligence.

In this relatively formal definition of the word “task,” the process of learning itself is not the task. Learning is our means of attaining the ability to perform the task. For example, if we want a robot to be able to walk, then walking is the task. We could program the robot to learn to walk, or we could attempt to directly write a program that specifies how to walk manually.

Machine learning tasks are usually described in terms of how the machine learning system should process an **example**. An example is a collection of **features** that have been quantitatively measured from some object or event that we want the machine learning system to process. We typically represent an example as a vector  $\mathbf{x} \in \mathbb{R}^n$  where each entry  $x_i$  of the vector is another feature. For example, the features of an image are usually the values of the pixels in the image.

Many kinds of tasks can be solved with machine learning. Some of the most common machine learning tasks include the following:

- **Classification:** In this type of task, the computer program is asked to specify which of  $k$  categories some input belongs to. To solve this task, the learning algorithm is usually asked to produce a function  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$ . When  $y = f(\mathbf{x})$ , the model assigns an input described by vector  $\mathbf{x}$  to a category identified by numeric code  $y$ . There are other variants of the classification task, for example, where  $f$  outputs a probability distribution over classes. An example of a classification task is object recognition, where the input is an image (usually described as a set of pixel brightness values), and the output is a numeric code identifying the object in the image. For example, the Willow Garage PR2 robot is able to act as a waiter that can recognize different kinds of drinks and deliver them to people on command (Goodfellow *et al.*, 2010). Modern object recognition is best accomplished with deep learning (Krizhevsky *et al.*, 2012; Ioffe and Szegedy, 2015). Object recognition is the same basic technology that allows computers to recognize faces (Taigman *et al.*, 2014), which can be used to automatically tag people in photo collections and allow computers to interact more naturally with their users.
- **Classification with missing inputs:** Classification becomes more challenging if the computer program is not guaranteed that every measurement in its input vector will always be provided. In order to solve the classification task, the learning algorithm only has to define a *single* function mapping from a vector input to a categorical output. When some of the inputs may be missing, rather than providing a single classification function, the learning algorithm must learn a *set* of functions. Each function corresponds to classifying  $\mathbf{x}$  with a different subset of its inputs missing. This kind of situation arises frequently in medical diagnosis, because many kinds of medical tests are expensive or invasive. One way to efficiently define such a large set of functions is to learn a probability distribution over all of the relevant variables, then solve the classification task by marginalizing out the missing variables. With  $n$  input variables, we can now obtain all  $2^n$  different classification functions needed for each possible set of missing inputs, but we only need to learn a single function describing the joint probability distribution. See Goodfellow *et al.* (2013b) for an example of a deep probabilistic model applied to such a task in this way. Many of the other tasks described in this section can also be generalized to work with missing inputs; classification with missing inputs is just one example of what machine learning can do.

- **Regression:** In this type of task, the computer program is asked to predict a numerical value given some input. To solve this task, the learning algorithm is asked to output a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . This type of task is similar to classification, except that the format of output is different. An example of a regression task is the prediction of the expected claim amount that an insured person will make (used to set insurance premiums), or the prediction of future prices of securities. These kinds of predictions are also used for algorithmic trading.
- **Transcription:** In this type of task, the machine learning system is asked to observe a relatively unstructured representation of some kind of data and transcribe it into discrete, textual form. For example, in optical character recognition, the computer program is shown a photograph containing an image of text and is asked to return this text in the form of a sequence of characters (e.g., in ASCII or Unicode format). Google Street View uses deep learning to process address numbers in this way (Goodfellow *et al.*, 2014d). Another example is speech recognition, where the computer program is provided an audio waveform and emits a sequence of characters or word ID codes describing the words that were spoken in the audio recording. Deep learning is a crucial component of modern speech recognition systems used at major companies including Microsoft, IBM and Google (Hinton *et al.*, 2012b).
- **Machine translation:** In a machine translation task, the input already consists of a sequence of symbols in some language, and the computer program must convert this into a sequence of symbols in another language. This is commonly applied to natural languages, such as translating from English to French. Deep learning has recently begun to have an important impact on this kind of task (Sutskever *et al.*, 2014; Bahdanau *et al.*, 2015).
- **Structured output:** Structured output tasks involve any task where the output is a vector (or other data structure containing multiple values) with important relationships between the different elements. This is a broad category, and subsumes the transcription and translation tasks described above, but also many other tasks. One example is parsing—mapping a natural language sentence into a tree that describes its grammatical structure and tagging nodes of the trees as being verbs, nouns, or adverbs, and so on. See Collobert (2011) for an example of deep learning applied to a parsing task. Another example is pixel-wise segmentation of images, where the computer program assigns every pixel in an image to a specific category. For

example, deep learning can be used to annotate the locations of roads in aerial photographs (Mnih and Hinton, 2010). The output need not have its form mirror the structure of the input as closely as in these annotation-style tasks. For example, in image captioning, the computer program observes an image and outputs a natural language sentence describing the image (Kiros *et al.*, 2014a,b; Mao *et al.*, 2015; Vinyals *et al.*, 2015b; Donahue *et al.*, 2014; Karpathy and Li, 2015; Fang *et al.*, 2015; Xu *et al.*, 2015). These tasks are called structured output tasks because the program must output several values that are all tightly inter-related. For example, the words produced by an image captioning program must form a valid sentence.

- **Anomaly detection:** In this type of task, the computer program sifts through a set of events or objects, and flags some of them as being unusual or atypical. An example of an anomaly detection task is credit card fraud detection. By modeling your purchasing habits, a credit card company can detect misuse of your cards. If a thief steals your credit card or credit card information, the thief's purchases will often come from a different probability distribution over purchase types than your own. The credit card company can prevent fraud by placing a hold on an account as soon as that card has been used for an uncharacteristic purchase. See Chandola *et al.* (2009) for a survey of anomaly detection methods.
- **Synthesis and sampling:** In this type of task, the machine learning algorithm is asked to generate new examples that are similar to those in the training data. Synthesis and sampling via machine learning can be useful for media applications where it can be expensive or boring for an artist to generate large volumes of content by hand. For example, video games can automatically generate textures for large objects or landscapes, rather than requiring an artist to manually label each pixel (Luo *et al.*, 2013). In some cases, we want the sampling or synthesis procedure to generate some specific kind of output given the input. For example, in a speech synthesis task, we provide a written sentence and ask the program to emit an audio waveform containing a spoken version of that sentence. This is a kind of structured output task, but with the added qualification that there is no single correct output for each input, and we explicitly desire a large amount of variation in the output, in order for the output to seem more natural and realistic.
- **Imputation of missing values:** In this type of task, the machine learning algorithm is given a new example  $\mathbf{x} \in \mathbb{R}^n$ , but with some entries  $x_i$  of  $\mathbf{x}$  missing. The algorithm must provide a prediction of the values of the missing entries.

- **Denoising:** In this type of task, the machine learning algorithm is given in input a *corrupted example*  $\tilde{\mathbf{x}} \in \mathbb{R}^n$  obtained by an unknown corruption process from a *clean example*  $\mathbf{x} \in \mathbb{R}^n$ . The learner must predict the clean example  $\mathbf{x}$  from its corrupted version  $\tilde{\mathbf{x}}$ , or more generally predict the conditional probability distribution  $p(\mathbf{x} \mid \tilde{\mathbf{x}})$ .
- **Density estimation or probability mass function estimation:** In the density estimation problem, the machine learning algorithm is asked to learn a function  $p_{\text{model}} : \mathbb{R}^n \rightarrow \mathbb{R}$ , where  $p_{\text{model}}(\mathbf{x})$  can be interpreted as a probability density function (if  $\mathbf{x}$  is continuous) or a probability mass function (if  $\mathbf{x}$  is discrete) on the space that the examples were drawn from. To do such a task well (we will specify exactly what that means when we discuss performance measures  $P$ ), the algorithm needs to learn the structure of the data it has seen. It must know where examples cluster tightly and where they are unlikely to occur. Most of the tasks described above require the learning algorithm to at least implicitly capture the structure of the probability distribution. Density estimation allows us to explicitly capture that distribution. In principle, we can then perform computations on that distribution in order to solve the other tasks as well. For example, if we have performed density estimation to obtain a probability distribution  $p(\mathbf{x})$ , we can use that distribution to solve the missing value imputation task. If a value  $x_i$  is missing and all of the other values, denoted  $\mathbf{x}_{-i}$ , are given, then we know the distribution over it is given by  $p(x_i \mid \mathbf{x}_{-i})$ . In practice, density estimation does not always allow us to solve all of these related tasks, because in many cases the required operations on  $p(\mathbf{x})$  are computationally intractable.

Of course, many other tasks and types of tasks are possible. The types of tasks we list here are intended only to provide examples of what machine learning can do, not to define a rigid taxonomy of tasks.

### 5.1.2 The Performance Measure, $P$

In order to evaluate the abilities of a machine learning algorithm, we must design a quantitative measure of its performance. Usually this performance measure  $P$  is specific to the task  $T$  being carried out by the system.

For tasks such as classification, classification with missing inputs, and transcription, we often measure the **accuracy** of the model. Accuracy is just the proportion of examples for which the model produces the correct output. We can

also obtain equivalent information by measuring the **error rate**, the proportion of examples for which the model produces an incorrect output. We often refer to the error rate as the expected 0-1 loss. The 0-1 loss on a particular example is 0 if it is correctly classified and 1 if it is not. For tasks such as density estimation, it does not make sense to measure accuracy, error rate, or any other kind of 0-1 loss. Instead, we must use a different performance metric that gives the model a continuous-valued score for each example. The most common approach is to report the average log-probability the model assigns to some examples.

Usually we are interested in how well the machine learning algorithm performs on data that it has not seen before, since this determines how well it will work when deployed in the real world. We therefore evaluate these performance measures using a **test set** of data that is separate from the data used for training the machine learning system.

The choice of performance measure may seem straightforward and objective, but it is often difficult to choose a performance measure that corresponds well to the desired behavior of the system.

In some cases, this is because it is difficult to decide what should be measured. For example, when performing a transcription task, should we measure the accuracy of the system at transcribing entire sequences, or should we use a more fine-grained performance measure that gives partial credit for getting some elements of the sequence correct? When performing a regression task, should we penalize the system more if it frequently makes medium-sized mistakes or if it rarely makes very large mistakes? These kinds of design choices depend on the application.

In other cases, we know what quantity we would ideally like to measure, but measuring it is impractical. For example, this arises frequently in the context of density estimation. Many of the best probabilistic models represent probability distributions only implicitly. Computing the actual probability value assigned to a specific point in space in many such models is intractable. In these cases, one must design an alternative criterion that still corresponds to the design objectives, or design a good approximation to the desired criterion.

### 5.1.3 The Experience, $E$

Machine learning algorithms can be broadly categorized as **unsupervised** or **supervised** by what kind of experience they are allowed to have during the learning process.

Most of the learning algorithms in this book can be understood as being allowed to experience an entire **dataset**. A dataset is a collection of many examples, as



defined in section 5.1.1. Sometimes we will also call examples **data points**.

One of the oldest datasets studied by statisticians and machine learning researchers is the Iris dataset (Fisher, 1936). It is a collection of measurements of different parts of 150 iris plants. Each individual plant corresponds to one example. The features within each example are the measurements of each of the parts of the plant: the sepal length, sepal width, petal length and petal width. The dataset also records which species each plant belonged to. Three different species are represented in the dataset.

**Unsupervised learning algorithms** experience a dataset containing many features, then learn useful properties of the structure of this dataset. In the context of deep learning, we usually want to learn the entire probability distribution that generated a dataset, whether explicitly as in density estimation or implicitly for tasks like synthesis or denoising. Some other unsupervised learning algorithms perform other roles, like clustering, which consists of dividing the dataset into clusters of similar examples.

**Supervised learning algorithms** experience a dataset containing features, but each example is also associated with a **label** or **target**. For example, the Iris dataset is annotated with the species of each iris plant. A supervised learning algorithm can study the Iris dataset and learn to classify iris plants into three different species based on their measurements.

Roughly speaking, unsupervised learning involves observing several examples of a random vector  $\mathbf{x}$ , and attempting to implicitly or explicitly learn the probability distribution  $p(\mathbf{x})$ , or some interesting properties of that distribution, while supervised learning involves observing several examples of a random vector  $\mathbf{x}$  and an associated value or vector  $\mathbf{y}$ , and learning to predict  $\mathbf{y}$  from  $\mathbf{x}$ , usually by estimating  $p(\mathbf{y} \mid \mathbf{x})$ . The term **supervised learning** originates from the view of the target  $\mathbf{y}$  being provided by an instructor or teacher who shows the machine learning system what to do. In unsupervised learning, there is no instructor or teacher, and the algorithm must learn to make sense of the data without this guide.

Unsupervised learning and supervised learning are not formally defined terms. The lines between them are often blurred. Many machine learning technologies can be used to perform both tasks. For example, the chain rule of probability states that for a vector  $\mathbf{x} \in \mathbb{R}^n$ , the joint distribution can be decomposed as

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i \mid x_1, \dots, x_{i-1}). \quad (5.1)$$

This decomposition means that we can solve the ostensibly unsupervised problem of modeling  $p(\mathbf{x})$  by splitting it into  $n$  supervised learning problems. Alternatively, we



can solve the supervised learning problem of learning  $p(y \mid \mathbf{x})$  by using traditional unsupervised learning technologies to learn the joint distribution  $p(\mathbf{x}, y)$  and inferring

$$p(y \mid \mathbf{x}) = \frac{p(\mathbf{x}, y)}{\sum_{y'} p(\mathbf{x}, y')}. \quad (5.2)$$

Though unsupervised learning and supervised learning are not completely formal or distinct concepts, they do help to roughly categorize some of the things we do with machine learning algorithms. Traditionally, people refer to regression, classification and structured output problems as supervised learning. Density estimation in support of other tasks is usually considered unsupervised learning.

Other variants of the learning paradigm are possible. For example, in semi-supervised learning, some examples include a supervision target but others do not. In multi-instance learning, an entire collection of examples is labeled as containing or not containing an example of a class, but the individual members of the collection are not labeled. For a recent example of multi-instance learning with deep models, see [Kotzias \*et al.\* \(2015\)](#).

Some machine learning algorithms do not just experience a fixed dataset. For example, **reinforcement learning** algorithms interact with an environment, so there is a feedback loop between the learning system and its experiences. Such algorithms are beyond the scope of this book. Please see [Sutton and Barto \(1998\)](#) or [Bertsekas and Tsitsiklis \(1996\)](#) for information about reinforcement learning, and [Mnih \*et al.\* \(2013\)](#) for the deep learning approach to reinforcement learning.

Most machine learning algorithms simply experience a dataset. A dataset can be described in many ways. In all cases, a dataset is a collection of examples, which are in turn collections of features.

One common way of describing a dataset is with a **design matrix**. A design matrix is a matrix containing a different example in each row. Each column of the matrix corresponds to a different feature. For instance, the Iris dataset contains 150 examples with four features for each example. This means we can represent the dataset with a design matrix  $\mathbf{X} \in \mathbb{R}^{150 \times 4}$ , where  $X_{i,1}$  is the sepal length of plant  $i$ ,  $X_{i,2}$  is the sepal width of plant  $i$ , etc. We will describe most of the learning algorithms in this book in terms of how they operate on design matrix datasets.

Of course, to describe a dataset as a design matrix, it must be possible to describe each example as a vector, and each of these vectors must be the same size. This is not always possible. For example, if you have a collection of photographs with different widths and heights, then different photographs will contain different numbers of pixels, so not all of the photographs may be described with the same length of vector. Section [9.7](#) and chapter [10](#) describe how to handle different

types of such heterogeneous data. In cases like these, rather than describing the dataset as a matrix with  $m$  rows, we will describe it as a set containing  $m$  elements:  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$ . This notation does not imply that any two example vectors  $\mathbf{x}^{(i)}$  and  $\mathbf{x}^{(j)}$  have the same size.

In the case of supervised learning, the example contains a label or target as well as a collection of features. For example, if we want to use a learning algorithm to perform object recognition from photographs, we need to specify which object appears in each of the photos. We might do this with a numeric code, with 0 signifying a person, 1 signifying a car, 2 signifying a cat, etc. Often when working with a dataset containing a design matrix of feature observations  $\mathbf{X}$ , we also provide a vector of labels  $\mathbf{y}$ , with  $y_i$  providing the label for example  $i$ .

Of course, sometimes the label may be more than just a single number. For example, if we want to train a speech recognition system to transcribe entire sentences, then the label for each example sentence is a sequence of words.

Just as there is no formal definition of supervised and unsupervised learning, there is no rigid taxonomy of datasets or experiences. The structures described here cover most cases, but it is always possible to design new ones for new applications.

### 5.1.4 Example: Linear Regression

Our definition of a machine learning algorithm as an algorithm that is capable of improving a computer program's performance at some task via experience is somewhat abstract. To make this more concrete, we present an example of a simple machine learning algorithm: **linear regression**. We will return to this example repeatedly as we introduce more machine learning concepts that help to understand its behavior.

As the name implies, linear regression solves a regression problem. In other words, the goal is to build a system that can take a vector  $\mathbf{x} \in \mathbb{R}^n$  as input and predict the value of a scalar  $y \in \mathbb{R}$  as its output. In the case of linear regression, the output is a linear function of the input. Let  $\hat{y}$  be the value that our model predicts  $y$  should take on. We define the output to be

$$\hat{y} = \mathbf{w}^\top \mathbf{x} \tag{5.3}$$

where  $\mathbf{w} \in \mathbb{R}^n$  is a vector of **parameters**.

Parameters are values that control the behavior of the system. In this case,  $w_i$  is the coefficient that we multiply by feature  $x_i$  before summing up the contributions from all the features. We can think of  $\mathbf{w}$  as a set of **weights** that determine how each feature affects the prediction. If a feature  $x_i$  receives a positive weight  $w_i$ ,

then increasing the value of that feature increases the value of our prediction  $\hat{y}$ . If a feature receives a negative weight, then increasing the value of that feature decreases the value of our prediction. If a feature's weight is large in magnitude, then it has a large effect on the prediction. If a feature's weight is zero, it has no effect on the prediction.

We thus have a definition of our task  $T$ : to predict  $y$  from  $\mathbf{x}$  by outputting  $\hat{y} = \mathbf{w}^\top \mathbf{x}$ . Next we need a definition of our performance measure,  $P$ .

Suppose that we have a design matrix of  $m$  example inputs that we will not use for training, only for evaluating how well the model performs. We also have a vector of regression targets providing the correct value of  $y$  for each of these examples. Because this dataset will only be used for evaluation, we call it the **test set**. We refer to the design matrix of inputs as  $\mathbf{X}^{(\text{test})}$  and the vector of regression targets as  $\mathbf{y}^{(\text{test})}$ .

One way of measuring the performance of the model is to compute the **mean squared error** of the model on the test set. If  $\hat{\mathbf{y}}^{(\text{test})}$  gives the predictions of the model on the test set, then the mean squared error is given by

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})})_i^2. \quad (5.4)$$

Intuitively, one can see that this error measure decreases to 0 when  $\hat{\mathbf{y}}^{(\text{test})} = \mathbf{y}^{(\text{test})}$ . We can also see that

$$\text{MSE}_{\text{test}} = \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})}\|_2^2, \quad (5.5)$$

so the error increases whenever the Euclidean distance between the predictions and the targets increases.

To make a machine learning algorithm, we need to design an algorithm that will improve the weights  $\mathbf{w}$  in a way that reduces  $\text{MSE}_{\text{test}}$  when the algorithm is allowed to gain experience by observing a training set  $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$ . One intuitive way of doing this (which we will justify later, in section 5.5.1) is just to minimize the mean squared error on the training set,  $\text{MSE}_{\text{train}}$ .

To minimize  $\text{MSE}_{\text{train}}$ , we can simply solve for where its gradient is  $\mathbf{0}$ :

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = \mathbf{0} \quad (5.6)$$

$$\Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 = \mathbf{0} \quad (5.7)$$

$$\Rightarrow \frac{1}{m} \nabla_{\mathbf{w}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = \mathbf{0} \quad (5.8)$$

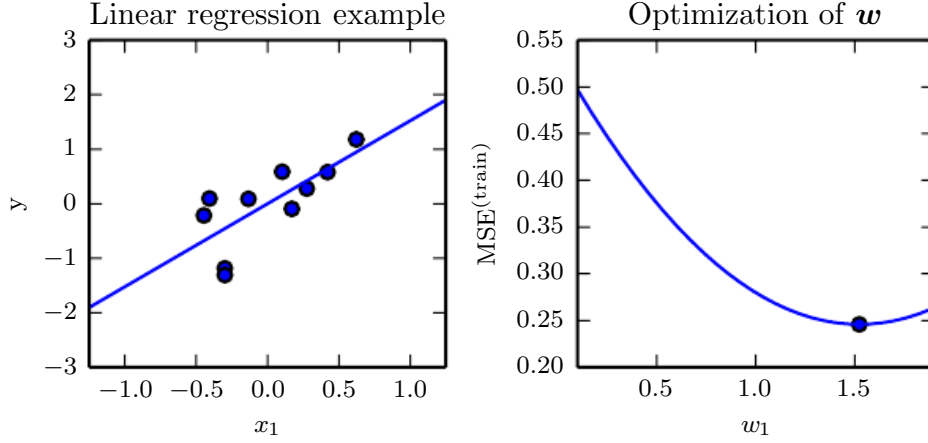


Figure 5.1: A linear regression problem, with a training set consisting of ten data points, each containing one feature. Because there is only one feature, the weight vector  $\mathbf{w}$  contains only a single parameter to learn,  $w_1$ . (Left) Observe that linear regression learns to set  $w_1$  such that the line  $y = w_1 x$  comes as close as possible to passing through all the training points. (Right) The plotted point indicates the value of  $w_1$  found by the normal equations, which we can see minimizes the mean squared error on the training set.

$$\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0 \quad (5.9)$$

$$\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0 \quad (5.10)$$

$$\Rightarrow 2\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0 \quad (5.11)$$

$$\Rightarrow \mathbf{w} = \left( \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \quad (5.12)$$

The system of equations whose solution is given by equation 5.12 is known as the **normal equations**. Evaluating equation 5.12 constitutes a simple learning algorithm. For an example of the linear regression learning algorithm in action, see figure 5.1.

It is worth noting that the term **linear regression** is often used to refer to a slightly more sophisticated model with one additional parameter—an intercept term  $b$ . In this model

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b \quad (5.13)$$

so the mapping from parameters to predictions is still a linear function but the mapping from features to predictions is now an affine function. This extension to affine functions means that the plot of the model's predictions still looks like a line, but it need not pass through the origin. Instead of adding the bias parameter

$b$ , one can continue to use the model with only weights but augment  $\mathbf{x}$  with an extra entry that is always set to 1. The weight corresponding to the extra 1 entry plays the role of the bias parameter. We will frequently use the term “linear” when referring to affine functions throughout this book.

The intercept term  $b$  is often called the **bias** parameter of the affine transformation. This terminology derives from the point of view that the output of the transformation is biased toward being  $b$  in the absence of any input. This term is different from the idea of a statistical bias, in which a statistical estimation algorithm’s expected estimate of a quantity is not equal to the true quantity.

Linear regression is of course an extremely simple and limited learning algorithm, but it provides an example of how a learning algorithm can work. In the subsequent sections we will describe some of the basic principles underlying learning algorithm design and demonstrate how these principles can be used to build more complicated learning algorithms.

## 5.2 Capacity, Overfitting and Underfitting

The central challenge in machine learning is that we must perform well on *new, previously unseen* inputs—not just those on which our model was trained. The ability to perform well on previously unobserved inputs is called **generalization**.

Typically, when training a machine learning model, we have access to a training set, we can compute some error measure on the training set called the **training error**, and we reduce this training error. So far, what we have described is simply an optimization problem. What separates machine learning from optimization is that we want the **generalization error**, also called the **test error**, to be low as well. The generalization error is defined as the expected value of the error on a new input. Here the expectation is taken across different possible inputs, drawn from the distribution of inputs we expect the system to encounter in practice.

We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.

In our linear regression example, we trained the model by minimizing the training error,

$$\frac{1}{m^{(\text{train})}} \|\mathbf{X}^{(\text{train})}\mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2, \quad (5.14)$$

but we actually care about the test error,  $\frac{1}{m^{(\text{test})}} \|\mathbf{X}^{(\text{test})}\mathbf{w} - \mathbf{y}^{(\text{test})}\|_2^2$ .

How can we affect performance on the test set when we get to observe only the

training set? The field of **statistical learning theory** provides some answers. If the training and the test set are collected arbitrarily, there is indeed little we can do. If we are allowed to make some assumptions about how the training and test set are collected, then we can make some progress.

The train and test data are generated by a probability distribution over datasets called the **data generating process**. We typically make a set of assumptions known collectively as the **i.i.d. assumptions**. These assumptions are that the examples in each dataset are **independent** from each other, and that the train set and test set are **identically distributed**, drawn from the same probability distribution as each other. This assumption allows us to describe the data generating process with a probability distribution over a single example. The same distribution is then used to generate every train example and every test example. We call that shared underlying distribution the **data generating distribution**, denoted  $p_{\text{data}}$ . This probabilistic framework and the i.i.d. assumptions allow us to mathematically study the relationship between training error and test error.

One immediate connection we can observe between the training and test error is that the expected training error of a randomly selected model is equal to the expected test error of that model. Suppose we have a probability distribution  $p(\mathbf{x}, y)$  and we sample from it repeatedly to generate the train set and the test set. For some fixed value  $\mathbf{w}$ , the expected training set error is exactly the same as the expected test set error, because both expectations are formed using the same dataset sampling process. The only difference between the two conditions is the name we assign to the dataset we sample.

Of course, when we use a machine learning algorithm, we do not fix the parameters ahead of time, then sample both datasets. We sample the training set, then use it to choose the parameters to reduce training set error, then sample the test set. Under this process, the expected test error is greater than or equal to the expected value of training error. The factors determining how well a machine learning algorithm will perform are its ability to:

1. Make the training error small.
2. Make the gap between training and test error small.

These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

We can control whether a model is more likely to overfit or underfit by altering its **capacity**. Informally, a model's capacity is its ability to fit a wide variety of

functions. Models with low capacity may struggle to fit the training set. Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.

One way to control the capacity of a learning algorithm is by choosing its **hypothesis space**, the set of functions that the learning algorithm is allowed to select as being the solution. For example, the linear regression algorithm has the set of all linear functions of its input as its hypothesis space. We can generalize linear regression to include polynomials, rather than just linear functions, in its hypothesis space. Doing so increases the model's capacity.

A polynomial of degree one gives us the linear regression model with which we are already familiar, with prediction

$$\hat{y} = b + wx. \quad (5.15)$$

By introducing  $x^2$  as another feature provided to the linear regression model, we can learn a model that is quadratic as a function of  $x$ :

$$\hat{y} = b + w_1x + w_2x^2. \quad (5.16)$$

Though this model implements a quadratic function of its *input*, the output is still a linear function of the *parameters*, so we can still use the normal equations to train the model in closed form. We can continue to add more powers of  $x$  as additional features, for example to obtain a polynomial of degree 9:

$$\hat{y} = b + \sum_{i=1}^9 w_i x^i. \quad (5.17)$$

Machine learning algorithms will generally perform best when their capacity is appropriate for the true complexity of the task they need to perform and the amount of training data they are provided with. Models with insufficient capacity are unable to solve complex tasks. Models with high capacity can solve complex tasks, but when their capacity is higher than needed to solve the present task they may overfit.

Figure 5.2 shows this principle in action. We compare a linear, quadratic and degree-9 predictor attempting to fit a problem where the true underlying function is quadratic. The linear function is unable to capture the curvature in the true underlying problem, so it underfits. The degree-9 predictor is capable of representing the correct function, but it is also capable of representing infinitely many other functions that pass exactly through the training points, because we

have more parameters than training examples. We have little chance of choosing a solution that generalizes well when so many wildly different solutions exist. In this example, the quadratic model is perfectly matched to the true structure of the task so it generalizes well to new data.

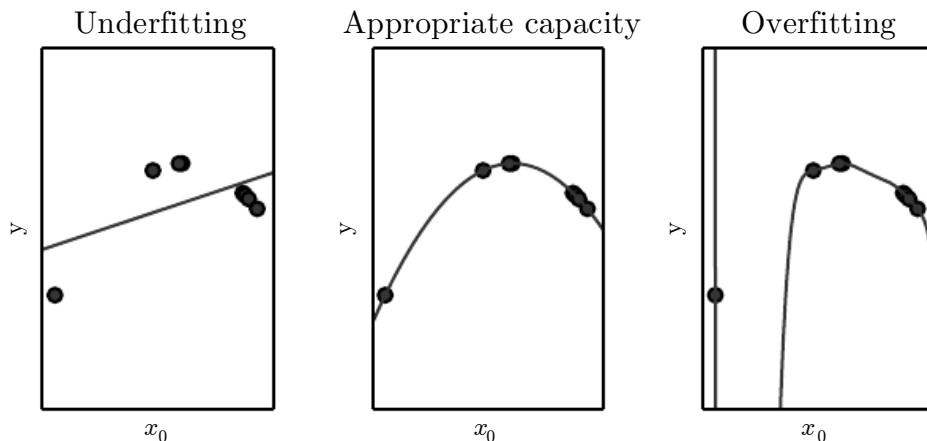


Figure 5.2: We fit three models to this example training set. The training data was generated synthetically, by randomly sampling  $x$  values and choosing  $y$  deterministically by evaluating a quadratic function. (Left) A linear function fit to the data suffers from underfitting—it cannot capture the curvature that is present in the data. (Center) A quadratic function fit to the data generalizes well to unseen points. It does not suffer from a significant amount of overfitting or underfitting. (Right) A polynomial of degree 9 fit to the data suffers from overfitting. Here we used the Moore-Penrose pseudoinverse to solve the underdetermined normal equations. The solution passes through all of the training points exactly, but we have not been lucky enough for it to extract the correct structure. It now has a deep valley in between two training points that does not appear in the true underlying function. It also increases sharply on the left side of the data, while the true function decreases in this area.

So far we have described only one way of changing a model’s capacity: by changing the number of input features it has, and simultaneously adding new parameters associated with those features. There are in fact many ways of changing a model’s capacity. Capacity is not determined only by the choice of model. The model specifies which family of functions the learning algorithm can choose from when varying the parameters in order to reduce a training objective. This is called the **representational capacity** of the model. In many cases, finding the best function within this family is a very difficult optimization problem. In practice, the learning algorithm does not actually find the best function, but merely one that significantly reduces the training error. These additional limitations, such as



the imperfection of the optimization algorithm, mean that the learning algorithm’s **effective capacity** may be less than the representational capacity of the model family.

Our modern ideas about improving the generalization of machine learning models are refinements of thought dating back to philosophers at least as early as Ptolemy. Many early scholars invoke a principle of parsimony that is now most widely known as **Occam’s razor** (c. 1287-1347). This principle states that among competing hypotheses that explain known observations equally well, one should choose the “simplest” one. This idea was formalized and made more precise in the 20th century by the founders of statistical learning theory (Vapnik and Chervonenkis, 1971; Vapnik, 1982; Blumer *et al.*, 1989; Vapnik, 1995).

Statistical learning theory provides various means of quantifying model capacity. Among these, the most well-known is the **Vapnik-Chervonenkis dimension**, or VC dimension. The VC dimension measures the capacity of a binary classifier. The VC dimension is defined as being the largest possible value of  $m$  for which there exists a training set of  $m$  different  $\mathbf{x}$  points that the classifier can label arbitrarily.

Quantifying the capacity of the model allows statistical learning theory to make quantitative predictions. The most important results in statistical learning theory show that the discrepancy between training error and generalization error is bounded from above by a quantity that grows as the model capacity grows but shrinks as the number of training examples increases (Vapnik and Chervonenkis, 1971; Vapnik, 1982; Blumer *et al.*, 1989; Vapnik, 1995). These bounds provide intellectual justification that machine learning algorithms can work, but they are rarely used in practice when working with deep learning algorithms. This is in part because the bounds are often quite loose and in part because it can be quite difficult to determine the capacity of deep learning algorithms. The problem of determining the capacity of a deep learning model is especially difficult because the effective capacity is limited by the capabilities of the optimization algorithm, and we have little theoretical understanding of the very general non-convex optimization problems involved in deep learning.

We must remember that while simpler functions are more likely to generalize (to have a small gap between training and test error) we must still choose a sufficiently complex hypothesis to achieve low training error. Typically, training error decreases until it asymptotes to the minimum possible error value as model capacity increases (assuming the error measure has a minimum value). Typically, generalization error has a U-shaped curve as a function of model capacity. This is illustrated in figure 5.3.

To reach the most extreme case of arbitrarily high capacity, we introduce

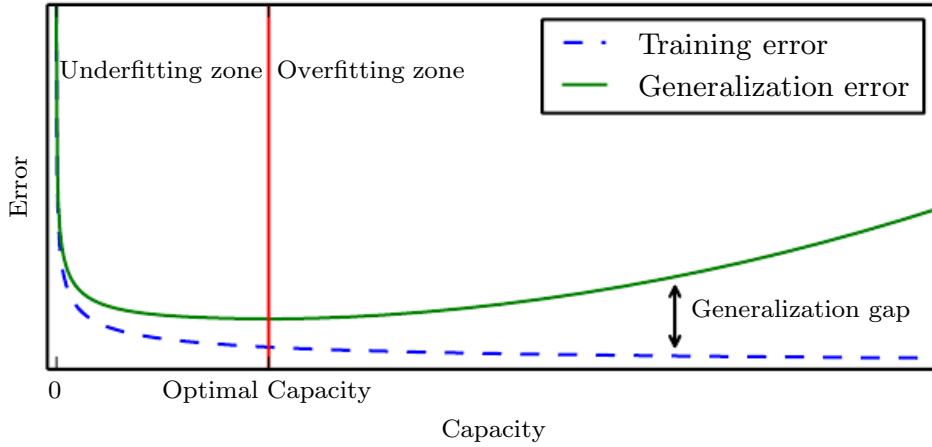


Figure 5.3: Typical relationship between capacity and error. Training and test error behave differently. At the left end of the graph, training error and generalization error are both high. This is the **underfitting regime**. As we increase capacity, training error decreases, but the gap between training and generalization error increases. Eventually, the size of this gap outweighs the decrease in training error, and we enter the **overfitting regime**, where capacity is too large, above the **optimal capacity**.

the concept of **non-parametric** models. So far, we have seen only parametric models, such as linear regression. Parametric models learn a function described by a parameter vector whose size is finite and fixed before any data is observed. Non-parametric models have no such limitation.

Sometimes, non-parametric models are just theoretical abstractions (such as an algorithm that searches over all possible probability distributions) that cannot be implemented in practice. However, we can also design practical non-parametric models by making their complexity a function of the training set size. One example of such an algorithm is **nearest neighbor regression**. Unlike linear regression, which has a fixed-length vector of weights, the nearest neighbor regression model simply stores the  $\mathbf{X}$  and  $\mathbf{y}$  from the training set. When asked to classify a test point  $\mathbf{x}$ , the model looks up the nearest entry in the training set and returns the associated regression target. In other words,  $\hat{y} = y_i$  where  $i = \arg \min \|\mathbf{X}_{i,:} - \mathbf{x}\|_2^2$ . The algorithm can also be generalized to distance metrics other than the  $L^2$  norm, such as learned distance metrics (Goldberger *et al.*, 2005). If the algorithm is allowed to break ties by averaging the  $y_i$  values for all  $\mathbf{X}_{i,:}$  that are tied for nearest, then this algorithm is able to achieve the minimum possible training error (which might be greater than zero, if two identical inputs are associated with different outputs) on any regression dataset.

Finally, we can also create a non-parametric learning algorithm by wrapping a

parametric learning algorithm inside another algorithm that increases the number of parameters as needed. For example, we could imagine an outer loop of learning that changes the degree of the polynomial learned by linear regression on top of a polynomial expansion of the input.

The ideal model is an oracle that simply knows the true probability distribution that generates the data. Even such a model will still incur some error on many problems, because there may still be some noise in the distribution. In the case of supervised learning, the mapping from  $\mathbf{x}$  to  $y$  may be inherently stochastic, or  $y$  may be a deterministic function that involves other variables besides those included in  $\mathbf{x}$ . The error incurred by an oracle making predictions from the true distribution  $p(\mathbf{x}, y)$  is called the **Bayes error**.

Training and generalization error vary as the size of the training set varies. Expected generalization error can never increase as the number of training examples increases. For non-parametric models, more data yields better generalization until the best possible error is achieved. Any fixed parametric model with less than optimal capacity will asymptote to an error value that exceeds the Bayes error. See figure 5.4 for an illustration. Note that it is possible for the model to have optimal capacity and yet still have a large gap between training and generalization error. In this situation, we may be able to reduce this gap by gathering more training examples.

### 5.2.1 The No Free Lunch Theorem

Learning theory claims that a machine learning algorithm can generalize well from a finite training set of examples. This seems to contradict some basic principles of logic. Inductive reasoning, or inferring general rules from a limited set of examples, is not logically valid. To logically infer a rule describing every member of a set, one must have information about every member of that set.

In part, machine learning avoids this problem by offering only probabilistic rules, rather than the entirely certain rules used in purely logical reasoning. Machine learning promises to find rules that are *probably* correct about *most* members of the set they concern.

Unfortunately, even this does not resolve the entire problem. The **no free lunch theorem** for machine learning (Wolpert, 1996) states that, averaged over all possible data generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points. In other words, in some sense, no machine learning algorithm is universally any better than any other. The most sophisticated algorithm we can conceive of has the same average

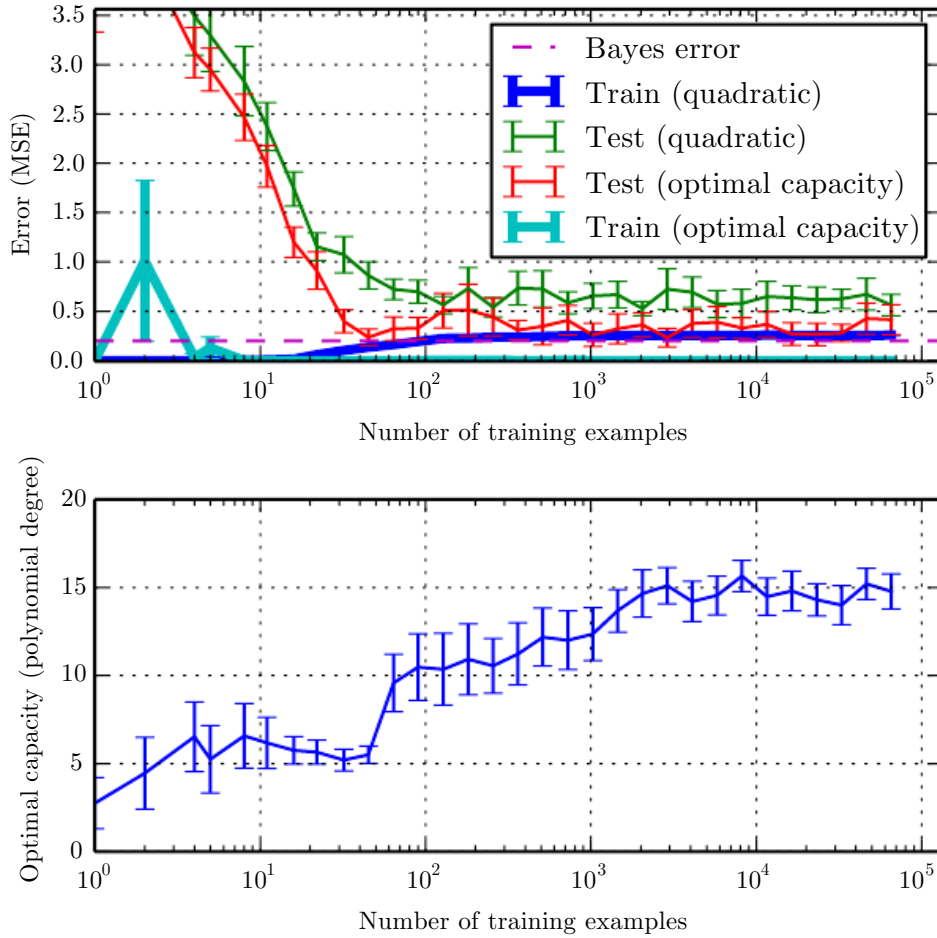


Figure 5.4: The effect of the training dataset size on the train and test error, as well as on the optimal model capacity. We constructed a synthetic regression problem based on adding a moderate amount of noise to a degree-5 polynomial, generated a single test set, and then generated several different sizes of training set. For each size, we generated 40 different training sets in order to plot error bars showing 95 percent confidence intervals. (*Top*) The MSE on the training and test set for two different models: a quadratic model, and a model with degree chosen to minimize the test error. Both are fit in closed form. For the quadratic model, the training error increases as the size of the training set increases. This is because larger datasets are harder to fit. Simultaneously, the test error decreases, because fewer incorrect hypotheses are consistent with the training data. The quadratic model does not have enough capacity to solve the task, so its test error asymptotes to a high value. The test error at optimal capacity asymptotes to the Bayes error. The training error can fall below the Bayes error, due to the ability of the training algorithm to memorize specific instances of the training set. As the training size increases to infinity, the training error of any fixed-capacity model (here, the quadratic model) must rise to at least the Bayes error. (*Bottom*) As the training set size increases, the optimal capacity (shown here as the degree of the optimal polynomial regressor) increases. The optimal capacity plateaus after reaching sufficient complexity to solve the task.

performance (over all possible tasks) as merely predicting that every point belongs to the same class.

Fortunately, these results hold only when we average over *all* possible data generating distributions. If we make assumptions about the kinds of probability distributions we encounter in real-world applications, then we can design learning algorithms that perform well on these distributions.

This means that the goal of machine learning research is not to seek a universal learning algorithm or the absolute best learning algorithm. Instead, our goal is to understand what kinds of distributions are relevant to the “real world” that an AI agent experiences, and what kinds of machine learning algorithms perform well on data drawn from the kinds of data generating distributions we care about.

### 5.2.2 Regularization

The no free lunch theorem implies that we must design our machine learning algorithms to perform well on a specific task. We do so by building a set of preferences into the learning algorithm. When these preferences are aligned with the learning problems we ask the algorithm to solve, it performs better.

So far, the only method of modifying a learning algorithm that we have discussed concretely is to increase or decrease the model’s representational capacity by adding or removing functions from the hypothesis space of solutions the learning algorithm is able to choose. We gave the specific example of increasing or decreasing the degree of a polynomial for a regression problem. The view we have described so far is oversimplified.

The behavior of our algorithm is strongly affected not just by how large we make the set of functions allowed in its hypothesis space, but by the specific identity of those functions. The learning algorithm we have studied so far, linear regression, has a hypothesis space consisting of the set of linear functions of its input. These linear functions can be very useful for problems where the relationship between inputs and outputs truly is close to linear. They are less useful for problems that behave in a very nonlinear fashion. For example, linear regression would not perform very well if we tried to use it to predict  $\sin(x)$  from  $x$ . We can thus control the performance of our algorithms by choosing what kind of functions we allow them to draw solutions from, as well as by controlling the amount of these functions.

We can also give a learning algorithm a preference for one solution in its hypothesis space to another. This means that both functions are eligible, but one is preferred. The unpreferred solution will be chosen only if it fits the training

data significantly better than the preferred solution.

For example, we can modify the training criterion for linear regression to include **weight decay**. To perform linear regression with weight decay, we minimize a sum comprising both the mean squared error on the training and a criterion  $J(\mathbf{w})$  that expresses a preference for the weights to have smaller squared  $L^2$  norm. Specifically,

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^\top \mathbf{w}, \quad (5.18)$$

where  $\lambda$  is a value chosen ahead of time that controls the strength of our preference for smaller weights. When  $\lambda = 0$ , we impose no preference, and larger  $\lambda$  forces the weights to become smaller. Minimizing  $J(\mathbf{w})$  results in a choice of weights that make a tradeoff between fitting the training data and being small. This gives us solutions that have a smaller slope, or put weight on fewer of the features. As an example of how we can control a model's tendency to overfit or underfit via weight decay, we can train a high-degree polynomial regression model with different values of  $\lambda$ . See figure 5.5 for the results.

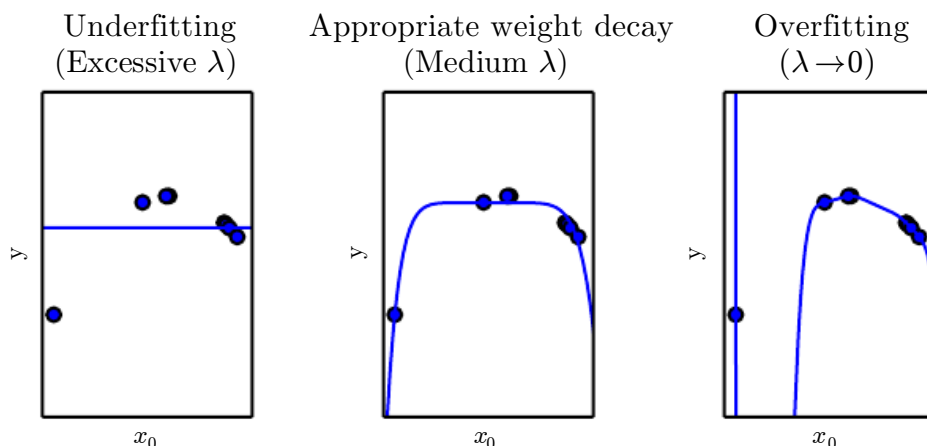


Figure 5.5: We fit a high-degree polynomial regression model to our example training set from figure 5.2. The true function is quadratic, but here we use only models with degree 9. We vary the amount of weight decay to prevent these high-degree models from overfitting. (Left) With very large  $\lambda$ , we can force the model to learn a function with no slope at all. This underfits because it can only represent a constant function. (Center) With a medium value of  $\lambda$ , the learning algorithm recovers a curve with the right general shape. Even though the model is capable of representing functions with much more complicated shape, weight decay has encouraged it to use a simpler function described by smaller coefficients. (Right) With weight decay approaching zero (i.e., using the Moore-Penrose pseudoinverse to solve the underdetermined problem with minimal regularization), the degree-9 polynomial overfits significantly, as we saw in figure 5.2.

More generally, we can regularize a model that learns a function  $f(\mathbf{x}; \boldsymbol{\theta})$  by adding a penalty called a **regularizer** to the cost function. In the case of weight decay, the regularizer is  $\Omega(\mathbf{w}) = \mathbf{w}^\top \mathbf{w}$ . In chapter 7, we will see that many other regularizers are possible.

Expressing preferences for one function over another is a more general way of controlling a model's capacity than including or excluding members from the hypothesis space. We can think of excluding a function from a hypothesis space as expressing an infinitely strong preference against that function.

In our weight decay example, we expressed our preference for linear functions defined with smaller weights explicitly, via an extra term in the criterion we minimize. There are many other ways of expressing preferences for different solutions, both implicitly and explicitly. Together, these different approaches are known as **regularization**. *Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.* Regularization is one of the central concerns of the field of machine learning, rivaled in its importance only by optimization.

The no free lunch theorem has made it clear that there is no best machine learning algorithm, and, in particular, no best form of regularization. Instead we must choose a form of regularization that is well-suited to the particular task we want to solve. The philosophy of deep learning in general and this book in particular is that a very wide range of tasks (such as all of the intellectual tasks that people can do) may all be solved effectively using very general-purpose forms of regularization.

## 5.3 Hyperparameters and Validation Sets

Most machine learning algorithms have several settings that we can use to control the behavior of the learning algorithm. These settings are called **hyperparameters**. The values of hyperparameters are not adapted by the learning algorithm itself (though we can design a nested learning procedure where one learning algorithm learns the best hyperparameters for another learning algorithm).

In the polynomial regression example we saw in figure 5.2, there is a single hyperparameter: the degree of the polynomial, which acts as a **capacity** hyperparameter. The  $\lambda$  value used to control the strength of weight decay is another example of a hyperparameter.

Sometimes a setting is chosen to be a hyperparameter that the learning algorithm does not learn because it is difficult to optimize. More frequently, the



setting must be a hyperparameter because it is not appropriate to learn that hyperparameter on the training set. This applies to all hyperparameters that control model capacity. If learned on the training set, such hyperparameters would always choose the maximum possible model capacity, resulting in overfitting (refer to figure 5.3). For example, we can always fit the training set better with a higher degree polynomial and a weight decay setting of  $\lambda = 0$  than we could with a lower degree polynomial and a positive weight decay setting.

To solve this problem, we need a **validation set** of examples that the training algorithm does not observe.

Earlier we discussed how a held-out test set, composed of examples coming from the same distribution as the training set, can be used to estimate the generalization error of a learner, after the learning process has completed. It is important that the test examples are not used in any way to make choices about the model, including its hyperparameters. For this reason, no example from the test set can be used in the validation set. Therefore, we always construct the validation set from the *training* data. Specifically, we split the training data into two disjoint subsets. One of these subsets is used to learn the parameters. The other subset is our validation set, used to estimate the generalization error during or after training, allowing for the hyperparameters to be updated accordingly. The subset of data used to learn the parameters is still typically called the training set, even though this may be confused with the larger pool of data used for the entire training process. The subset of data used to guide the selection of hyperparameters is called the validation set. Typically, one uses about 80% of the training data for training and 20% for validation. Since the validation set is used to “train” the hyperparameters, the validation set error will underestimate the generalization error, though typically by a smaller amount than the training error. After all hyperparameter optimization is complete, the generalization error may be estimated using the test set.

In practice, when the same test set has been used repeatedly to evaluate performance of different algorithms over many years, and especially if we consider all the attempts from the scientific community at beating the reported state-of-the-art performance on that test set, we end up having optimistic evaluations with the test set as well. Benchmarks can thus become stale and then do not reflect the true field performance of a trained system. Thankfully, the community tends to move on to new (and usually more ambitious and larger) benchmark datasets.



### 5.3.1 Cross-Validation

Dividing the dataset into a fixed training set and a fixed test set can be problematic if it results in the test set being small. A small test set implies statistical uncertainty around the estimated average test error, making it difficult to claim that algorithm  $A$  works better than algorithm  $B$  on the given task.

When the dataset has hundreds of thousands of examples or more, this is not a serious issue. When the dataset is too small, alternative procedures enable one to use all of the examples in the estimation of the mean test error, at the price of increased computational cost. These procedures are based on the idea of repeating the training and testing computation on different randomly chosen subsets or splits of the original dataset. The most common of these is the  $k$ -fold cross-validation procedure, shown in algorithm 5.1, in which a partition of the dataset is formed by splitting it into  $k$  non-overlapping subsets. The test error may then be estimated by taking the average test error across  $k$  trials. On trial  $i$ , the  $i$ -th subset of the data is used as the test set and the rest of the data is used as the training set. One problem is that there exist no unbiased estimators of the variance of such average error estimators (Bengio and Grandvalet, 2004), but approximations are typically used.

## 5.4 Estimators, Bias and Variance

The field of statistics gives us many tools that can be used to achieve the machine learning goal of solving a task not only on the training set but also to generalize. Foundational concepts such as parameter estimation, bias and variance are useful to formally characterize notions of generalization, underfitting and overfitting.

### 5.4.1 Point Estimation

Point estimation is the attempt to provide the single “best” prediction of some quantity of interest. In general the quantity of interest can be a single parameter or a vector of parameters in some parametric model, such as the weights in our linear regression example in section 5.1.4, but it can also be a whole function.

In order to distinguish estimates of parameters from their true value, our convention will be to denote a point estimate of a parameter  $\theta$  by  $\hat{\theta}$ .

Let  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  be a set of  $m$  independent and identically distributed

---

**Algorithm 5.1** The  $k$ -fold cross-validation algorithm. It can be used to estimate generalization error of a learning algorithm  $A$  when the given dataset  $\mathbb{D}$  is too small for a simple train/test or train/valid split to yield accurate estimation of generalization error, because the mean of a loss  $L$  on a small test set may have too high variance. The dataset  $\mathbb{D}$  contains as elements the abstract examples  $\mathbf{z}^{(i)}$  (for the  $i$ -th example), which could stand for an (input,target) pair  $\mathbf{z}^{(i)} = (\mathbf{x}^{(i)}, y^{(i)})$  in the case of supervised learning, or for just an input  $\mathbf{z}^{(i)} = \mathbf{x}^{(i)}$  in the case of unsupervised learning. The algorithm returns the vector of errors  $\mathbf{e}$  for each example in  $\mathbb{D}$ , whose mean is the estimated generalization error. The errors on individual examples can be used to compute a confidence interval around the mean (equation 5.47). While these confidence intervals are not well-justified after the use of cross-validation, it is still common practice to use them to declare that algorithm  $A$  is better than algorithm  $B$  only if the confidence interval of the error of algorithm  $A$  lies below and does not intersect the confidence interval of algorithm  $B$ .

---

**Define**  $\text{KFoldXV}(\mathbb{D}, A, L, k)$ :

**Require:**  $\mathbb{D}$ , the given dataset, with elements  $\mathbf{z}^{(i)}$

**Require:**  $A$ , the learning algorithm, seen as a function that takes a dataset as input and outputs a learned function

**Require:**  $L$ , the loss function, seen as a function from a learned function  $f$  and an example  $\mathbf{z}^{(i)} \in \mathbb{D}$  to a scalar  $\in \mathbb{R}$

**Require:**  $k$ , the number of folds

Split  $\mathbb{D}$  into  $k$  mutually exclusive subsets  $\mathbb{D}_i$ , whose union is  $\mathbb{D}$ .

**for**  $i$  from 1 to  $k$  **do**

$f_i = A(\mathbb{D} \setminus \mathbb{D}_i)$

**for**  $\mathbf{z}^{(j)}$  in  $\mathbb{D}_i$  **do**

$e_j = L(f_i, \mathbf{z}^{(j)})$

**end for**

**end for**

**Return**  $\mathbf{e}$

---

(i.i.d.) data points. A **point estimator** or **statistic** is any function of the data:

$$\hat{\boldsymbol{\theta}}_m = g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}). \quad (5.19)$$

The definition does not require that  $g$  return a value that is close to the true  $\boldsymbol{\theta}$  or even that the range of  $g$  is the same as the set of allowable values of  $\boldsymbol{\theta}$ . This definition of a point estimator is very general and allows the designer of an estimator great flexibility. While almost any function thus qualifies as an estimator, a good estimator is a function whose output is close to the true underlying  $\boldsymbol{\theta}$  that generated the training data.

For now, we take the frequentist perspective on statistics. That is, we assume that the true parameter value  $\boldsymbol{\theta}$  is fixed but unknown, while the point estimate  $\hat{\boldsymbol{\theta}}$  is a function of the data. Since the data is drawn from a random process, any function of the data is random. Therefore  $\hat{\boldsymbol{\theta}}$  is a random variable.

Point estimation can also refer to the estimation of the relationship between input and target variables. We refer to these types of point estimates as function estimators.

**Function Estimation** As we mentioned above, sometimes we are interested in performing function estimation (or function approximation). Here we are trying to predict a variable  $\mathbf{y}$  given an input vector  $\mathbf{x}$ . We assume that there is a function  $f(\mathbf{x})$  that describes the approximate relationship between  $\mathbf{y}$  and  $\mathbf{x}$ . For example, we may assume that  $\mathbf{y} = f(\mathbf{x}) + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon}$  stands for the part of  $\mathbf{y}$  that is not predictable from  $\mathbf{x}$ . In function estimation, we are interested in approximating  $f$  with a model or estimate  $\hat{f}$ . Function estimation is really just the same as estimating a parameter  $\boldsymbol{\theta}$ ; the function estimator  $\hat{f}$  is simply a point estimator in function space. The linear regression example (discussed above in section 5.1.4) and the polynomial regression example (discussed in section 5.2) are both examples of scenarios that may be interpreted either as estimating a parameter  $\mathbf{w}$  or estimating a function  $\hat{f}$  mapping from  $\mathbf{x}$  to  $y$ .

We now review the most commonly studied properties of point estimators and discuss what they tell us about these estimators.

### 5.4.2 Bias

The bias of an estimator is defined as:

$$\text{bias}(\hat{\boldsymbol{\theta}}_m) = \mathbb{E}(\hat{\boldsymbol{\theta}}_m) - \boldsymbol{\theta} \quad (5.20)$$

where the expectation is over the data (seen as samples from a random variable) and  $\theta$  is the true underlying value of  $\theta$  used to define the data generating distribution. An estimator  $\hat{\theta}_m$  is said to be **unbiased** if  $\text{bias}(\hat{\theta}_m) = \mathbf{0}$ , which implies that  $\mathbb{E}(\hat{\theta}_m) = \theta$ . An estimator  $\hat{\theta}_m$  is said to be **asymptotically unbiased** if  $\lim_{m \rightarrow \infty} \text{bias}(\hat{\theta}_m) = \mathbf{0}$ , which implies that  $\lim_{m \rightarrow \infty} \mathbb{E}(\hat{\theta}_m) = \theta$ .

**Example: Bernoulli Distribution** Consider a set of samples  $\{x^{(1)}, \dots, x^{(m)}\}$  that are independently and identically distributed according to a Bernoulli distribution with mean  $\theta$ :

$$P(x^{(i)}; \theta) = \theta^{x^{(i)}} (1 - \theta)^{(1-x^{(i)})}. \quad (5.21)$$

A common estimator for the  $\theta$  parameter of this distribution is the mean of the training samples:

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}. \quad (5.22)$$

To determine whether this estimator is biased, we can substitute equation 5.22 into equation 5.20:

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}[\hat{\theta}_m] - \theta \quad (5.23)$$

$$= \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m x^{(i)} \right] - \theta \quad (5.24)$$

$$= \frac{1}{m} \sum_{i=1}^m \mathbb{E} [x^{(i)}] - \theta \quad (5.25)$$

$$= \frac{1}{m} \sum_{i=1}^m \sum_{x^{(i)}=0}^1 \left( x^{(i)} \theta^{x^{(i)}} (1 - \theta)^{(1-x^{(i)})} \right) - \theta \quad (5.26)$$

$$= \frac{1}{m} \sum_{i=1}^m (\theta) - \theta \quad (5.27)$$

$$= \theta - \theta = 0 \quad (5.28)$$

Since  $\text{bias}(\hat{\theta}) = 0$ , we say that our estimator  $\hat{\theta}$  is unbiased.

**Example: Gaussian Distribution Estimator of the Mean** Now, consider a set of samples  $\{x^{(1)}, \dots, x^{(m)}\}$  that are independently and identically distributed according to a Gaussian distribution  $p(x^{(i)}) = \mathcal{N}(x^{(i)}; \mu, \sigma^2)$ , where  $i \in \{1, \dots, m\}$ .

Recall that the Gaussian probability density function is given by

$$p(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}\right). \quad (5.29)$$

A common estimator of the Gaussian mean parameter is known as the **sample mean**:

$$\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad (5.30)$$

To determine the bias of the sample mean, we are again interested in calculating its expectation:

$$\text{bias}(\hat{\mu}_m) = \mathbb{E}[\hat{\mu}_m] - \mu \quad (5.31)$$

$$= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right] - \mu \quad (5.32)$$

$$= \left(\frac{1}{m} \sum_{i=1}^m \mathbb{E}[x^{(i)}]\right) - \mu \quad (5.33)$$

$$= \left(\frac{1}{m} \sum_{i=1}^m \mu\right) - \mu \quad (5.34)$$

$$= \mu - \mu = 0 \quad (5.35)$$

Thus we find that the sample mean is an unbiased estimator of Gaussian mean parameter.

**Example: Estimators of the Variance of a Gaussian Distribution** As an example, we compare two different estimators of the variance parameter  $\sigma^2$  of a Gaussian distribution. We are interested in knowing if either estimator is biased.

The first estimator of  $\sigma^2$  we consider is known as the **sample variance**:

$$\hat{\sigma}_m^2 = \frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \hat{\mu}_m\right)^2, \quad (5.36)$$

where  $\hat{\mu}_m$  is the sample mean, defined above. More formally, we are interested in computing

$$\text{bias}(\hat{\sigma}_m^2) = \mathbb{E}[\hat{\sigma}_m^2] - \sigma^2 \quad (5.37)$$

We begin by evaluating the term  $\mathbb{E}[\hat{\sigma}_m^2]$ :

$$\mathbb{E}[\hat{\sigma}_m^2] = \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \left( x^{(i)} - \hat{\mu}_m \right)^2 \right] \quad (5.38)$$

$$= \frac{m-1}{m} \sigma^2 \quad (5.39)$$

Returning to equation 5.37, we conclude that the bias of  $\hat{\sigma}_m^2$  is  $-\sigma^2/m$ . Therefore, the sample variance is a biased estimator.

The **unbiased sample variance** estimator

$$\tilde{\sigma}_m^2 = \frac{1}{m-1} \sum_{i=1}^m \left( x^{(i)} - \hat{\mu}_m \right)^2 \quad (5.40)$$

provides an alternative approach. As the name suggests this estimator is unbiased. That is, we find that  $\mathbb{E}[\tilde{\sigma}_m^2] = \sigma^2$ :

$$\mathbb{E}[\tilde{\sigma}_m^2] = \mathbb{E} \left[ \frac{1}{m-1} \sum_{i=1}^m \left( x^{(i)} - \hat{\mu}_m \right)^2 \right] \quad (5.41)$$

$$= \frac{m}{m-1} \mathbb{E}[\hat{\sigma}_m^2] \quad (5.42)$$

$$= \frac{m}{m-1} \left( \frac{m-1}{m} \sigma^2 \right) \quad (5.43)$$

$$= \sigma^2. \quad (5.44)$$

We have two estimators: one is biased and the other is not. While unbiased estimators are clearly desirable, they are not always the “best” estimators. As we will see we often use biased estimators that possess other important properties.

### 5.4.3 Variance and Standard Error

Another property of the estimator that we might want to consider is how much we expect it to vary as a function of the data sample. Just as we computed the expectation of the estimator to determine its bias, we can compute its variance. The **variance** of an estimator is simply the variance

$$\text{Var}(\hat{\theta}) \quad (5.45)$$

where the random variable is the training set. Alternately, the square root of the variance is called the **standard error**, denoted  $\text{SE}(\hat{\theta})$ .

The variance or the standard error of an estimator provides a measure of how we would expect the estimate we compute from data to vary as we independently resample the dataset from the underlying data generating process. Just as we might like an estimator to exhibit low bias we would also like it to have relatively low variance.

When we compute any statistic using a finite number of samples, our estimate of the true underlying parameter is uncertain, in the sense that we could have obtained other samples from the same distribution and their statistics would have been different. The expected degree of variation in any estimator is a source of error that we want to quantify.

The standard error of the mean is given by

$$\text{SE}(\hat{\mu}_m) = \sqrt{\text{Var} \left[ \frac{1}{m} \sum_{i=1}^m x^{(i)} \right]} = \frac{\sigma}{\sqrt{m}}, \quad (5.46)$$

where  $\sigma^2$  is the true variance of the samples  $x^i$ . The standard error is often estimated by using an estimate of  $\sigma$ . Unfortunately, neither the square root of the sample variance nor the square root of the unbiased estimator of the variance provide an unbiased estimate of the standard deviation. Both approaches tend to underestimate the true standard deviation, but are still used in practice. The square root of the unbiased estimator of the variance is less of an underestimate. For large  $m$ , the approximation is quite reasonable.

The standard error of the mean is very useful in machine learning experiments. We often estimate the generalization error by computing the sample mean of the error on the test set. The number of examples in the test set determines the accuracy of this estimate. Taking advantage of the central limit theorem, which tells us that the mean will be approximately distributed with a normal distribution, we can use the standard error to compute the probability that the true expectation falls in any chosen interval. For example, the 95% confidence interval centered on the mean  $\hat{\mu}_m$  is

$$(\hat{\mu}_m - 1.96\text{SE}(\hat{\mu}_m), \hat{\mu}_m + 1.96\text{SE}(\hat{\mu}_m)), \quad (5.47)$$

under the normal distribution with mean  $\hat{\mu}_m$  and variance  $\text{SE}(\hat{\mu}_m)^2$ . In machine learning experiments, it is common to say that algorithm  $A$  is better than algorithm  $B$  if the upper bound of the 95% confidence interval for the error of algorithm  $A$  is less than the lower bound of the 95% confidence interval for the error of algorithm  $B$ .

**Example: Bernoulli Distribution** We once again consider a set of samples  $\{x^{(1)}, \dots, x^{(m)}\}$  drawn independently and identically from a Bernoulli distribution (recall  $P(x^{(i)}; \theta) = \theta^{x^{(i)}}(1 - \theta)^{(1-x^{(i)})}$ ). This time we are interested in computing the variance of the estimator  $\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ .

$$\text{Var}(\hat{\theta}_m) = \text{Var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) \quad (5.48)$$

$$= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(x^{(i)}) \quad (5.49)$$

$$= \frac{1}{m^2} \sum_{i=1}^m \theta(1 - \theta) \quad (5.50)$$

$$= \frac{1}{m^2} m \theta(1 - \theta) \quad (5.51)$$

$$= \frac{1}{m} \theta(1 - \theta) \quad (5.52)$$

The variance of the estimator decreases as a function of  $m$ , the number of examples in the dataset. This is a common property of popular estimators that we will return to when we discuss consistency (see section 5.4.5).

#### 5.4.4 Trading off Bias and Variance to Minimize Mean Squared Error

Bias and variance measure two different sources of error in an estimator. Bias measures the expected deviation from the true value of the function or parameter. Variance on the other hand, provides a measure of the deviation from the expected estimator value that any particular sampling of the data is likely to cause.

What happens when we are given a choice between two estimators, one with more bias and one with more variance? How do we choose between them? For example, imagine that we are interested in approximating the function shown in figure 5.2 and we are only offered the choice between a model with large bias and one that suffers from large variance. How do we choose between them?

The most common way to negotiate this trade-off is to use cross-validation. Empirically, cross-validation is highly successful on many real-world tasks. Alternatively, we can also compare the **mean squared error** (MSE) of the estimates:

$$\text{MSE} = \mathbb{E}[(\hat{\theta}_m - \theta)^2] \quad (5.53)$$

$$= \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m) \quad (5.54)$$



The MSE measures the overall expected deviation—in a squared error sense—between the estimator and the true value of the parameter  $\theta$ . As is clear from equation 5.54, evaluating the MSE incorporates both the bias and the variance. Desirable estimators are those with small MSE and these are estimators that manage to keep both their bias and variance somewhat in check.

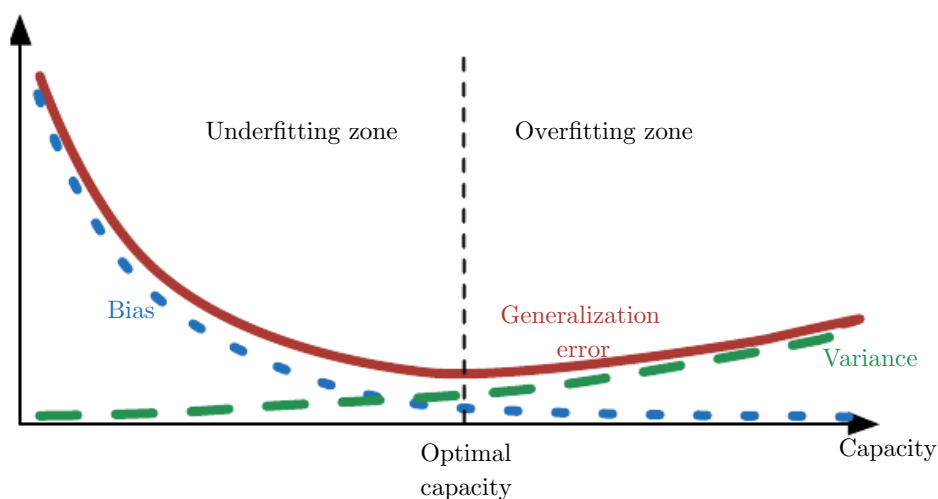


Figure 5.6: As capacity increases ( $x$ -axis), bias (dotted) tends to decrease and variance (dashed) tends to increase, yielding another U-shaped curve for generalization error (bold curve). If we vary capacity along one axis, there is an optimal capacity, with underfitting when the capacity is below this optimum and overfitting when it is above. This relationship is similar to the relationship between capacity, underfitting, and overfitting, discussed in section 5.2 and figure 5.3.

The relationship between bias and variance is tightly linked to the machine learning concepts of capacity, underfitting and overfitting. In the case where generalization error is measured by the MSE (where bias and variance are meaningful components of generalization error), increasing capacity tends to increase variance and decrease bias. This is illustrated in figure 5.6, where we see again the U-shaped curve of generalization error as a function of capacity.

### 5.4.5 Consistency

So far we have discussed the properties of various estimators for a training set of fixed size. Usually, we are also concerned with the behavior of an estimator as the amount of training data grows. In particular, we usually wish that, as the number of data points  $m$  in our dataset increases, our point estimates converge to the true

value of the corresponding parameters. More formally, we would like that

$$\text{plim}_{m \rightarrow \infty} \hat{\theta}_m = \theta. \quad (5.55)$$

The symbol  $\text{plim}$  indicates convergence in probability, meaning that for any  $\epsilon > 0$ ,  $P(|\hat{\theta}_m - \theta| > \epsilon) \rightarrow 0$  as  $m \rightarrow \infty$ . The condition described by equation 5.55 is known as **consistency**. It is sometimes referred to as weak consistency, with strong consistency referring to the **almost sure** convergence of  $\hat{\theta}$  to  $\theta$ . **Almost sure convergence** of a sequence of random variables  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$  to a value  $\mathbf{x}$  occurs when  $p(\lim_{m \rightarrow \infty} \mathbf{x}^{(m)} = \mathbf{x}) = 1$ .

Consistency ensures that the bias induced by the estimator diminishes as the number of data examples grows. However, the reverse is not true—asymptotic unbiasedness does not imply consistency. For example, consider estimating the mean parameter  $\mu$  of a normal distribution  $\mathcal{N}(x; \mu, \sigma^2)$ , with a dataset consisting of  $m$  samples:  $\{x^{(1)}, \dots, x^{(m)}\}$ . We could use the first sample  $x^{(1)}$  of the dataset as an unbiased estimator:  $\hat{\theta} = x^{(1)}$ . In that case,  $\mathbb{E}(\hat{\theta}_m) = \theta$  so the estimator is unbiased no matter how many data points are seen. This, of course, implies that the estimate is asymptotically unbiased. However, this is not a consistent estimator as it is *not* the case that  $\hat{\theta}_m \rightarrow \theta$  as  $m \rightarrow \infty$ .

## 5.5 Maximum Likelihood Estimation

Previously, we have seen some definitions of common estimators and analyzed their properties. But where did these estimators come from? Rather than guessing that some function might make a good estimator and then analyzing its bias and variance, we would like to have some principle from which we can derive specific functions that are good estimators for different models.

The most common such principle is the maximum likelihood principle.

Consider a set of  $m$  examples  $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  drawn independently from the true but unknown data generating distribution  $p_{\text{data}}(\mathbf{x})$ .

Let  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  be a parametric family of probability distributions over the same space indexed by  $\boldsymbol{\theta}$ . In other words,  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  maps any configuration  $\mathbf{x}$  to a real number estimating the true probability  $p_{\text{data}}(\mathbf{x})$ .

The maximum likelihood estimator for  $\boldsymbol{\theta}$  is then defined as

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) \quad (5.56)$$

$$= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad (5.57)$$

This product over many probabilities can be inconvenient for a variety of reasons. For example, it is prone to numerical underflow. To obtain a more convenient but equivalent optimization problem, we observe that taking the logarithm of the likelihood does not change its  $\arg \max$  but does conveniently transform a product into a sum:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}). \quad (5.58)$$

Because the  $\arg \max$  does not change when we rescale the cost function, we can divide by  $m$  to obtain a version of the criterion that is expressed as an expectation with respect to the empirical distribution  $\hat{p}_{\text{data}}$  defined by the training data:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}). \quad (5.59)$$

One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution  $\hat{p}_{\text{data}}$  defined by the training set and the model distribution, with the degree of dissimilarity between the two measured by the KL divergence. The KL divergence is given by

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]. \quad (5.60)$$

The term on the left is a function only of the data generating process, not the model. This means when we train the model to minimize the KL divergence, we need only minimize

$$- \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})] \quad (5.61)$$

which is of course the same as the maximization in equation 5.59.

Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions. Many authors use the term “cross-entropy” to identify specifically the negative log-likelihood of a Bernoulli or softmax distribution, but that is a misnomer. Any loss consisting of a negative log-likelihood is a cross-entropy between the empirical distribution defined by the training set and the probability distribution defined by model. For example, mean squared error is the cross-entropy between the empirical distribution and a Gaussian model.

We can thus see maximum likelihood as an attempt to make the model distribution match the empirical distribution  $\hat{p}_{\text{data}}$ . Ideally, we would like to match the true data generating distribution  $p_{\text{data}}$ , but we have no direct access to this distribution.

While the optimal  $\boldsymbol{\theta}$  is the same regardless of whether we are maximizing the likelihood or minimizing the KL divergence, the values of the objective functions

are different. In software, we often phrase both as minimizing a cost function. Maximum likelihood thus becomes minimization of the negative log-likelihood (NLL), or equivalently, minimization of the cross entropy. The perspective of maximum likelihood as minimum KL divergence becomes helpful in this case because the KL divergence has a known minimum value of zero. The negative log-likelihood can actually become negative when  $\mathbf{x}$  is real-valued.

### 5.5.1 Conditional Log-Likelihood and Mean Squared Error

The maximum likelihood estimator can readily be generalized to the case where our goal is to estimate a conditional probability  $P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$  in order to predict  $\mathbf{y}$  given  $\mathbf{x}$ . This is actually the most common situation because it forms the basis for most supervised learning. If  $\mathbf{X}$  represents all our inputs and  $\mathbf{Y}$  all our observed targets, then the conditional maximum likelihood estimator is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta}). \quad (5.62)$$

If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}). \quad (5.63)$$

**Example: Linear Regression as Maximum Likelihood** Linear regression, introduced earlier in section 5.1.4, may be justified as a maximum likelihood procedure. Previously, we motivated linear regression as an algorithm that learns to take an input  $\mathbf{x}$  and produce an output value  $\hat{y}$ . The mapping from  $\mathbf{x}$  to  $\hat{y}$  is chosen to minimize mean squared error, a criterion that we introduced more or less arbitrarily. We now revisit linear regression from the point of view of maximum likelihood estimation. Instead of producing a single prediction  $\hat{y}$ , we now think of the model as producing a conditional distribution  $p(y \mid \mathbf{x})$ . We can imagine that with an infinitely large training set, we might see several training examples with the same input value  $\mathbf{x}$  but different values of  $y$ . The goal of the learning algorithm is now to fit the distribution  $p(y \mid \mathbf{x})$  to all of those different  $y$  values that are all compatible with  $\mathbf{x}$ . To derive the same linear regression algorithm we obtained before, we define  $p(y \mid \mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$ . The function  $\hat{y}(\mathbf{x}; \mathbf{w})$  gives the prediction of the mean of the Gaussian. In this example, we assume that the variance is fixed to some constant  $\sigma^2$  chosen by the user. We will see that this choice of the functional form of  $p(y \mid \mathbf{x})$  causes the maximum likelihood estimation procedure to yield the same learning algorithm as we developed before. Since the

examples are assumed to be i.i.d., the conditional log-likelihood (equation 5.63) is given by

$$\sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad (5.64)$$

$$= -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}, \quad (5.65)$$

where  $\hat{y}^{(i)}$  is the output of the linear regression on the  $i$ -th input  $\mathbf{x}^{(i)}$  and  $m$  is the number of the training examples. Comparing the log-likelihood with the mean squared error,

$$\text{MSE}_{\text{train}} = \frac{1}{m} \sum_{i=1}^m \|\hat{y}^{(i)} - y^{(i)}\|^2, \quad (5.66)$$

we immediately see that maximizing the log-likelihood with respect to  $\mathbf{w}$  yields the same estimate of the parameters  $\mathbf{w}$  as does minimizing the mean squared error. The two criteria have different values but the same location of the optimum. This justifies the use of the MSE as a maximum likelihood estimation procedure. As we will see, the maximum likelihood estimator has several desirable properties.

### 5.5.2 Properties of Maximum Likelihood

The main appeal of the maximum likelihood estimator is that it can be shown to be the best estimator asymptotically, as the number of examples  $m \rightarrow \infty$ , in terms of its rate of convergence as  $m$  increases.

Under appropriate conditions, the maximum likelihood estimator has the property of consistency (see section 5.4.5 above), meaning that as the number of training examples approaches infinity, the maximum likelihood estimate of a parameter converges to the true value of the parameter. These conditions are:

- The true distribution  $p_{\text{data}}$  must lie within the model family  $p_{\text{model}}(\cdot; \boldsymbol{\theta})$ . Otherwise, no estimator can recover  $p_{\text{data}}$ .
- The true distribution  $p_{\text{data}}$  must correspond to exactly one value of  $\boldsymbol{\theta}$ . Otherwise, maximum likelihood can recover the correct  $p_{\text{data}}$ , but will not be able to determine which value of  $\boldsymbol{\theta}$  was used by the data generating processing.

There are other inductive principles besides the maximum likelihood estimator, many of which share the property of being consistent estimators. However,

consistent estimators can differ in their **statistic efficiency**, meaning that one consistent estimator may obtain lower generalization error for a fixed number of samples  $m$ , or equivalently, may require fewer examples to obtain a fixed level of generalization error.

Statistical efficiency is typically studied in the **parametric case** (like in linear regression) where our goal is to estimate the value of a parameter (and assuming it is possible to identify the true parameter), not the value of a function. A way to measure how close we are to the true parameter is by the expected mean squared error, computing the squared difference between the estimated and true parameter values, where the expectation is over  $m$  training samples from the data generating distribution. That parametric mean squared error decreases as  $m$  increases, and for  $m$  large, the Cramér-Rao lower bound (Rao, 1945; Cramér, 1946) shows that no consistent estimator has a lower mean squared error than the maximum likelihood estimator.

For these reasons (consistency and efficiency), maximum likelihood is often considered the preferred estimator to use for machine learning. When the number of examples is small enough to yield overfitting behavior, regularization strategies such as weight decay may be used to obtain a biased version of maximum likelihood that has less variance when training data is limited.

## 5.6 Bayesian Statistics

So far we have discussed **frequentist statistics** and approaches based on estimating a single value of  $\theta$ , then making all predictions thereafter based on that one estimate. Another approach is to consider all possible values of  $\theta$  when making a prediction. The latter is the domain of **Bayesian statistics**.

As discussed in section 5.4.1, the frequentist perspective is that the true parameter value  $\theta$  is fixed but unknown, while the point estimate  $\hat{\theta}$  is a random variable on account of it being a function of the dataset (which is seen as random).

The Bayesian perspective on statistics is quite different. The Bayesian uses probability to reflect degrees of certainty of states of knowledge. The dataset is directly observed and so is not random. On the other hand, the true parameter  $\theta$  is unknown or uncertain and thus is represented as a random variable.

Before observing the data, we represent our knowledge of  $\theta$  using the **prior probability distribution**,  $p(\theta)$  (sometimes referred to as simply “the prior”). Generally, the machine learning practitioner selects a prior distribution that is quite broad (i.e. with high entropy) to reflect a high degree of uncertainty in the

value of  $\theta$  before observing any data. For example, one might assume *a priori* that  $\theta$  lies in some finite range or volume, with a uniform distribution. Many priors instead reflect a preference for “simpler” solutions (such as smaller magnitude coefficients, or a function that is closer to being constant).

Now consider that we have a set of data samples  $\{x^{(1)}, \dots, x^{(m)}\}$ . We can recover the effect of data on our belief about  $\theta$  by combining the data likelihood  $p(x^{(1)}, \dots, x^{(m)} | \theta)$  with the prior via Bayes’ rule:

$$p(\theta | x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} | \theta)p(\theta)}{p(x^{(1)}, \dots, x^{(m)})} \quad (5.67)$$

In the scenarios where Bayesian estimation is typically used, the prior begins as a relatively uniform or Gaussian distribution with high entropy, and the observation of the data usually causes the posterior to lose entropy and concentrate around a few highly likely values of the parameters.

Relative to maximum likelihood estimation, Bayesian estimation offers two important differences. First, unlike the maximum likelihood approach that makes predictions using a point estimate of  $\theta$ , the Bayesian approach is to make predictions using a full distribution over  $\theta$ . For example, after observing  $m$  examples, the predicted distribution over the next data sample,  $x^{(m+1)}$ , is given by

$$p(x^{(m+1)} | x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} | \theta)p(\theta | x^{(1)}, \dots, x^{(m)}) d\theta. \quad (5.68)$$

Here each value of  $\theta$  with positive probability density contributes to the prediction of the next example, with the contribution weighted by the posterior density itself. After having observed  $\{x^{(1)}, \dots, x^{(m)}\}$ , if we are still quite uncertain about the value of  $\theta$ , then this uncertainty is incorporated directly into any predictions we might make.

In section 5.4, we discussed how the frequentist approach addresses the uncertainty in a given point estimate of  $\theta$  by evaluating its variance. The variance of the estimator is an assessment of how the estimate might change with alternative samplings of the observed data. The Bayesian answer to the question of how to deal with the uncertainty in the estimator is to simply integrate over it, which tends to protect well against overfitting. This integral is of course just an application of the laws of probability, making the Bayesian approach simple to justify, while the frequentist machinery for constructing an estimator is based on the rather ad hoc decision to summarize all knowledge contained in the dataset with a single point estimate.

The second important difference between the Bayesian approach to estimation and the maximum likelihood approach is due to the contribution of the Bayesian

prior distribution. The prior has an influence by shifting probability mass density towards regions of the parameter space that are preferred *a priori*. In practice, the prior often expresses a preference for models that are simpler or more smooth. Critics of the Bayesian approach identify the prior as a source of subjective human judgment impacting the predictions.

Bayesian methods typically generalize much better when limited training data is available, but typically suffer from high computational cost when the number of training examples is large.

**Example: Bayesian Linear Regression** Here we consider the Bayesian estimation approach to learning the linear regression parameters. In linear regression, we learn a linear mapping from an input vector  $\mathbf{x} \in \mathbb{R}^n$  to predict the value of a scalar  $y \in \mathbb{R}$ . The prediction is parametrized by the vector  $\mathbf{w} \in \mathbb{R}^n$ :

$$\hat{y} = \mathbf{w}^\top \mathbf{x}. \quad (5.69)$$

Given a set of  $m$  training samples  $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$ , we can express the prediction of  $y$  over the entire training set as:

$$\hat{\mathbf{y}}^{(\text{train})} = \mathbf{X}^{(\text{train})} \mathbf{w}. \quad (5.70)$$

Expressed as a Gaussian conditional distribution on  $\mathbf{y}^{(\text{train})}$ , we have

$$p(\mathbf{y}^{(\text{train})} \mid \mathbf{X}^{(\text{train})}, \mathbf{w}) = \mathcal{N}(\mathbf{y}^{(\text{train})}; \mathbf{X}^{(\text{train})} \mathbf{w}, \mathbf{I}) \quad (5.71)$$

$$\propto \exp \left( -\frac{1}{2} (\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w})^\top (\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w}) \right), \quad (5.72)$$

where we follow the standard MSE formulation in assuming that the Gaussian variance on  $y$  is one. In what follows, to reduce the notational burden, we refer to  $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$  as simply  $(\mathbf{X}, \mathbf{y})$ .

To determine the posterior distribution over the model parameter vector  $\mathbf{w}$ , we first need to specify a prior distribution. The prior should reflect our naive belief about the value of these parameters. While it is sometimes difficult or unnatural to express our prior beliefs in terms of the parameters of the model, in practice we typically assume a fairly broad distribution expressing a high degree of uncertainty about  $\boldsymbol{\theta}$ . For real-valued parameters it is common to use a Gaussian as a prior distribution:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0) \propto \exp \left( -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Lambda}_0^{-1} (\mathbf{w} - \boldsymbol{\mu}_0) \right), \quad (5.73)$$



where  $\boldsymbol{\mu}_0$  and  $\boldsymbol{\Lambda}_0$  are the prior distribution mean vector and covariance matrix respectively.<sup>1</sup>

With the prior thus specified, we can now proceed in determining the **posterior** distribution over the model parameters.

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w})p(\mathbf{w}) \quad (5.74)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\right) \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Lambda}_0^{-1}(\mathbf{w} - \boldsymbol{\mu}_0)\right) \quad (5.75)$$

$$\propto \exp\left(-\frac{1}{2}\left(-2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \boldsymbol{\Lambda}_0^{-1}\mathbf{w} - 2\boldsymbol{\mu}_0^\top \boldsymbol{\Lambda}_0^{-1}\mathbf{w}\right)\right). \quad (5.76)$$

We now define  $\boldsymbol{\Lambda}_m = (\mathbf{X}^\top \mathbf{X} + \boldsymbol{\Lambda}_0^{-1})^{-1}$  and  $\boldsymbol{\mu}_m = \boldsymbol{\Lambda}_m (\mathbf{X}^\top \mathbf{y} + \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0)$ . Using these new variables, we find that the posterior may be rewritten as a Gaussian distribution:

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top \boldsymbol{\Lambda}_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m) + \frac{1}{2}\boldsymbol{\mu}_m^\top \boldsymbol{\Lambda}_m^{-1} \boldsymbol{\mu}_m\right) \quad (5.77)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top \boldsymbol{\Lambda}_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m)\right). \quad (5.78)$$

All terms that do not include the parameter vector  $\mathbf{w}$  have been omitted; they are implied by the fact that the distribution must be normalized to integrate to 1. Equation 3.23 shows how to normalize a multivariate Gaussian distribution.

Examining this posterior distribution allows us to gain some intuition for the effect of Bayesian inference. In most situations, we set  $\boldsymbol{\mu}_0$  to  $\mathbf{0}$ . If we set  $\boldsymbol{\Lambda}_0 = \frac{1}{\alpha} \mathbf{I}$ , then  $\boldsymbol{\mu}_m$  gives the same estimate of  $\mathbf{w}$  as does frequentist linear regression with a weight decay penalty of  $\alpha \mathbf{w}^\top \mathbf{w}$ . One difference is that the Bayesian estimate is undefined if  $\alpha$  is set to zero—we are not allowed to begin the Bayesian learning process with an infinitely wide prior on  $\mathbf{w}$ . The more important difference is that the Bayesian estimate provides a covariance matrix, showing how likely all the different values of  $\mathbf{w}$  are, rather than providing only the estimate  $\boldsymbol{\mu}_m$ .

### 5.6.1 Maximum *A Posteriori* (MAP) Estimation

While the most principled approach is to make predictions using the full Bayesian posterior distribution over the parameter  $\boldsymbol{\theta}$ , it is still often desirable to have a

---

<sup>1</sup> Unless there is a reason to assume a particular covariance structure, we typically assume a diagonal covariance matrix  $\boldsymbol{\Lambda}_0 = \text{diag}(\boldsymbol{\lambda}_0)$ .

single point estimate. One common reason for desiring a point estimate is that most operations involving the Bayesian posterior for most interesting models are intractable, and a point estimate offers a tractable approximation. Rather than simply returning to the maximum likelihood estimate, we can still gain some of the benefit of the Bayesian approach by allowing the prior to influence the choice of the point estimate. One rational way to do this is to choose the **maximum a posteriori** (MAP) point estimate. The MAP estimate chooses the point of maximal posterior probability (or maximal probability density in the more common case of continuous  $\theta$ ):

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta \mid \mathbf{x}) = \arg \max_{\theta} \log p(\mathbf{x} \mid \theta) + \log p(\theta). \quad (5.79)$$

We recognize, above on the right hand side,  $\log p(\mathbf{x} \mid \theta)$ , i.e. the standard log-likelihood term, and  $\log p(\theta)$ , corresponding to the prior distribution.

As an example, consider a linear regression model with a Gaussian prior on the weights  $\mathbf{w}$ . If this prior is given by  $\mathcal{N}(\mathbf{w}; \mathbf{0}, \frac{1}{\lambda} \mathbf{I}^2)$ , then the log-prior term in equation 5.79 is proportional to the familiar  $\lambda \mathbf{w}^\top \mathbf{w}$  weight decay penalty, plus a term that does not depend on  $\mathbf{w}$  and does not affect the learning process. MAP Bayesian inference with a Gaussian prior on the weights thus corresponds to weight decay.

As with full Bayesian inference, MAP Bayesian inference has the advantage of leveraging information that is brought by the prior and cannot be found in the training data. This additional information helps to reduce the variance in the MAP point estimate (in comparison to the ML estimate). However, it does so at the price of increased bias.

Many regularized estimation strategies, such as maximum likelihood learning regularized with weight decay, can be interpreted as making the MAP approximation to Bayesian inference. This view applies when the regularization consists of adding an extra term to the objective function that corresponds to  $\log p(\theta)$ . Not all regularization penalties correspond to MAP Bayesian inference. For example, some regularizer terms may not be the logarithm of a probability distribution. Other regularization terms depend on the data, which of course a prior probability distribution is not allowed to do.

MAP Bayesian inference provides a straightforward way to design complicated yet interpretable regularization terms. For example, a more complicated penalty term can be derived by using a mixture of Gaussians, rather than a single Gaussian distribution, as the prior (Nowlan and Hinton, 1992).

## 5.7 Supervised Learning Algorithms

Recall from section 5.1.3 that supervised learning algorithms are, roughly speaking, learning algorithms that learn to associate some input with some output, given a training set of examples of inputs  $\mathbf{x}$  and outputs  $\mathbf{y}$ . In many cases the outputs  $\mathbf{y}$  may be difficult to collect automatically and must be provided by a human “supervisor,” but the term still applies even when the training set targets were collected automatically.

### 5.7.1 Probabilistic Supervised Learning

Most supervised learning algorithms in this book are based on estimating a probability distribution  $p(y \mid \mathbf{x})$ . We can do this simply by using maximum likelihood estimation to find the best parameter vector  $\boldsymbol{\theta}$  for a parametric family of distributions  $p(y \mid \mathbf{x}; \boldsymbol{\theta})$ .

We have already seen that linear regression corresponds to the family

$$p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y; \boldsymbol{\theta}^\top \mathbf{x}, I). \quad (5.80)$$

We can generalize linear regression to the classification scenario by defining a different family of probability distributions. If we have two classes, class 0 and class 1, then we need only specify the probability of one of these classes. The probability of class 1 determines the probability of class 0, because these two values must add up to 1.

The normal distribution over real-valued numbers that we used for linear regression is parametrized in terms of a mean. Any value we supply for this mean is valid. A distribution over a binary variable is slightly more complicated, because its mean must always be between 0 and 1. One way to solve this problem is to use the logistic sigmoid function to squash the output of the linear function into the interval  $(0, 1)$  and interpret that value as a probability:

$$p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x}). \quad (5.81)$$

This approach is known as **logistic regression** (a somewhat strange name since we use the model for classification rather than regression).

In the case of linear regression, we were able to find the optimal weights by solving the normal equations. Logistic regression is somewhat more difficult. There is no closed-form solution for its optimal weights. Instead, we must search for them by maximizing the log-likelihood. We can do this by minimizing the negative log-likelihood (NLL) using gradient descent.

This same strategy can be applied to essentially any supervised learning problem, by writing down a parametric family of conditional probability distributions over the right kind of input and output variables.

### 5.7.2 Support Vector Machines

One of the most influential approaches to supervised learning is the support vector machine (Boser *et al.*, 1992; Cortes and Vapnik, 1995). This model is similar to logistic regression in that it is driven by a linear function  $\mathbf{w}^\top \mathbf{x} + b$ . Unlike logistic regression, the support vector machine does not provide probabilities, but only outputs a class identity. The SVM predicts that the positive class is present when  $\mathbf{w}^\top \mathbf{x} + b$  is positive. Likewise, it predicts that the negative class is present when  $\mathbf{w}^\top \mathbf{x} + b$  is negative.

One key innovation associated with support vector machines is the **kernel trick**. The kernel trick consists of observing that many machine learning algorithms can be written exclusively in terms of dot products between examples. For example, it can be shown that the linear function used by the support vector machine can be re-written as

$$\mathbf{w}^\top \mathbf{x} + b = b + \sum_{i=1}^m \alpha_i \mathbf{x}^\top \mathbf{x}^{(i)} \quad (5.82)$$

where  $\mathbf{x}^{(i)}$  is a training example and  $\boldsymbol{\alpha}$  is a vector of coefficients. Rewriting the learning algorithm this way allows us to replace  $\mathbf{x}$  by the output of a given feature function  $\phi(\mathbf{x})$  and the dot product with a function  $k(\mathbf{x}, \mathbf{x}^{(i)}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}^{(i)})$  called a **kernel**. The  $\cdot$  operator represents an inner product analogous to  $\phi(\mathbf{x})^\top \phi(\mathbf{x}^{(i)})$ . For some feature spaces, we may not use literally the vector inner product. In some infinite dimensional spaces, we need to use other kinds of inner products, for example, inner products based on integration rather than summation. A complete development of these kinds of inner products is beyond the scope of this book.

After replacing dot products with kernel evaluations, we can make predictions using the function

$$f(\mathbf{x}) = b + \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)}). \quad (5.83)$$

This function is nonlinear with respect to  $\mathbf{x}$ , but the relationship between  $\phi(\mathbf{x})$  and  $f(\mathbf{x})$  is linear. Also, the relationship between  $\boldsymbol{\alpha}$  and  $f(\mathbf{x})$  is linear. The kernel-based function is exactly equivalent to preprocessing the data by applying  $\phi(\mathbf{x})$  to all inputs, then learning a linear model in the new transformed space.

The kernel trick is powerful for two reasons. First, it allows us to learn models that are nonlinear as a function of  $\mathbf{x}$  using convex optimization techniques that are

guaranteed to converge efficiently. This is possible because we consider  $\phi$  fixed and optimize only  $\alpha$ , i.e., the optimization algorithm can view the decision function as being linear in a different space. Second, the kernel function  $k$  often admits an implementation that is significantly more computationally efficient than naively constructing two  $\phi(\mathbf{x})$  vectors and explicitly taking their dot product.

In some cases,  $\phi(\mathbf{x})$  can even be infinite dimensional, which would result in an infinite computational cost for the naive, explicit approach. In many cases,  $k(\mathbf{x}, \mathbf{x}')$  is a nonlinear, tractable function of  $\mathbf{x}$  even when  $\phi(\mathbf{x})$  is intractable. As an example of an infinite-dimensional feature space with a tractable kernel, we construct a feature mapping  $\phi(x)$  over the non-negative integers  $x$ . Suppose that this mapping returns a vector containing  $x$  ones followed by infinitely many zeros. We can write a kernel function  $k(x, x^{(i)}) = \min(x, x^{(i)})$  that is exactly equivalent to the corresponding infinite-dimensional dot product.

The most commonly used kernel is the **Gaussian kernel**

$$k(\mathbf{u}, \mathbf{v}) = \mathcal{N}(\mathbf{u} - \mathbf{v}; 0, \sigma^2 \mathbf{I}) \quad (5.84)$$

where  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the standard normal density. This kernel is also known as the **radial basis function** (RBF) kernel, because its value decreases along lines in  $\mathbf{v}$  space radiating outward from  $\mathbf{u}$ . The Gaussian kernel corresponds to a dot product in an infinite-dimensional space, but the derivation of this space is less straightforward than in our example of the min kernel over the integers.

We can think of the Gaussian kernel as performing a kind of **template matching**. A training example  $\mathbf{x}$  associated with training label  $y$  becomes a template for class  $y$ . When a test point  $\mathbf{x}'$  is near  $\mathbf{x}$  according to Euclidean distance, the Gaussian kernel has a large response, indicating that  $\mathbf{x}'$  is very similar to the  $\mathbf{x}$  template. The model then puts a large weight on the associated training label  $y$ . Overall, the prediction will combine many such training labels weighted by the similarity of the corresponding training examples.

Support vector machines are not the only algorithm that can be enhanced using the kernel trick. Many other linear models can be enhanced in this way. The category of algorithms that employ the kernel trick is known as **kernel machines** or **kernel methods** (Williams and Rasmussen, 1996; Schölkopf *et al.*, 1999).

A major drawback to kernel machines is that the cost of evaluating the decision function is linear in the number of training examples, because the  $i$ -th example contributes a term  $\alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$  to the decision function. Support vector machines are able to mitigate this by learning an  $\alpha$  vector that contains mostly zeros. Classifying a new example then requires evaluating the kernel function only for the training examples that have non-zero  $\alpha_i$ . These training examples are known

as **support vectors**.

Kernel machines also suffer from a high computational cost of training when the dataset is large. We will revisit this idea in section 5.9. Kernel machines with generic kernels struggle to generalize well. We will explain why in section 5.11. The modern incarnation of deep learning was designed to overcome these limitations of kernel machines. The current deep learning renaissance began when Hinton *et al.* (2006) demonstrated that a neural network could outperform the RBF kernel SVM on the MNIST benchmark.

### 5.7.3 Other Simple Supervised Learning Algorithms

We have already briefly encountered another non-probabilistic supervised learning algorithm, nearest neighbor regression. More generally,  $k$ -nearest neighbors is a family of techniques that can be used for classification or regression. As a non-parametric learning algorithm,  $k$ -nearest neighbors is not restricted to a fixed number of parameters. We usually think of the  $k$ -nearest neighbors algorithm as not having any parameters, but rather implementing a simple function of the training data. In fact, there is not even really a training stage or learning process. Instead, at test time, when we want to produce an output  $y$  for a new test input  $\mathbf{x}$ , we find the  $k$ -nearest neighbors to  $\mathbf{x}$  in the training data  $\mathbf{X}$ . We then return the average of the corresponding  $y$  values in the training set. This works for essentially any kind of supervised learning where we can define an average over  $y$  values. In the case of classification, we can average over one-hot code vectors  $\mathbf{c}$  with  $c_y = 1$  and  $c_i = 0$  for all other values of  $i$ . We can then interpret the average over these one-hot codes as giving a probability distribution over classes. As a non-parametric learning algorithm,  $k$ -nearest neighbor can achieve very high capacity. For example, suppose we have a multiclass classification task and measure performance with 0-1 loss. In this setting, 1-nearest neighbor converges to double the Bayes error as the number of training examples approaches infinity. The error in excess of the Bayes error results from choosing a single neighbor by breaking ties between equally distant neighbors randomly. When there is infinite training data, all test points  $\mathbf{x}$  will have infinitely many training set neighbors at distance zero. If we allow the algorithm to use all of these neighbors to vote, rather than randomly choosing one of them, the procedure converges to the Bayes error rate. The high capacity of  $k$ -nearest neighbors allows it to obtain high accuracy given a large training set. However, it does so at high computational cost, and it may generalize very badly given a small, finite training set. One weakness of  $k$ -nearest neighbors is that it cannot learn that one feature is more discriminative than another. For example, imagine we have a regression task with  $\mathbf{x} \in \mathbb{R}^{100}$  drawn from an isotropic Gaussian

distribution, but only a single variable  $x_1$  is relevant to the output. Suppose further that this feature simply encodes the output directly, i.e. that  $y = x_1$  in all cases. Nearest neighbor regression will not be able to detect this simple pattern. The nearest neighbor of most points  $\mathbf{x}$  will be determined by the large number of features  $x_2$  through  $x_{100}$ , not by the lone feature  $x_1$ . Thus the output on small training sets will essentially be random.

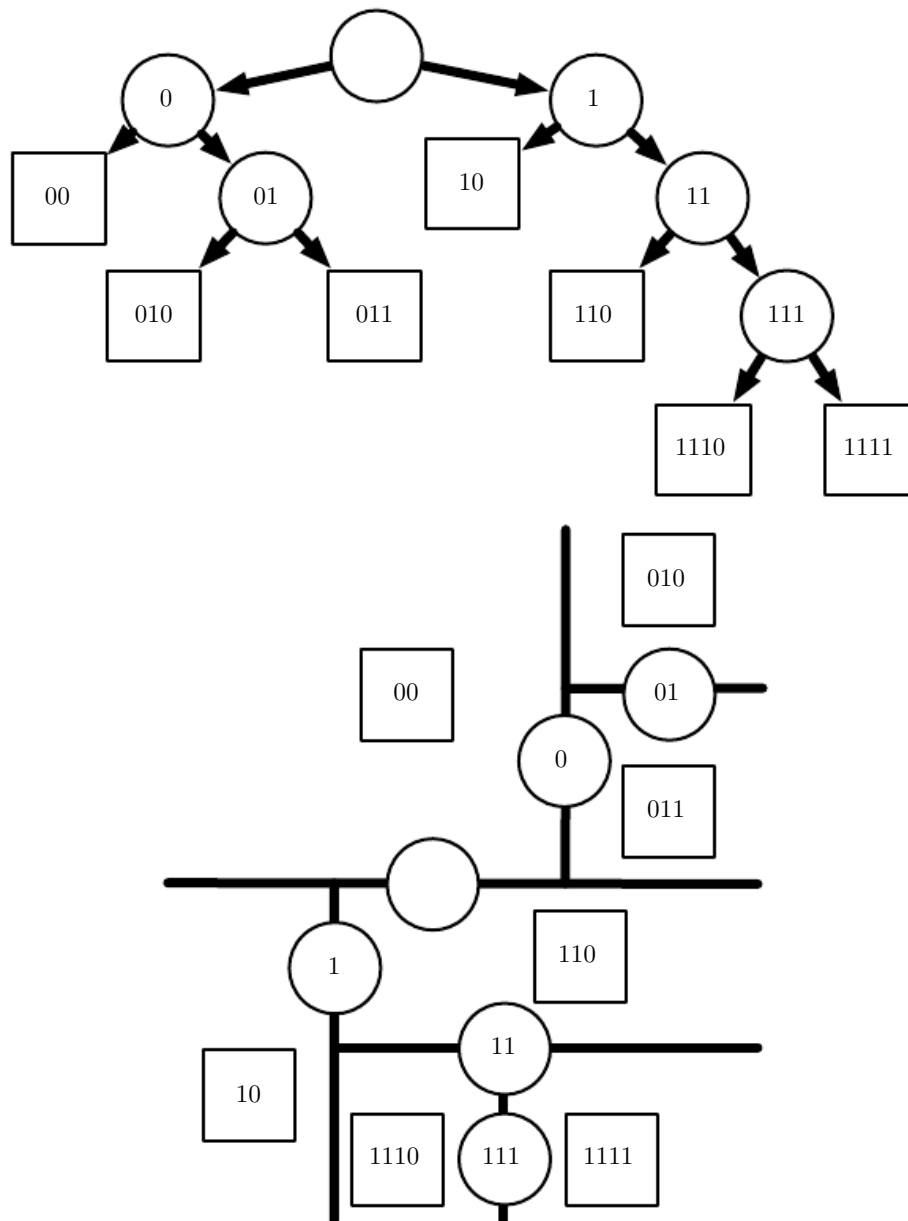


Figure 5.7: Diagrams describing how a decision tree works. (*Top*) Each node of the tree chooses to send the input example to the child node on the left (0) or to the child node on the right (1). Internal nodes are drawn as circles and leaf nodes as squares. Each node is displayed with a binary string identifier corresponding to its position in the tree, obtained by appending a bit to its parent identifier (0=choose left or top, 1=choose right or bottom). (*Bottom*) The tree divides space into regions. The 2D plane shows how a decision tree might divide  $\mathbb{R}^2$ . The nodes of the tree are plotted in this plane, with each internal node drawn along the dividing line it uses to categorize examples, and leaf nodes drawn in the center of the region of examples they receive. The result is a piecewise-constant function, with one piece per leaf. Each leaf requires at least one training example to define, so it is not possible for the decision tree to learn a function that has more local maxima than the number of training examples.



Another type of learning algorithm that also breaks the input space into regions and has separate parameters for each region is the **decision tree** (Breiman *et al.*, 1984) and its many variants. As shown in figure 5.7, each node of the decision tree is associated with a region in the input space, and internal nodes break that region into one sub-region for each child of the node (typically using an axis-aligned cut). Space is thus sub-divided into non-overlapping regions, with a one-to-one correspondence between leaf nodes and input regions. Each leaf node usually maps every point in its input region to the same output. Decision trees are usually trained with specialized algorithms that are beyond the scope of this book. The learning algorithm can be considered non-parametric if it is allowed to learn a tree of arbitrary size, though decision trees are usually regularized with size constraints that turn them into parametric models in practice. Decision trees as they are typically used, with axis-aligned splits and constant outputs within each node, struggle to solve some problems that are easy even for logistic regression. For example, if we have a two-class problem and the positive class occurs wherever  $x_2 > x_1$ , the decision boundary is not axis-aligned. The decision tree will thus need to approximate the decision boundary with many nodes, implementing a step function that constantly walks back and forth across the true decision function with axis-aligned steps.

As we have seen, nearest neighbor predictors and decision trees have many limitations. Nonetheless, they are useful learning algorithms when computational resources are constrained. We can also build intuition for more sophisticated learning algorithms by thinking about the similarities and differences between sophisticated algorithms and  $k$ -NN or decision tree baselines.

See Murphy (2012), Bishop (2006), Hastie *et al.* (2001) or other machine learning textbooks for more material on traditional supervised learning algorithms.

## 5.8 Unsupervised Learning Algorithms

Recall from section 5.1.3 that unsupervised algorithms are those that experience only “features” but not a supervision signal. The distinction between supervised and unsupervised algorithms is not formally and rigidly defined because there is no objective test for distinguishing whether a value is a feature or a target provided by a supervisor. Informally, unsupervised learning refers to most attempts to extract information from a distribution that do not require human labor to annotate examples. The term is usually associated with density estimation, learning to draw samples from a distribution, learning to denoise data from some distribution, finding a manifold that the data lies near, or clustering the data into groups of

related examples.

A classic unsupervised learning task is to find the “best” representation of the data. By ‘best’ we can mean different things, but generally speaking we are looking for a representation that preserves as much information about  $\mathbf{x}$  as possible while obeying some penalty or constraint aimed at keeping the representation *simpler* or more accessible than  $\mathbf{x}$  itself.

There are multiple ways of defining a *simpler* representation. Three of the most common include lower dimensional representations, sparse representations and independent representations. Low-dimensional representations attempt to compress as much information about  $x$  as possible in a smaller representation. Sparse representations (Barlow, 1989; Olshausen and Field, 1996; Hinton and Ghahramani, 1997) embed the dataset into a representation whose entries are mostly zeroes for most inputs. The use of sparse representations typically requires increasing the dimensionality of the representation, so that the representation becoming mostly zeroes does not discard too much information. This results in an overall structure of the representation that tends to distribute data along the axes of the representation space. Independent representations attempt to *disentangle* the sources of variation underlying the data distribution such that the dimensions of the representation are statistically independent.

Of course these three criteria are certainly not mutually exclusive. Low-dimensional representations often yield elements that have fewer or weaker dependencies than the original high-dimensional data. This is because one way to reduce the size of a representation is to find and remove redundancies. Identifying and removing more redundancy allows the dimensionality reduction algorithm to achieve more compression while discarding less information.

The notion of representation is one of the central themes of deep learning and therefore one of the central themes in this book. In this section, we develop some simple examples of representation learning algorithms. Together, these example algorithms show how to operationalize all three of the criteria above. Most of the remaining chapters introduce additional representation learning algorithms that develop these criteria in different ways or introduce other criteria.

### 5.8.1 Principal Components Analysis

In section 2.12, we saw that the principal components analysis algorithm provides a means of compressing data. We can also view PCA as an unsupervised learning algorithm that learns a representation of data. This representation is based on two of the criteria for a simple representation described above. PCA learns a

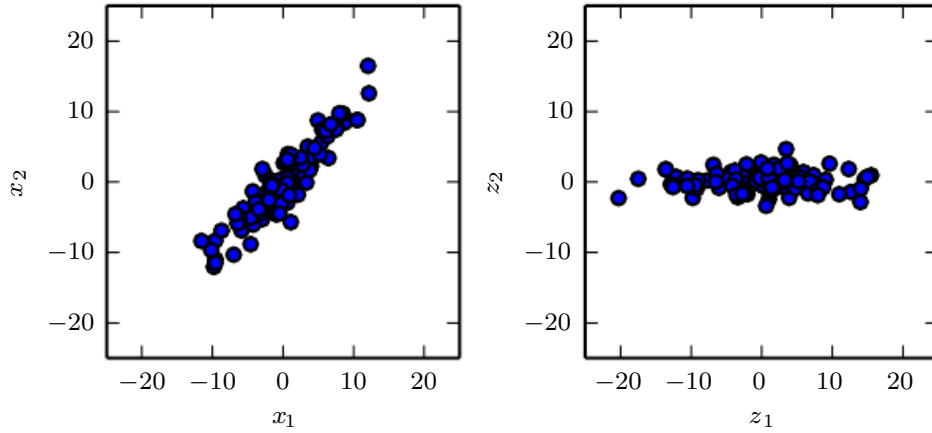


Figure 5.8: PCA learns a linear projection that aligns the direction of greatest variance with the axes of the new space. *(Left)* The original data consists of samples of  $\mathbf{x}$ . In this space, the variance might occur along directions that are not axis-aligned. *(Right)* The transformed data  $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$  now varies most along the axis  $z_1$ . The direction of second most variance is now along  $z_2$ .

representation that has lower dimensionality than the original input. It also learns a representation whose elements have no linear correlation with each other. This is a first step toward the criterion of learning representations whose elements are statistically independent. To achieve full independence, a representation learning algorithm must also remove the nonlinear relationships between variables.

PCA learns an orthogonal, linear transformation of the data that projects an input  $\mathbf{x}$  to a representation  $\mathbf{z}$  as shown in figure 5.8. In section 2.12, we saw that we could learn a one-dimensional representation that best reconstructs the original data (in the sense of mean squared error) and that this representation actually corresponds to the first principal component of the data. Thus we can use PCA as a simple and effective dimensionality reduction method that preserves as much of the information in the data as possible (again, as measured by least-squares reconstruction error). In the following, we will study how the PCA representation decorrelates the original data representation  $\mathbf{X}$ .

Let us consider the  $m \times n$ -dimensional design matrix  $\mathbf{X}$ . We will assume that the data has a mean of zero,  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ . If this is not the case, the data can easily be centered by subtracting the mean from all examples in a preprocessing step.

The unbiased sample covariance matrix associated with  $\mathbf{X}$  is given by:

$$\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^\top \mathbf{X}. \quad (5.85)$$

PCA finds a representation (through linear transformation)  $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$  where  $\text{Var}[\mathbf{z}]$  is diagonal.

In section 2.12, we saw that the principal components of a design matrix  $\mathbf{X}$  are given by the eigenvectors of  $\mathbf{X}^\top \mathbf{X}$ . From this view,

$$\mathbf{X}^\top \mathbf{X} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^\top. \quad (5.86)$$

In this section, we exploit an alternative derivation of the principal components. The principal components may also be obtained via the singular value decomposition. Specifically, they are the right singular vectors of  $\mathbf{X}$ . To see this, let  $\mathbf{W}$  be the right singular vectors in the decomposition  $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top$ . We then recover the original eigenvector equation with  $\mathbf{W}$  as the eigenvector basis:

$$\mathbf{X}^\top \mathbf{X} = (\mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top)^\top \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top = \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^\top. \quad (5.87)$$

The SVD is helpful to show that PCA results in a diagonal  $\text{Var}[\mathbf{z}]$ . Using the SVD of  $\mathbf{X}$ , we can express the variance of  $\mathbf{X}$  as:

$$\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^\top \mathbf{X} \quad (5.88)$$

$$= \frac{1}{m-1} (\mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top)^\top \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top \quad (5.89)$$

$$= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top \quad (5.90)$$

$$= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^\top, \quad (5.91)$$

where we use the fact that  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$  because the  $\mathbf{U}$  matrix of the singular value decomposition is defined to be orthogonal. This shows that if we take  $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$ , we can ensure that the covariance of  $\mathbf{z}$  is diagonal as required:

$$\text{Var}[\mathbf{z}] = \frac{1}{m-1} \mathbf{Z}^\top \mathbf{Z} \quad (5.92)$$

$$= \frac{1}{m-1} \mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} \quad (5.93)$$

$$= \frac{1}{m-1} \mathbf{W}^\top \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^\top \mathbf{W} \quad (5.94)$$

$$= \frac{1}{m-1} \mathbf{\Sigma}^2, \quad (5.95)$$

where this time we use the fact that  $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$ , again from the definition of the SVD.

The above analysis shows that when we project the data  $\mathbf{x}$  to  $\mathbf{z}$ , via the linear transformation  $\mathbf{W}$ , the resulting representation has a diagonal covariance matrix (as given by  $\Sigma^2$ ) which immediately implies that the individual elements of  $\mathbf{z}$  are mutually uncorrelated.

This ability of PCA to transform data into a representation where the elements are mutually uncorrelated is a very important property of PCA. It is a simple example of a representation that attempts to *disentangle the unknown factors of variation* underlying the data. In the case of PCA, this disentangling takes the form of finding a rotation of the input space (described by  $\mathbf{W}$ ) that aligns the principal axes of variance with the basis of the new representation space associated with  $\mathbf{z}$ .

While correlation is an important category of dependency between elements of the data, we are also interested in learning representations that disentangle more complicated forms of feature dependencies. For this, we will need more than what can be done with a simple linear transformation.

### 5.8.2 $k$ -means Clustering

Another example of a simple representation learning algorithm is  $k$ -means clustering. The  $k$ -means clustering algorithm divides the training set into  $k$  different clusters of examples that are near each other. We can thus think of the algorithm as providing a  $k$ -dimensional one-hot code vector  $\mathbf{h}$  representing an input  $\mathbf{x}$ . If  $\mathbf{x}$  belongs to cluster  $i$ , then  $h_i = 1$  and all other entries of the representation  $\mathbf{h}$  are zero.

The one-hot code provided by  $k$ -means clustering is an example of a sparse representation, because the majority of its entries are zero for every input. Later, we will develop other algorithms that learn more flexible sparse representations, where more than one entry can be non-zero for each input  $\mathbf{x}$ . One-hot codes are an extreme example of sparse representations that lose many of the benefits of a distributed representation. The one-hot code still confers some statistical advantages (it naturally conveys the idea that all examples in the same cluster are similar to each other) and it confers the computational advantage that the entire representation may be captured by a single integer.

The  $k$ -means algorithm works by initializing  $k$  different centroids  $\{\boldsymbol{\mu}^{(1)}, \dots, \boldsymbol{\mu}^{(k)}\}$  to different values, then alternating between two different steps until convergence. In one step, each training example is assigned to cluster  $i$ , where  $i$  is the index of the nearest centroid  $\boldsymbol{\mu}^{(i)}$ . In the other step, each centroid  $\boldsymbol{\mu}^{(i)}$  is updated to the mean of all training examples  $\mathbf{x}^{(j)}$  assigned to cluster  $i$ .

One difficulty pertaining to clustering is that the clustering problem is inherently ill-posed, in the sense that there is no single criterion that measures how well a clustering of the data corresponds to the real world. We can measure properties of the clustering such as the average Euclidean distance from a cluster centroid to the members of the cluster. This allows us to tell how well we are able to reconstruct the training data from the cluster assignments. We do not know how well the cluster assignments correspond to properties of the real world. Moreover, there may be many different clusterings that all correspond well to some property of the real world. We may hope to find a clustering that relates to one feature but obtain a different, equally valid clustering that is not relevant to our task. For example, suppose that we run two clustering algorithms on a dataset consisting of images of red trucks, images of red cars, images of gray trucks, and images of gray cars. If we ask each clustering algorithm to find two clusters, one algorithm may find a cluster of cars and a cluster of trucks, while another may find a cluster of red vehicles and a cluster of gray vehicles. Suppose we also run a third clustering algorithm, which is allowed to determine the number of clusters. This may assign the examples to four clusters, red cars, red trucks, gray cars, and gray trucks. This new clustering now at least captures information about both attributes, but it has lost information about similarity. Red cars are in a different cluster from gray cars, just as they are in a different cluster from gray trucks. The output of the clustering algorithm does not tell us that red cars are more similar to gray cars than they are to gray trucks. They are different from both things, and that is all we know.

These issues illustrate some of the reasons that we may prefer a distributed representation to a one-hot representation. A distributed representation could have two attributes for each vehicle—one representing its color and one representing whether it is a car or a truck. It is still not entirely clear what the optimal distributed representation is (how can the learning algorithm know whether the two attributes we are interested in are color and car-versus-truck rather than manufacturer and age?) but having many attributes reduces the burden on the algorithm to guess which single attribute we care about, and allows us to measure similarity between objects in a fine-grained way by comparing many attributes instead of just testing whether one attribute matches.

## 5.9 Stochastic Gradient Descent

Nearly all of deep learning is powered by one very important algorithm: **stochastic gradient descent** or SGD. Stochastic gradient descent is an extension of the

gradient descent algorithm introduced in section 4.3.

A recurring problem in machine learning is that large training sets are necessary for good generalization, but large training sets are also more computationally expensive.

The cost function used by a machine learning algorithm often decomposes as a sum over training examples of some per-example loss function. For example, the negative conditional log-likelihood of the training data can be written as

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} L(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}) \quad (5.96)$$

where  $L$  is the per-example loss  $L(\mathbf{x}, y, \boldsymbol{\theta}) = -\log p(y \mid \mathbf{x}; \boldsymbol{\theta})$ .

For these additive cost functions, gradient descent requires computing

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}). \quad (5.97)$$

The computational cost of this operation is  $O(m)$ . As the training set size grows to billions of examples, the time to take a single gradient step becomes prohibitively long.

The insight of stochastic gradient descent is that the gradient is an expectation. The expectation may be approximately estimated using a small set of samples. Specifically, on each step of the algorithm, we can sample a **minibatch** of examples  $\mathbb{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$  drawn uniformly from the training set. The minibatch size  $m'$  is typically chosen to be a relatively small number of examples, ranging from 1 to a few hundred. Crucially,  $m'$  is usually held fixed as the training set size  $m$  grows. We may fit a training set with billions of examples using updates computed on only a hundred examples.

The estimate of the gradient is formed as

$$\mathbf{g} = \frac{1}{m'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}). \quad (5.98)$$

using examples from the minibatch  $\mathbb{B}$ . The stochastic gradient descent algorithm then follows the estimated gradient downhill:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g}, \quad (5.99)$$

where  $\epsilon$  is the learning rate.

Gradient descent in general has often been regarded as slow or unreliable. In the past, the application of gradient descent to non-convex optimization problems was regarded as foolhardy or unprincipled. Today, we know that the machine learning models described in part II work very well when trained with gradient descent. The optimization algorithm may not be guaranteed to arrive at even a local minimum in a reasonable amount of time, but it often finds a very low value of the cost function quickly enough to be useful.

Stochastic gradient descent has many important uses outside the context of deep learning. It is the main way to train large linear models on very large datasets. For a fixed model size, the cost per SGD update does not depend on the training set size  $m$ . In practice, we often use a larger model as the training set size increases, but we are not forced to do so. The number of updates required to reach convergence usually increases with training set size. However, as  $m$  approaches infinity, the model will eventually converge to its best possible test error before SGD has sampled every example in the training set. Increasing  $m$  further will not extend the amount of training time needed to reach the model's best possible test error. From this point of view, one can argue that the asymptotic cost of training a model with SGD is  $O(1)$  as a function of  $m$ .

Prior to the advent of deep learning, the main way to learn nonlinear models was to use the kernel trick in combination with a linear model. Many kernel learning algorithms require constructing an  $m \times m$  matrix  $G_{i,j} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ . Constructing this matrix has computational cost  $O(m^2)$ , which is clearly undesirable for datasets with billions of examples. In academia, starting in 2006, deep learning was initially interesting because it was able to generalize to new examples better than competing algorithms when trained on medium-sized datasets with tens of thousands of examples. Soon after, deep learning garnered additional interest in industry, because it provided a scalable way of training nonlinear models on large datasets.

Stochastic gradient descent and many enhancements to it are described further in chapter 8.

## 5.10 Building a Machine Learning Algorithm

Nearly all deep learning algorithms can be described as particular instances of a fairly simple recipe: combine a specification of a dataset, a cost function, an optimization procedure and a model.

For example, the linear regression algorithm combines a dataset consisting of



$\mathbf{X}$  and  $\mathbf{y}$ , the cost function

$$J(\mathbf{w}, b) = -\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y \mid \mathbf{x}), \quad (5.100)$$

the model specification  $p_{\text{model}}(y \mid \mathbf{x}) = \mathcal{N}(y; \mathbf{x}^\top \mathbf{w} + b, 1)$ , and, in most cases, the optimization algorithm defined by solving for where the gradient of the cost is zero using the normal equations.

By realizing that we can replace any of these components mostly independently from the others, we can obtain a very wide variety of algorithms.

The cost function typically includes at least one term that causes the learning process to perform statistical estimation. The most common cost function is the negative log-likelihood, so that minimizing the cost function causes maximum likelihood estimation.

The cost function may also include additional terms, such as regularization terms. For example, we can add weight decay to the linear regression cost function to obtain

$$J(\mathbf{w}, b) = \lambda \|\mathbf{w}\|_2^2 - \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y \mid \mathbf{x}). \quad (5.101)$$

This still allows closed-form optimization.

If we change the model to be nonlinear, then most cost functions can no longer be optimized in closed form. This requires us to choose an iterative numerical optimization procedure, such as gradient descent.

The recipe for constructing a learning algorithm by combining models, costs, and optimization algorithms supports both supervised and unsupervised learning. The linear regression example shows how to support supervised learning. Unsupervised learning can be supported by defining a dataset that contains only  $\mathbf{X}$  and providing an appropriate unsupervised cost and model. For example, we can obtain the first PCA vector by specifying that our loss function is

$$J(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \|\mathbf{x} - r(\mathbf{x}; \mathbf{w})\|_2^2 \quad (5.102)$$

while our model is defined to have  $\mathbf{w}$  with norm one and reconstruction function  $r(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \mathbf{w}$ .

In some cases, the cost function may be a function that we cannot actually evaluate, for computational reasons. In these cases, we can still approximately minimize it using iterative numerical optimization so long as we have some way of approximating its gradients.

Most machine learning algorithms make use of this recipe, though it may not immediately be obvious. If a machine learning algorithm seems especially unique or

hand-designed, it can usually be understood as using a special-case optimizer. Some models such as decision trees or  $k$ -means require special-case optimizers because their cost functions have flat regions that make them inappropriate for minimization by gradient-based optimizers. Recognizing that most machine learning algorithms can be described using this recipe helps to see the different algorithms as part of a taxonomy of methods for doing related tasks that work for similar reasons, rather than as a long list of algorithms that each have separate justifications.

## 5.11 Challenges Motivating Deep Learning

The simple machine learning algorithms described in this chapter work very well on a wide variety of important problems. However, they have not succeeded in solving the central problems in AI, such as recognizing speech or recognizing objects.

The development of deep learning was motivated in part by the failure of traditional algorithms to generalize well on such AI tasks.

This section is about how the challenge of generalizing to new examples becomes exponentially more difficult when working with high-dimensional data, and how the mechanisms used to achieve generalization in traditional machine learning are insufficient to learn complicated functions in high-dimensional spaces. Such spaces also often impose high computational costs. Deep learning was designed to overcome these and other obstacles.

### 5.11.1 The Curse of Dimensionality

Many machine learning problems become exceedingly difficult when the number of dimensions in the data is high. This phenomenon is known as the **curse of dimensionality**. Of particular concern is that the number of possible distinct configurations of a set of variables increases exponentially as the number of variables increases.

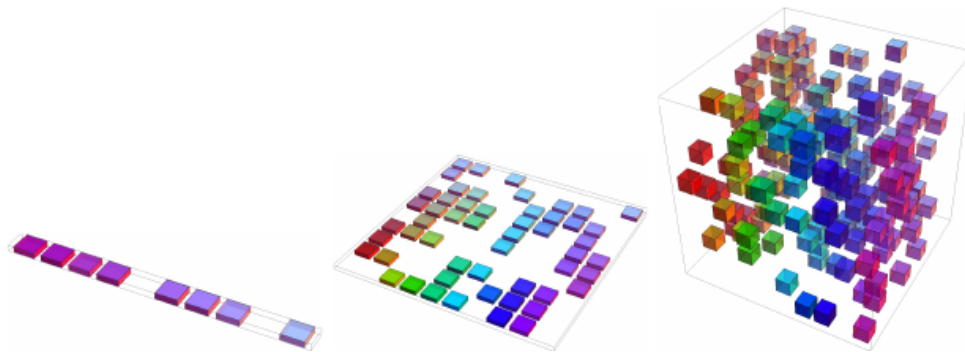


Figure 5.9: As the number of relevant dimensions of the data increases (from left to right), the number of configurations of interest may grow exponentially. *(Left)* In this one-dimensional example, we have one variable for which we only care to distinguish 10 regions of interest. With enough examples falling within each of these regions (each region corresponds to a cell in the illustration), learning algorithms can easily generalize correctly. A straightforward way to generalize is to estimate the value of the target function within each region (and possibly interpolate between neighboring regions). *(Center)* With 2 dimensions it is more difficult to distinguish 10 different values of each variable. We need to keep track of up to  $10 \times 10 = 100$  regions, and we need at least that many examples to cover all those regions. *(Right)* With 3 dimensions this grows to  $10^3 = 1000$  regions and at least that many examples. For  $d$  dimensions and  $v$  values to be distinguished along each axis, we seem to need  $O(v^d)$  regions and examples. This is an instance of the curse of dimensionality. Figure graciously provided by Nicolas Chapados.

The curse of dimensionality arises in many places in computer science, and especially so in machine learning.

One challenge posed by the curse of dimensionality is a statistical challenge. As illustrated in figure 5.9, a statistical challenge arises because the number of possible configurations of  $\mathbf{x}$  is much larger than the number of training examples. To understand the issue, let us consider that the input space is organized into a grid, like in the figure. We can describe low-dimensional space with a low number of grid cells that are mostly occupied by the data. When generalizing to a new data point, we can usually tell what to do simply by inspecting the training examples that lie in the same cell as the new input. For example, if estimating the probability density at some point  $\mathbf{x}$ , we can just return the number of training examples in the same unit volume cell as  $\mathbf{x}$ , divided by the total number of training examples. If we wish to classify an example, we can return the most common class of training examples in the same cell. If we are doing regression we can average the target values observed over the examples in that cell. But what about the cells for which we have seen no example? Because in high-dimensional spaces the number of configurations is huge, much larger than our number of examples, a typical grid cell has no training example associated with it. How could we possibly say something

meaningful about these new configurations? Many traditional machine learning algorithms simply assume that the output at a new point should be approximately the same as the output at the nearest training point.

### 5.11.2 Local Constancy and Smoothness Regularization

In order to generalize well, machine learning algorithms need to be guided by prior beliefs about what kind of function they should learn. Previously, we have seen these priors incorporated as explicit beliefs in the form of probability distributions over parameters of the model. More informally, we may also discuss prior beliefs as directly influencing the *function* itself and only indirectly acting on the parameters via their effect on the function. Additionally, we informally discuss prior beliefs as being expressed implicitly, by choosing algorithms that are biased toward choosing some class of functions over another, even though these biases may not be expressed (or even possible to express) in terms of a probability distribution representing our degree of belief in various functions.

Among the most widely used of these implicit “priors” is the **smoothness prior** or **local constancy prior**. This prior states that the function we learn should not change very much within a small region.

Many simpler algorithms rely exclusively on this prior to generalize well, and as a result they fail to scale to the statistical challenges involved in solving AI-level tasks. Throughout this book, we will describe how deep learning introduces additional (explicit and implicit) priors in order to reduce the generalization error on sophisticated tasks. Here, we explain why the smoothness prior alone is insufficient for these tasks.

There are many different ways to implicitly or explicitly express a prior belief that the learned function should be smooth or locally constant. All of these different methods are designed to encourage the learning process to learn a function  $f^*$  that satisfies the condition

$$f^*(\mathbf{x}) \approx f^*(\mathbf{x} + \epsilon) \quad (5.103)$$

for most configurations  $\mathbf{x}$  and small change  $\epsilon$ . In other words, if we know a good answer for an input  $\mathbf{x}$  (for example, if  $\mathbf{x}$  is a labeled training example) then that answer is probably good in the neighborhood of  $\mathbf{x}$ . If we have several good answers in some neighborhood we would combine them (by some form of averaging or interpolation) to produce an answer that agrees with as many of them as much as possible.

An extreme example of the local constancy approach is the  $k$ -nearest neighbors family of learning algorithms. These predictors are literally constant over each

region containing all the points  $\mathbf{x}$  that have the same set of  $k$  nearest neighbors in the training set. For  $k = 1$ , the number of distinguishable regions cannot be more than the number of training examples.

While the  $k$ -nearest neighbors algorithm copies the output from nearby training examples, most kernel machines interpolate between training set outputs associated with nearby training examples. An important class of kernels is the family of **local kernels** where  $k(\mathbf{u}, \mathbf{v})$  is large when  $\mathbf{u} = \mathbf{v}$  and decreases as  $\mathbf{u}$  and  $\mathbf{v}$  grow farther apart from each other. A local kernel can be thought of as a similarity function that performs template matching, by measuring how closely a test example  $\mathbf{x}$  resembles each training example  $\mathbf{x}^{(i)}$ . Much of the modern motivation for deep learning is derived from studying the limitations of local template matching and how deep models are able to succeed in cases where local template matching fails (Bengio *et al.*, 2006b).

Decision trees also suffer from the limitations of exclusively smoothness-based learning because they break the input space into as many regions as there are leaves and use a separate parameter (or sometimes many parameters for extensions of decision trees) in each region. If the target function requires a tree with at least  $n$  leaves to be represented accurately, then at least  $n$  training examples are required to fit the tree. A multiple of  $n$  is needed to achieve some level of statistical confidence in the predicted output.

In general, to distinguish  $O(k)$  regions in input space, all of these methods require  $O(k)$  examples. Typically there are  $O(k)$  parameters, with  $O(1)$  parameters associated with each of the  $O(k)$  regions. The case of a nearest neighbor scenario, where each training example can be used to define at most one region, is illustrated in figure 5.10.

Is there a way to represent a complex function that has many more regions to be distinguished than the number of training examples? Clearly, assuming only smoothness of the underlying function will not allow a learner to do that. For example, imagine that the target function is a kind of checkerboard. A checkerboard contains many variations but there is a simple structure to them. Imagine what happens when the number of training examples is substantially smaller than the number of black and white squares on the checkerboard. Based on only local generalization and the smoothness or local constancy prior, we would be guaranteed to correctly guess the color of a new point if it lies within the same checkerboard square as a training example. There is no guarantee that the learner could correctly extend the checkerboard pattern to points lying in squares that do not contain training examples. With this prior alone, the only information that an example tells us is the color of its square, and the only way to get the colors of the

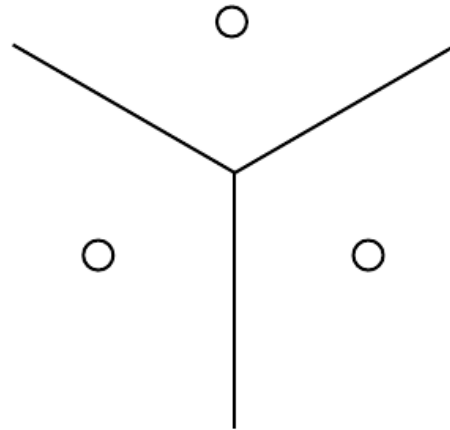


Figure 5.10: Illustration of how the nearest neighbor algorithm breaks up the input space into regions. An example (represented here by a circle) within each region defines the region boundary (represented here by the lines). The  $y$  value associated with each example defines what the output should be for all points within the corresponding region. The regions defined by nearest neighbor matching form a geometric pattern called a Voronoi diagram. The number of these contiguous regions cannot grow faster than the number of training examples. While this figure illustrates the behavior of the nearest neighbor algorithm specifically, other machine learning algorithms that rely exclusively on the local smoothness prior for generalization exhibit similar behaviors: each training example only informs the learner about how to generalize in some neighborhood immediately surrounding that example.

entire checkerboard right is to cover each of its cells with at least one example.

The smoothness assumption and the associated non-parametric learning algorithms work extremely well so long as there are enough examples for the learning algorithm to observe high points on most peaks and low points on most valleys of the true underlying function to be learned. This is generally true when the function to be learned is smooth enough and varies in few enough dimensions. In high dimensions, even a very smooth function can change smoothly but in a different way along each dimension. If the function additionally behaves differently in different regions, it can become extremely complicated to describe with a set of training examples. If the function is complicated (we want to distinguish a huge number of regions compared to the number of examples), is there any hope to generalize well?

The answer to both of these questions—whether it is possible to represent a complicated function efficiently, and whether it is possible for the estimated function to generalize well to new inputs—is yes. The key insight is that a very large number of regions, e.g.,  $O(2^k)$ , can be defined with  $O(k)$  examples, so long as we introduce some dependencies between the regions via additional assumptions about the underlying data generating distribution. In this way, we can actually generalize non-locally (Bengio and Monperrus, 2005; Bengio *et al.*, 2006c). Many different deep learning algorithms provide implicit or explicit assumptions that are reasonable for a broad range of AI tasks in order to capture these advantages.

Other approaches to machine learning often make stronger, task-specific assumptions. For example, we could easily solve the checkerboard task by providing the assumption that the target function is periodic. Usually we do not include such strong, task-specific assumptions into neural networks so that they can generalize to a much wider variety of structures. AI tasks have structure that is much too complex to be limited to simple, manually specified properties such as periodicity, so we want learning algorithms that embody more general-purpose assumptions. The core idea in deep learning is that we assume that the data was generated by the *composition of factors* or features, potentially at multiple levels in a hierarchy. Many other similarly generic assumptions can further improve deep learning algorithms. These apparently mild assumptions allow an exponential gain in the relationship between the number of examples and the number of regions that can be distinguished. These exponential gains are described more precisely in sections 6.4.1, 15.4 and 15.5. The exponential advantages conferred by the use of deep, distributed representations counter the exponential challenges posed by the curse of dimensionality.

### 5.11.3 Manifold Learning

An important concept underlying many ideas in machine learning is that of a manifold.

A **manifold** is a connected region. Mathematically, it is a set of points, associated with a neighborhood around each point. From any given point, the manifold locally appears to be a Euclidean space. In everyday life, we experience the surface of the world as a 2-D plane, but it is in fact a spherical manifold in 3-D space.

The definition of a neighborhood surrounding each point implies the existence of transformations that can be applied to move on the manifold from one position to a neighboring one. In the example of the world’s surface as a manifold, one can walk north, south, east, or west.

Although there is a formal mathematical meaning to the term “manifold,” in machine learning it tends to be used more loosely to designate a connected set of points that can be approximated well by considering only a small number of degrees of freedom, or dimensions, embedded in a higher-dimensional space. Each dimension corresponds to a local direction of variation. See figure 5.11 for an example of training data lying near a one-dimensional manifold embedded in two-dimensional space. In the context of machine learning, we allow the dimensionality of the manifold to vary from one point to another. This often happens when a manifold intersects itself. For example, a figure eight is a manifold that has a single dimension in most places but two dimensions at the intersection at the center.

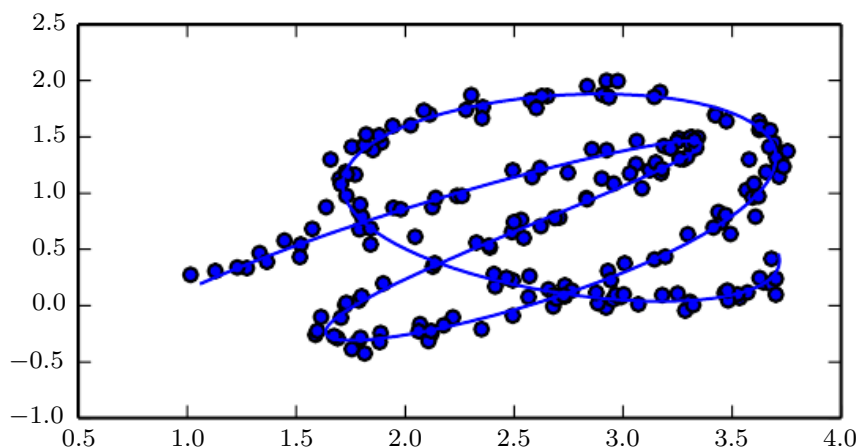


Figure 5.11: Data sampled from a distribution in a two-dimensional space that is actually concentrated near a one-dimensional manifold, like a twisted string. The solid line indicates the underlying manifold that the learner should infer.



Many machine learning problems seem hopeless if we expect the machine learning algorithm to learn functions with interesting variations across all of  $\mathbb{R}^n$ . **Manifold learning** algorithms surmount this obstacle by assuming that most of  $\mathbb{R}^n$  consists of invalid inputs, and that interesting inputs occur only along a collection of manifolds containing a small subset of points, with interesting variations in the output of the learned function occurring only along directions that lie on the manifold, or with interesting variations happening only when we move from one manifold to another. Manifold learning was introduced in the case of continuous-valued data and the unsupervised learning setting, although this probability concentration idea can be generalized to both discrete data and the supervised learning setting: the key assumption remains that probability mass is highly concentrated.

The assumption that the data lies along a low-dimensional manifold may not always be correct or useful. We argue that in the context of AI tasks, such as those that involve processing images, sounds, or text, the manifold assumption is at least approximately correct. The evidence in favor of this assumption consists of two categories of observations.

The first observation in favor of the **manifold hypothesis** is that the probability distribution over images, text strings, and sounds that occur in real life is highly concentrated. Uniform noise essentially never resembles structured inputs from these domains. Figure 5.12 shows how, instead, uniformly sampled points look like the patterns of static that appear on analog television sets when no signal is available. Similarly, if you generate a document by picking letters uniformly at random, what is the probability that you will get a meaningful English-language text? Almost zero, again, because most of the long sequences of letters do not correspond to a natural language sequence: the distribution of natural language sequences occupies a very small volume in the total space of sequences of letters.

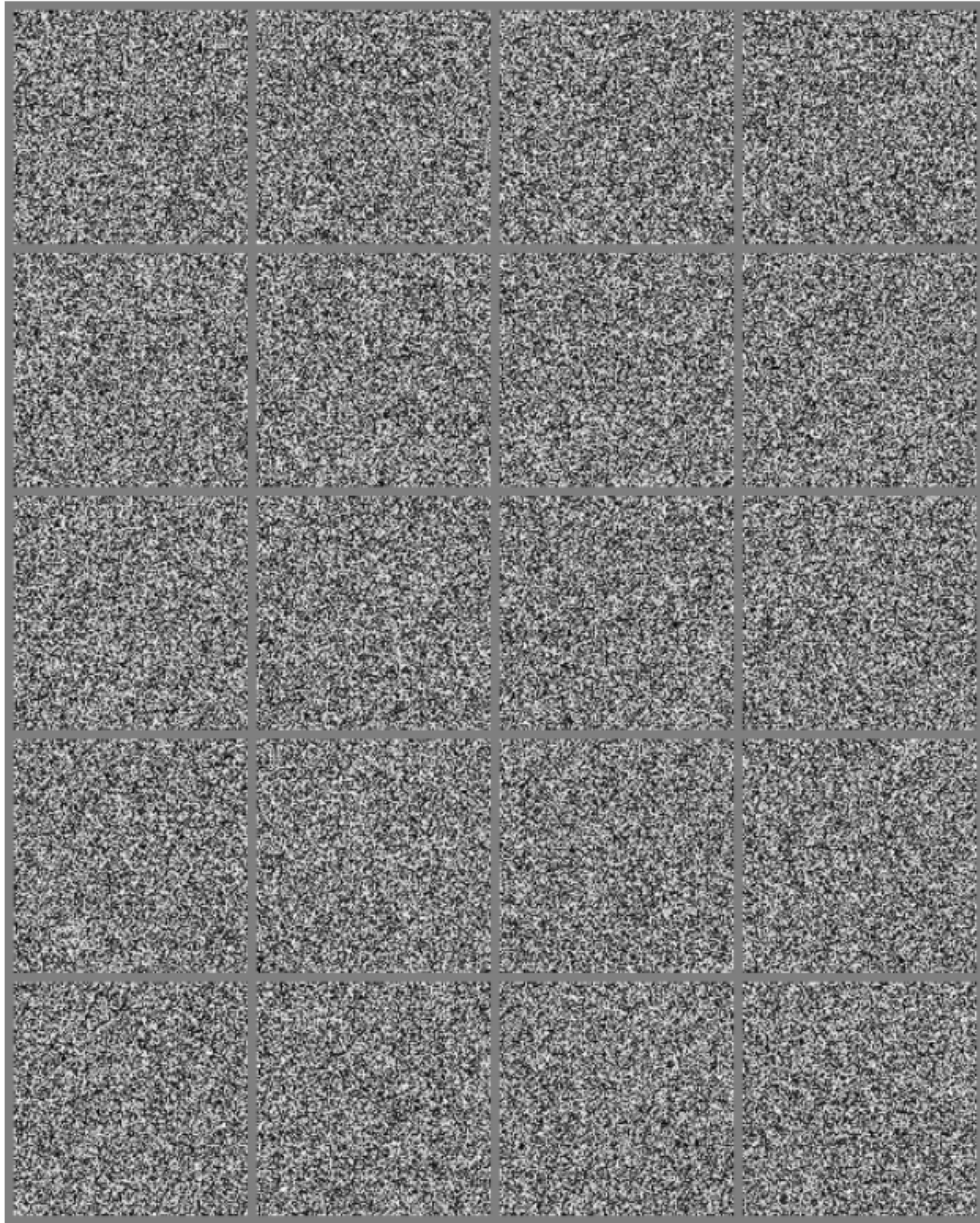


Figure 5.12: Sampling images uniformly at random (by randomly picking each pixel according to a uniform distribution) gives rise to noisy images. Although there is a non-zero probability to generate an image of a face or any other object frequently encountered in AI applications, we never actually observe this happening in practice. This suggests that the images encountered in AI applications occupy a negligible proportion of the volume of image space.

Of course, concentrated probability distributions are not sufficient to show that the data lies on a reasonably small number of manifolds. We must also establish that the examples we encounter are connected to each other by other

examples, with each example surrounded by other highly similar examples that may be reached by applying transformations to traverse the manifold. The second argument in favor of the manifold hypothesis is that we can also imagine such neighborhoods and transformations, at least informally. In the case of images, we can certainly think of many possible transformations that allow us to trace out a manifold in image space: we can gradually dim or brighten the lights, gradually move or rotate objects in the image, gradually alter the colors on the surfaces of objects, etc. It remains likely that there are multiple manifolds involved in most applications. For example, the manifold of images of human faces may not be connected to the manifold of images of cat faces.

These thought experiments supporting the manifold hypotheses convey some intuitive reasons supporting it. More rigorous experiments (Cayton, 2005; Narayanan and Mitter, 2010; Schölkopf *et al.*, 1998; Roweis and Saul, 2000; Tenenbaum *et al.*, 2000; Brand, 2003; Belkin and Niyogi, 2003; Donoho and Grimes, 2003; Weinberger and Saul, 2004) clearly support the hypothesis for a large class of datasets of interest in AI.

When the data lies on a low-dimensional manifold, it can be most natural for machine learning algorithms to represent the data in terms of coordinates on the manifold, rather than in terms of coordinates in  $\mathbb{R}^n$ . In everyday life, we can think of roads as 1-D manifolds embedded in 3-D space. We give directions to specific addresses in terms of address numbers along these 1-D roads, not in terms of coordinates in 3-D space. Extracting these manifold coordinates is challenging, but holds the promise to improve many machine learning algorithms. This general principle is applied in many contexts. Figure 5.13 shows the manifold structure of a dataset consisting of faces. By the end of this book, we will have developed the methods necessary to learn such a manifold structure. In figure 20.6, we will see how a machine learning algorithm can successfully accomplish this goal.

This concludes part I, which has provided the basic concepts in mathematics and machine learning which are employed throughout the remaining parts of the book. You are now prepared to embark upon your study of deep learning.



Figure 5.13: Training examples from the QMUL Multiview Face Dataset ([Gong \*et al.\*, 2000](#)) for which the subjects were asked to move in such a way as to cover the two-dimensional manifold corresponding to two angles of rotation. We would like learning algorithms to be able to discover and disentangle such manifold coordinates. Figure [20.6](#) illustrates such a feat.

## Part II

# Deep Networks: Modern Practices

---

This part of the book summarizes the state of modern deep learning as it is used to solve practical applications.

Deep learning has a long history and many aspirations. Several approaches have been proposed that have yet to entirely bear fruit. Several ambitious goals have yet to be realized. These less-developed branches of deep learning appear in the final part of the book.

This part focuses only on those approaches that are essentially working technologies that are already used heavily in industry.

Modern deep learning provides a very powerful framework for supervised learning. By adding more layers and more units within a layer, a deep network can represent functions of increasing complexity. Most tasks that consist of mapping an input vector to an output vector, and that are easy for a person to do rapidly, can be accomplished via deep learning, given sufficiently large models and sufficiently large datasets of labeled training examples. Other tasks, that can not be described as associating one vector to another, or that are difficult enough that a person would require time to think and reflect in order to accomplish the task, remain beyond the scope of deep learning for now.

This part of the book describes the core parametric function approximation technology that is behind nearly all modern practical applications of deep learning. We begin by describing the feedforward deep network model that is used to represent these functions. Next, we present advanced techniques for regularization and optimization of such models. Scaling these models to large inputs such as high resolution images or long temporal sequences requires specialization. We introduce the convolutional network for scaling to large images and the recurrent neural network for processing temporal sequences. Finally, we present general guidelines for the practical methodology involved in designing, building, and configuring an application involving deep learning, and review some of the applications of deep learning.

These chapters are the most important for a practitioner—someone who wants to begin implementing and using deep learning algorithms to solve real-world problems today.

## Chapter 6

# Deep Feedforward Networks

**Deep feedforward networks**, also often called **feedforward neural networks**, or **multilayer perceptrons** (MLPs), are the quintessential deep learning models. The goal of a feedforward network is to approximate some function  $f^*$ . For example, for a classifier,  $y = f^*(\mathbf{x})$  maps an input  $\mathbf{x}$  to a category  $y$ . A feedforward network defines a mapping  $\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta})$  and learns the value of the parameters  $\boldsymbol{\theta}$  that result in the best function approximation.

These models are called **feedforward** because information flows through the function being evaluated from  $\mathbf{x}$ , through the intermediate computations used to define  $f$ , and finally to the output  $\mathbf{y}$ . There are no **feedback** connections in which outputs of the model are fed back into itself. When feedforward neural networks are extended to include feedback connections, they are called **recurrent neural networks**, presented in chapter 10.

Feedforward networks are of extreme importance to machine learning practitioners. They form the basis of many important commercial applications. For example, the convolutional networks used for object recognition from photos are a specialized kind of feedforward network. Feedforward networks are a conceptual stepping stone on the path to recurrent networks, which power many natural language applications.

Feedforward neural networks are called **networks** because they are typically represented by composing together many different functions. The model is associated with a directed acyclic graph describing how the functions are composed together. For example, we might have three functions  $f^{(1)}$ ,  $f^{(2)}$ , and  $f^{(3)}$  connected in a chain, to form  $f(\mathbf{x}) = f^{(3)}(f^{(2)}(f^{(1)}(\mathbf{x})))$ . These chain structures are the most commonly used structures of neural networks. In this case,  $f^{(1)}$  is called the **first layer** of the network,  $f^{(2)}$  is called the **second layer**, and so on. The overall

length of the chain gives the **depth** of the model. It is from this terminology that the name “deep learning” arises. The final layer of a feedforward network is called the **output layer**. During neural network training, we drive  $f(\mathbf{x})$  to match  $f^*(\mathbf{x})$ . The training data provides us with noisy, approximate examples of  $f^*(\mathbf{x})$  evaluated at different training points. Each example  $\mathbf{x}$  is accompanied by a label  $y \approx f^*(\mathbf{x})$ . The training examples specify directly what the output layer must do at each point  $\mathbf{x}$ ; it must produce a value that is close to  $y$ . The behavior of the other layers is not directly specified by the training data. The learning algorithm must decide how to use those layers to produce the desired output, but the training data does not say what each individual layer should do. Instead, the learning algorithm must decide how to use these layers to best implement an approximation of  $f^*$ . Because the training data does not show the desired output for each of these layers, these layers are called **hidden layers**.

Finally, these networks are called *neural* because they are loosely inspired by neuroscience. Each hidden layer of the network is typically vector-valued. The dimensionality of these hidden layers determines the **width** of the model. Each element of the vector may be interpreted as playing a role analogous to a neuron. Rather than thinking of the layer as representing a single vector-to-vector function, we can also think of the layer as consisting of many **units** that act in parallel, each representing a vector-to-scalar function. Each unit resembles a neuron in the sense that it receives input from many other units and computes its own activation value. The idea of using many layers of vector-valued representation is drawn from neuroscience. The choice of the functions  $f^{(i)}(\mathbf{x})$  used to compute these representations is also loosely guided by neuroscientific observations about the functions that biological neurons compute. However, modern neural network research is guided by many mathematical and engineering disciplines, and the goal of neural networks is not to perfectly model the brain. It is best to think of feedforward networks as function approximation machines that are designed to achieve statistical generalization, occasionally drawing some insights from what we know about the brain, rather than as models of brain function.

One way to understand feedforward networks is to begin with linear models and consider how to overcome their limitations. Linear models, such as logistic regression and linear regression, are appealing because they may be fit efficiently and reliably, either in closed form or with convex optimization. Linear models also have the obvious defect that the model capacity is limited to linear functions, so the model cannot understand the interaction between any two input variables.

To extend linear models to represent nonlinear functions of  $\mathbf{x}$ , we can apply the linear model not to  $\mathbf{x}$  itself but to a transformed input  $\phi(\mathbf{x})$ , where  $\phi$  is a



nonlinear transformation. Equivalently, we can apply the kernel trick described in section 5.7.2, to obtain a nonlinear learning algorithm based on implicitly applying the  $\phi$  mapping. We can think of  $\phi$  as providing a set of features describing  $\mathbf{x}$ , or as providing a new representation for  $\mathbf{x}$ .

The question is then how to choose the mapping  $\phi$ .

1. One option is to use a very generic  $\phi$ , such as the infinite-dimensional  $\phi$  that is implicitly used by kernel machines based on the RBF kernel. If  $\phi(\mathbf{x})$  is of high enough dimension, we can always have enough capacity to fit the training set, but generalization to the test set often remains poor. Very generic feature mappings are usually based only on the principle of local smoothness and do not encode enough prior information to solve advanced problems.
2. Another option is to manually engineer  $\phi$ . Until the advent of deep learning, this was the dominant approach. This approach requires decades of human effort for each separate task, with practitioners specializing in different domains such as speech recognition or computer vision, and with little transfer between domains.
3. The strategy of deep learning is to learn  $\phi$ . In this approach, we have a model  $y = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \phi(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{w}$ . We now have parameters  $\boldsymbol{\theta}$  that we use to learn  $\phi$  from a broad class of functions, and parameters  $\mathbf{w}$  that map from  $\phi(\mathbf{x})$  to the desired output. This is an example of a deep feedforward network, with  $\phi$  defining a hidden layer. This approach is the only one of the three that gives up on the convexity of the training problem, but the benefits outweigh the harms. In this approach, we parametrize the representation as  $\phi(\mathbf{x}; \boldsymbol{\theta})$  and use the optimization algorithm to find the  $\boldsymbol{\theta}$  that corresponds to a good representation. If we wish, this approach can capture the benefit of the first approach by being highly generic—we do so by using a very broad family  $\phi(\mathbf{x}; \boldsymbol{\theta})$ . This approach can also capture the benefit of the second approach. Human practitioners can encode their knowledge to help generalization by designing families  $\phi(\mathbf{x}; \boldsymbol{\theta})$  that they expect will perform well. The advantage is that the human designer only needs to find the right general function family rather than finding precisely the right function.

This general principle of improving models by learning features extends beyond the feedforward networks described in this chapter. It is a recurring theme of deep learning that applies to all of the kinds of models described throughout this book. Feedforward networks are the application of this principle to learning deterministic

mappings from  $\mathbf{x}$  to  $\mathbf{y}$  that lack feedback connections. Other models presented later will apply these principles to learning stochastic mappings, learning functions with feedback, and learning probability distributions over a single vector.

We begin this chapter with a simple example of a feedforward network. Next, we address each of the design decisions needed to deploy a feedforward network. First, training a feedforward network requires making many of the same design decisions as are necessary for a linear model: choosing the optimizer, the cost function, and the form of the output units. We review these basics of gradient-based learning, then proceed to confront some of the design decisions that are unique to feedforward networks. Feedforward networks have introduced the concept of a hidden layer, and this requires us to choose the **activation functions** that will be used to compute the hidden layer values. We must also design the architecture of the network, including how many layers the network should contain, how these layers should be connected to each other, and how many units should be in each layer. Learning in deep neural networks requires computing the gradients of complicated functions. We present the **back-propagation** algorithm and its modern generalizations, which can be used to efficiently compute these gradients. Finally, we close with some historical perspective.

## 6.1 Example: Learning XOR

To make the idea of a feedforward network more concrete, we begin with an example of a fully functioning feedforward network on a very simple task: learning the XOR function.

The XOR function (“exclusive or”) is an operation on two binary values,  $x_1$  and  $x_2$ . When exactly one of these binary values is equal to 1, the XOR function returns 1. Otherwise, it returns 0. The XOR function provides the target function  $y = f^*(\mathbf{x})$  that we want to learn. Our model provides a function  $y = f(\mathbf{x}; \boldsymbol{\theta})$  and our learning algorithm will adapt the parameters  $\boldsymbol{\theta}$  to make  $f$  as similar as possible to  $f^*$ .

In this simple example, we will not be concerned with statistical generalization. We want our network to perform correctly on the four points  $\mathbb{X} = \{[0, 0]^\top, [0, 1]^\top, [1, 0]^\top, \text{ and } [1, 1]^\top\}$ . We will train the network on all four of these points. The only challenge is to fit the training set.

We can treat this problem as a regression problem and use a mean squared error loss function. We choose this loss function to simplify the math for this example as much as possible. In practical applications, MSE is usually not an

appropriate cost function for modeling binary data. More appropriate approaches are described in section 6.2.2.2.

Evaluated on our whole training set, the MSE loss function is

$$J(\boldsymbol{\theta}) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \boldsymbol{\theta}))^2. \quad (6.1)$$

Now we must choose the form of our model,  $f(\mathbf{x}; \boldsymbol{\theta})$ . Suppose that we choose a linear model, with  $\boldsymbol{\theta}$  consisting of  $\mathbf{w}$  and  $b$ . Our model is defined to be

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{x}^\top \mathbf{w} + b. \quad (6.2)$$

We can minimize  $J(\boldsymbol{\theta})$  in closed form with respect to  $\mathbf{w}$  and  $b$  using the normal equations.

After solving the normal equations, we obtain  $\mathbf{w} = \mathbf{0}$  and  $b = \frac{1}{2}$ . The linear model simply outputs 0.5 everywhere. Why does this happen? Figure 6.1 shows how a linear model is not able to represent the XOR function. One way to solve this problem is to use a model that learns a different feature space in which a linear model is able to represent the solution.

Specifically, we will introduce a very simple feedforward network with one hidden layer containing two hidden units. See figure 6.2 for an illustration of this model. This feedforward network has a vector of hidden units  $\mathbf{h}$  that are computed by a function  $f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c})$ . The values of these hidden units are then used as the input for a second layer. The second layer is the output layer of the network. The output layer is still just a linear regression model, but now it is applied to  $\mathbf{h}$  rather than to  $\mathbf{x}$ . The network now contains two functions chained together:  $\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c})$  and  $y = f^{(2)}(\mathbf{h}; \mathbf{w}, b)$ , with the complete model being  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = f^{(2)}(f^{(1)}(\mathbf{x}))$ .

What function should  $f^{(1)}$  compute? Linear models have served us well so far, and it may be tempting to make  $f^{(1)}$  be linear as well. Unfortunately, if  $f^{(1)}$  were linear, then the feedforward network as a whole would remain a linear function of its input. Ignoring the intercept terms for the moment, suppose  $f^{(1)}(\mathbf{x}) = \mathbf{W}^\top \mathbf{x}$  and  $f^{(2)}(\mathbf{h}) = \mathbf{h}^\top \mathbf{w}$ . Then  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{W}^\top \mathbf{x}$ . We could represent this function as  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}'$  where  $\mathbf{w}' = \mathbf{W}\mathbf{w}$ .

Clearly, we must use a nonlinear function to describe the features. Most neural networks do so using an affine transformation controlled by learned parameters, followed by a fixed, nonlinear function called an activation function. We use that strategy here, by defining  $\mathbf{h} = g(\mathbf{W}^\top \mathbf{x} + \mathbf{c})$ , where  $\mathbf{W}$  provides the weights of a linear transformation and  $\mathbf{c}$  the biases. Previously, to describe a linear regression

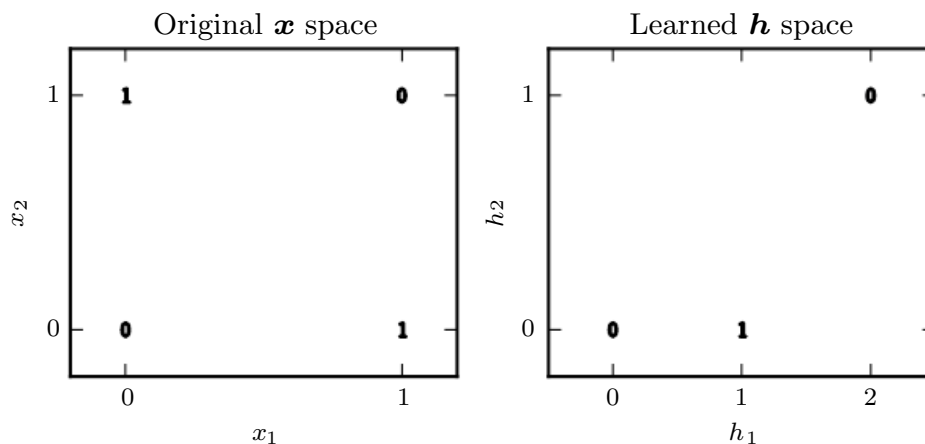


Figure 6.1: Solving the XOR problem by learning a representation. The bold numbers printed on the plot indicate the value that the learned function must output at each point. *(Left)* A linear model applied directly to the original input cannot implement the XOR function. When  $x_1 = 0$ , the model's output must increase as  $x_2$  increases. When  $x_1 = 1$ , the model's output must decrease as  $x_2$  increases. A linear model must apply a fixed coefficient  $w_2$  to  $x_2$ . The linear model therefore cannot use the value of  $x_1$  to change the coefficient on  $x_2$  and cannot solve this problem. *(Right)* In the transformed space represented by the features extracted by a neural network, a linear model can now solve the problem. In our example solution, the two points that must have output 1 have been collapsed into a single point in feature space. In other words, the nonlinear features have mapped both  $\mathbf{x} = [1, 0]^\top$  and  $\mathbf{x} = [0, 1]^\top$  to a single point in feature space,  $\mathbf{h} = [1, 0]^\top$ . The linear model can now describe the function as increasing in  $h_1$  and decreasing in  $h_2$ . In this example, the motivation for learning the feature space is only to make the model capacity greater so that it can fit the training set. In more realistic applications, learned representations can also help the model to generalize.

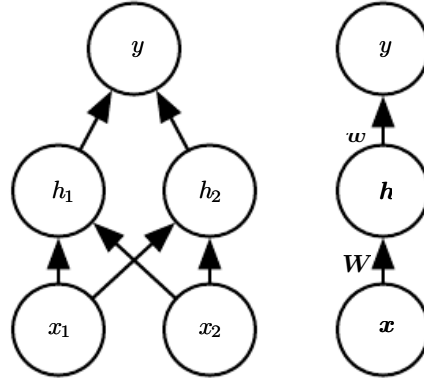


Figure 6.2: An example of a feedforward network, drawn in two different styles. Specifically, this is the feedforward network we use to solve the XOR example. It has a single hidden layer containing two units. *(Left)* In this style, we draw every unit as a node in the graph. This style is very explicit and unambiguous but for networks larger than this example it can consume too much space. *(Right)* In this style, we draw a node in the graph for each entire vector representing a layer’s activations. This style is much more compact. Sometimes we annotate the edges in this graph with the name of the parameters that describe the relationship between two layers. Here, we indicate that a matrix  $\mathbf{W}$  describes the mapping from  $\mathbf{x}$  to  $\mathbf{h}$ , and a vector  $\mathbf{w}$  describes the mapping from  $\mathbf{h}$  to  $y$ . We typically omit the intercept parameters associated with each layer when labeling this kind of drawing.

model, we used a vector of weights and a scalar bias parameter to describe an affine transformation from an input vector to an output scalar. Now, we describe an affine transformation from a vector  $\mathbf{x}$  to a vector  $\mathbf{h}$ , so an entire vector of bias parameters is needed. The activation function  $g$  is typically chosen to be a function that is applied element-wise, with  $h_i = g(\mathbf{x}^\top \mathbf{W}_{:,i} + c_i)$ . In modern neural networks, the default recommendation is to use the **rectified linear unit** or ReLU (Jarrett *et al.*, 2009; Nair and Hinton, 2010; Glorot *et al.*, 2011a) defined by the activation function  $g(z) = \max\{0, z\}$  depicted in figure 6.3.

We can now specify our complete network as

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b. \quad (6.3)$$

We can now specify a solution to the XOR problem. Let

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad (6.4)$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad (6.5)$$

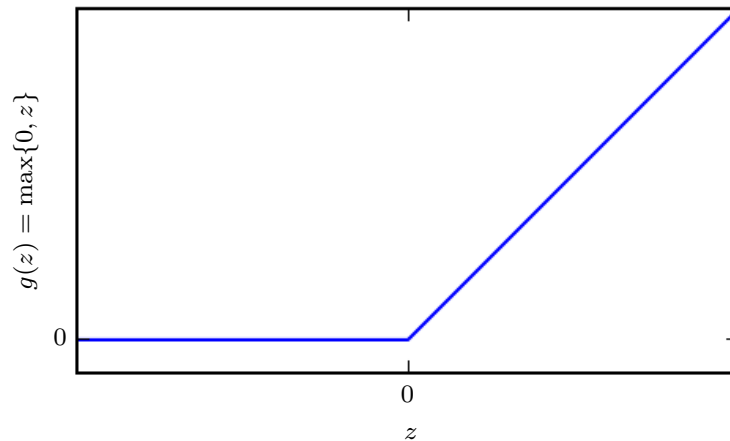


Figure 6.3: The rectified linear activation function. This activation function is the default activation function recommended for use with most feedforward neural networks. Applying this function to the output of a linear transformation yields a nonlinear transformation. However, the function remains very close to linear, in the sense that it is a piecewise linear function with two linear pieces. Because rectified linear units are nearly linear, they preserve many of the properties that make linear models easy to optimize with gradient-based methods. They also preserve many of the properties that make linear models generalize well. A common principle throughout computer science is that we can build complicated systems from minimal components. Much as a Turing machine’s memory needs only to be able to store 0 or 1 states, we can build a universal function approximator from rectified linear functions.

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad (6.6)$$

and  $b = 0$ .

We can now walk through the way that the model processes a batch of inputs. Let  $\mathbf{X}$  be the design matrix containing all four points in the binary input space, with one example per row:

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}. \quad (6.7)$$

The first step in the neural network is to multiply the input matrix by the first layer's weight matrix:

$$\mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}. \quad (6.8)$$

Next, we add the bias vector  $\mathbf{c}$ , to obtain

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}. \quad (6.9)$$

In this space, all of the examples lie along a line with slope 1. As we move along this line, the output needs to begin at 0, then rise to 1, then drop back down to 0. A linear model cannot implement such a function. To finish computing the value of  $\mathbf{h}$  for each example, we apply the rectified linear transformation:

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}. \quad (6.10)$$

This transformation has changed the relationship between the examples. They no longer lie on a single line. As shown in figure 6.1, they now lie in a space where a linear model can solve the problem.

We finish by multiplying by the weight vector  $\mathbf{w}$ :

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}. \quad (6.11)$$

The neural network has obtained the correct answer for every example in the batch.

In this example, we simply specified the solution, then showed that it obtained zero error. In a real situation, there might be billions of model parameters and billions of training examples, so one cannot simply guess the solution as we did here. Instead, a gradient-based optimization algorithm can find parameters that produce very little error. The solution we described to the XOR problem is at a global minimum of the loss function, so gradient descent could converge to this point. There are other equivalent solutions to the XOR problem that gradient descent could also find. The convergence point of gradient descent depends on the initial values of the parameters. In practice, gradient descent would usually not find clean, easily understood, integer-valued solutions like the one we presented here.

## 6.2 Gradient-Based Learning

Designing and training a neural network is not much different from training any other machine learning model with gradient descent. In section 5.10, we described how to build a machine learning algorithm by specifying an optimization procedure, a cost function, and a model family.

The largest difference between the linear models we have seen so far and neural networks is that the nonlinearity of a neural network causes most interesting loss functions to become non-convex. This means that neural networks are usually trained by using iterative, gradient-based optimizers that merely drive the cost function to a very low value, rather than the linear equation solvers used to train linear regression models or the convex optimization algorithms with global convergence guarantees used to train logistic regression or SVMs. Convex optimization converges starting from any initial parameters (in theory—in practice it is very robust but can encounter numerical problems). Stochastic gradient descent applied to non-convex loss functions has no such convergence guarantee, and is sensitive to the values of the initial parameters. For feedforward neural networks, it is important to initialize all weights to small random values. The biases may be initialized to zero or to small positive values. The iterative gradient-based optimization algorithms used to train feedforward networks and almost all other deep models will be described in detail in chapter 8, with parameter initialization in particular discussed in section 8.4. For the moment, it suffices to understand that the training algorithm is almost always based on using the gradient to descend the cost function in one way or another. The specific algorithms are improvements and refinements on the ideas of gradient descent, introduced in section 4.3, and,



more specifically, are most often improvements of the stochastic gradient descent algorithm, introduced in section 5.9.

We can of course, train models such as linear regression and support vector machines with gradient descent too, and in fact this is common when the training set is extremely large. From this point of view, training a neural network is not much different from training any other model. Computing the gradient is slightly more complicated for a neural network, but can still be done efficiently and exactly. Section 6.5 will describe how to obtain the gradient using the back-propagation algorithm and modern generalizations of the back-propagation algorithm.

As with other machine learning models, to apply gradient-based learning we must choose a cost function, and we must choose how to represent the output of the model. We now revisit these design considerations with special emphasis on the neural networks scenario.

### 6.2.1 Cost Functions

An important aspect of the design of a deep neural network is the choice of the cost function. Fortunately, the cost functions for neural networks are more or less the same as those for other parametric models, such as linear models.

In most cases, our parametric model defines a distribution  $p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$  and we simply use the principle of maximum likelihood. This means we use the cross-entropy between the training data and the model's predictions as the cost function.

Sometimes, we take a simpler approach, where rather than predicting a complete probability distribution over  $\mathbf{y}$ , we merely predict some statistic of  $\mathbf{y}$  conditioned on  $\mathbf{x}$ . Specialized loss functions allow us to train a predictor of these estimates.

The total cost function used to train a neural network will often combine one of the primary cost functions described here with a regularization term. We have already seen some simple examples of regularization applied to linear models in section 5.2.2. The weight decay approach used for linear models is also directly applicable to deep neural networks and is among the most popular regularization strategies. More advanced regularization strategies for neural networks will be described in chapter 7.

#### 6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

Most modern neural networks are trained using maximum likelihood. This means that the cost function is simply the negative log-likelihood, equivalently described

as the cross-entropy between the training data and the model distribution. This cost function is given by

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x}). \quad (6.12)$$

The specific form of the cost function changes from model to model, depending on the specific form of  $\log p_{\text{model}}$ . The expansion of the above equation typically yields some terms that do not depend on the model parameters and may be discarded. For example, as we saw in section 5.5.1, if  $p_{\text{model}}(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y}; f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{I})$ , then we recover the mean squared error cost,

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \boldsymbol{\theta})\|^2 + \text{const}, \quad (6.13)$$

up to a scaling factor of  $\frac{1}{2}$  and a term that does not depend on  $\boldsymbol{\theta}$ . The discarded constant is based on the variance of the Gaussian distribution, which in this case we chose not to parametrize. Previously, we saw that the equivalence between maximum likelihood estimation with an output distribution and minimization of mean squared error holds for a linear model, but in fact, the equivalence holds regardless of the  $f(\mathbf{x}; \boldsymbol{\theta})$  used to predict the mean of the Gaussian.

An advantage of this approach of deriving the cost function from maximum likelihood is that it removes the burden of designing cost functions for each model. Specifying a model  $p(\mathbf{y} \mid \mathbf{x})$  automatically determines a cost function  $\log p(\mathbf{y} \mid \mathbf{x})$ .

One recurring theme throughout neural network design is that the gradient of the cost function must be large and predictable enough to serve as a good guide for the learning algorithm. Functions that saturate (become very flat) undermine this objective because they make the gradient become very small. In many cases this happens because the activation functions used to produce the output of the hidden units or the output units saturate. The negative log-likelihood helps to avoid this problem for many models. Many output units involve an exp function that can saturate when its argument is very negative. The log function in the negative log-likelihood cost function undoes the exp of some output units. We will discuss the interaction between the cost function and the choice of output unit in section 6.2.2.

One unusual property of the cross-entropy cost used to perform maximum likelihood estimation is that it usually does not have a minimum value when applied to the models commonly used in practice. For discrete output variables, most models are parametrized in such a way that they cannot represent a probability of zero or one, but can come arbitrarily close to doing so. Logistic regression is an example of such a model. For real-valued output variables, if the model

can control the density of the output distribution (for example, by learning the variance parameter of a Gaussian output distribution) then it becomes possible to assign extremely high density to the correct training set outputs, resulting in cross-entropy approaching negative infinity. Regularization techniques described in chapter 7 provide several different ways of modifying the learning problem so that the model cannot reap unlimited reward in this way.

### 6.2.1.2 Learning Conditional Statistics

Instead of learning a full probability distribution  $p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$  we often want to learn just one conditional statistic of  $\mathbf{y}$  given  $\mathbf{x}$ .

For example, we may have a predictor  $f(\mathbf{x}; \boldsymbol{\theta})$  that we wish to predict the mean of  $\mathbf{y}$ .

If we use a sufficiently powerful neural network, we can think of the neural network as being able to represent any function  $f$  from a wide class of functions, with this class being limited only by features such as continuity and boundedness rather than by having a specific parametric form. From this point of view, we can view the cost function as being a **functional** rather than just a function. A functional is a mapping from functions to real numbers. We can thus think of learning as choosing a function rather than merely choosing a set of parameters. We can design our cost functional to have its minimum occur at some specific function we desire. For example, we can design the cost functional to have its minimum lie on the function that maps  $\mathbf{x}$  to the expected value of  $\mathbf{y}$  given  $\mathbf{x}$ . Solving an optimization problem with respect to a function requires a mathematical tool called **calculus of variations**, described in section 19.4.2. It is not necessary to understand calculus of variations to understand the content of this chapter. At the moment, it is only necessary to understand that calculus of variations may be used to derive the following two results.

Our first result derived using calculus of variations is that solving the optimization problem

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|^2 \quad (6.14)$$

yields

$$f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y}|\mathbf{x})}[\mathbf{y}], \quad (6.15)$$

so long as this function lies within the class we optimize over. In other words, if we could train on infinitely many samples from the true data generating distribution, minimizing the mean squared error cost function gives a function that predicts the mean of  $\mathbf{y}$  for each value of  $\mathbf{x}$ .

Different cost functions give different statistics. A second result derived using calculus of variations is that

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|_1 \quad (6.16)$$

yields a function that predicts the *median* value of  $\mathbf{y}$  for each  $\mathbf{x}$ , so long as such a function may be described by the family of functions we optimize over. This cost function is commonly called **mean absolute error**.

Unfortunately, mean squared error and mean absolute error often lead to poor results when used with gradient-based optimization. Some output units that saturate produce very small gradients when combined with these cost functions. This is one reason that the cross-entropy cost function is more popular than mean squared error or mean absolute error, even when it is not necessary to estimate an entire distribution  $p(\mathbf{y} \mid \mathbf{x})$ .

## 6.2.2 Output Units

The choice of cost function is tightly coupled with the choice of output unit. Most of the time, we simply use the cross-entropy between the data distribution and the model distribution. The choice of how to represent the output then determines the form of the cross-entropy function.

Any kind of neural network unit that may be used as an output can also be used as a hidden unit. Here, we focus on the use of these units as outputs of the model, but in principle they can be used internally as well. We revisit these units with additional detail about their use as hidden units in section 6.3.

Throughout this section, we suppose that the feedforward network provides a set of hidden features defined by  $\mathbf{h} = f(\mathbf{x}; \boldsymbol{\theta})$ . The role of the output layer is then to provide some additional transformation from the features to complete the task that the network must perform.

### 6.2.2.1 Linear Units for Gaussian Output Distributions

One simple kind of output unit is an output unit based on an affine transformation with no nonlinearity. These are often just called linear units.

Given features  $\mathbf{h}$ , a layer of linear output units produces a vector  $\hat{\mathbf{y}} = \mathbf{W}^\top \mathbf{h} + \mathbf{b}$ .

Linear output layers are often used to produce the mean of a conditional Gaussian distribution:

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \mathbf{I}). \quad (6.17)$$

Maximizing the log-likelihood is then equivalent to minimizing the mean squared error.

The maximum likelihood framework makes it straightforward to learn the covariance of the Gaussian too, or to make the covariance of the Gaussian be a function of the input. However, the covariance must be constrained to be a positive definite matrix for all inputs. It is difficult to satisfy such constraints with a linear output layer, so typically other output units are used to parametrize the covariance. Approaches to modeling the covariance are described shortly, in section 6.2.2.4.

Because linear units do not saturate, they pose little difficulty for gradient-based optimization algorithms and may be used with a wide variety of optimization algorithms.

### 6.2.2.2 Sigmoid Units for Bernoulli Output Distributions

Many tasks require predicting the value of a binary variable  $y$ . Classification problems with two classes can be cast in this form.

The maximum-likelihood approach is to define a Bernoulli distribution over  $y$  conditioned on  $\mathbf{x}$ .

A Bernoulli distribution is defined by just a single number. The neural net needs to predict only  $P(y = 1 \mid \mathbf{x})$ . For this number to be a valid probability, it must lie in the interval  $[0, 1]$ .

Satisfying this constraint requires some careful design effort. Suppose we were to use a linear unit, and threshold its value to obtain a valid probability:

$$P(y = 1 \mid \mathbf{x}) = \max \left\{ 0, \min \left\{ 1, \mathbf{w}^\top \mathbf{h} + b \right\} \right\}. \quad (6.18)$$

This would indeed define a valid conditional distribution, but we would not be able to train it very effectively with gradient descent. Any time that  $\mathbf{w}^\top \mathbf{h} + b$  strayed outside the unit interval, the gradient of the output of the model with respect to its parameters would be  $\mathbf{0}$ . A gradient of  $\mathbf{0}$  is typically problematic because the learning algorithm no longer has a guide for how to improve the corresponding parameters.

Instead, it is better to use a different approach that ensures there is always a strong gradient whenever the model has the wrong answer. This approach is based on using sigmoid output units combined with maximum likelihood.

A sigmoid output unit is defined by

$$\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b) \quad (6.19)$$

where  $\sigma$  is the logistic sigmoid function described in section 3.10.

We can think of the sigmoid output unit as having two components. First, it uses a linear layer to compute  $z = \mathbf{w}^\top \mathbf{h} + b$ . Next, it uses the sigmoid activation function to convert  $z$  into a probability.

We omit the dependence on  $\mathbf{x}$  for the moment to discuss how to define a probability distribution over  $y$  using the value  $z$ . The sigmoid can be motivated by constructing an unnormalized probability distribution  $\tilde{P}(y)$ , which does not sum to 1. We can then divide by an appropriate constant to obtain a valid probability distribution. If we begin with the assumption that the unnormalized log probabilities are linear in  $y$  and  $z$ , we can exponentiate to obtain the unnormalized probabilities. We then normalize to see that this yields a Bernoulli distribution controlled by a sigmoidal transformation of  $z$ :

$$\log \tilde{P}(y) = yz \quad (6.20)$$

$$\tilde{P}(y) = \exp(yz) \quad (6.21)$$

$$P(y) = \frac{\exp(yz)}{\sum_{y'=0}^1 \exp(y'z)} \quad (6.22)$$

$$P(y) = \sigma((2y - 1)z). \quad (6.23)$$

Probability distributions based on exponentiation and normalization are common throughout the statistical modeling literature. The  $z$  variable defining such a distribution over binary variables is called a **logit**.

This approach to predicting the probabilities in log-space is natural to use with maximum likelihood learning. Because the cost function used with maximum likelihood is  $-\log P(y | \mathbf{x})$ , the log in the cost function undoes the exp of the sigmoid. Without this effect, the saturation of the sigmoid could prevent gradient-based learning from making good progress. The loss function for maximum likelihood learning of a Bernoulli parametrized by a sigmoid is

$$J(\boldsymbol{\theta}) = -\log P(y | \mathbf{x}) \quad (6.24)$$

$$= -\log \sigma((2y - 1)z) \quad (6.25)$$

$$= \zeta((1 - 2y)z). \quad (6.26)$$

This derivation makes use of some properties from section 3.10. By rewriting the loss in terms of the softplus function, we can see that it saturates only when  $(1 - 2y)z$  is very negative. Saturation thus occurs only when the model already has the right answer—when  $y = 1$  and  $z$  is very positive, or  $y = 0$  and  $z$  is very negative. When  $z$  has the wrong sign, the argument to the softplus function,

$(1 - 2y)z$ , may be simplified to  $|z|$ . As  $|z|$  becomes large while  $z$  has the wrong sign, the softplus function asymptotes toward simply returning its argument  $|z|$ . The derivative with respect to  $z$  asymptotes to  $\text{sign}(z)$ , so, in the limit of extremely incorrect  $z$ , the softplus function does not shrink the gradient at all. This property is very useful because it means that gradient-based learning can act to quickly correct a mistaken  $z$ .

When we use other loss functions, such as mean squared error, the loss can saturate anytime  $\sigma(z)$  saturates. The sigmoid activation function saturates to 0 when  $z$  becomes very negative and saturates to 1 when  $z$  becomes very positive. The gradient can shrink too small to be useful for learning whenever this happens, whether the model has the correct answer or the incorrect answer. For this reason, maximum likelihood is almost always the preferred approach to training sigmoid output units.

Analytically, the logarithm of the sigmoid is always defined and finite, because the sigmoid returns values restricted to the open interval  $(0, 1)$ , rather than using the entire closed interval of valid probabilities  $[0, 1]$ . In software implementations, to avoid numerical problems, it is best to write the negative log-likelihood as a function of  $z$ , rather than as a function of  $\hat{y} = \sigma(z)$ . If the sigmoid function underflows to zero, then taking the logarithm of  $\hat{y}$  yields negative infinity.

### 6.2.2.3 Softmax Units for Multinoulli Output Distributions

Any time we wish to represent a probability distribution over a discrete variable with  $n$  possible values, we may use the softmax function. This can be seen as a generalization of the sigmoid function which was used to represent a probability distribution over a binary variable.

Softmax functions are most often used as the output of a classifier, to represent the probability distribution over  $n$  different classes. More rarely, softmax functions can be used inside the model itself, if we wish the model to choose between one of  $n$  different options for some internal variable.

In the case of binary variables, we wished to produce a single number

$$\hat{y} = P(y = 1 \mid \mathbf{x}). \quad (6.27)$$

Because this number needed to lie between 0 and 1, and because we wanted the logarithm of the number to be well-behaved for gradient-based optimization of the log-likelihood, we chose to instead predict a number  $z = \log \tilde{P}(y = 1 \mid \mathbf{x})$ . Exponentiating and normalizing gave us a Bernoulli distribution controlled by the sigmoid function.

To generalize to the case of a discrete variable with  $n$  values, we now need to produce a vector  $\hat{\mathbf{y}}$ , with  $\hat{y}_i = P(y = i \mid \mathbf{x})$ . We require not only that each element of  $\hat{\mathbf{y}}$  be between 0 and 1, but also that the entire vector sums to 1 so that it represents a valid probability distribution. The same approach that worked for the Bernoulli distribution generalizes to the multinoulli distribution. First, a linear layer predicts unnormalized log probabilities:

$$\mathbf{z} = \mathbf{W}^\top \mathbf{h} + \mathbf{b}, \quad (6.28)$$

where  $z_i = \log \tilde{P}(y = i \mid \mathbf{x})$ . The softmax function can then exponentiate and normalize  $\mathbf{z}$  to obtain the desired  $\hat{\mathbf{y}}$ . Formally, the softmax function is given by

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}. \quad (6.29)$$

As with the logistic sigmoid, the use of the  $\exp$  function works very well when training the softmax to output a target value  $y$  using maximum log-likelihood. In this case, we wish to maximize  $\log P(y = i; \mathbf{z}) = \log \text{softmax}(\mathbf{z})_i$ . Defining the softmax in terms of  $\exp$  is natural because the  $\log$  in the log-likelihood can undo the  $\exp$  of the softmax:

$$\log \text{softmax}(\mathbf{z})_i = z_i - \log \sum_j \exp(z_j). \quad (6.30)$$

The first term of equation 6.30 shows that the input  $z_i$  always has a direct contribution to the cost function. Because this term cannot saturate, we know that learning can proceed, even if the contribution of  $z_i$  to the second term of equation 6.30 becomes very small. When maximizing the log-likelihood, the first term encourages  $z_i$  to be pushed up, while the second term encourages all of  $\mathbf{z}$  to be pushed down. To gain some intuition for the second term,  $\log \sum_j \exp(z_j)$ , observe that this term can be roughly approximated by  $\max_j z_j$ . This approximation is based on the idea that  $\exp(z_k)$  is insignificant for any  $z_k$  that is noticeably less than  $\max_j z_j$ . The intuition we can gain from this approximation is that the negative log-likelihood cost function always strongly penalizes the most active incorrect prediction. If the correct answer already has the largest input to the softmax, then the  $-z_i$  term and the  $\log \sum_j \exp(z_j) \approx \max_j z_j = z_i$  terms will roughly cancel. This example will then contribute little to the overall training cost, which will be dominated by other examples that are not yet correctly classified.

So far we have discussed only a single example. Overall, unregularized maximum likelihood will drive the model to learn parameters that drive the softmax to predict