

# 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>28</b>
<b>2 Linear Algebra</b>	<b>30</b>
2.1 Scalars, Vectors, Matrices and Tensors . . . . .	30
2.2 Multiplying Matrices and Vectors . . . . .	33
2.3 Identity and Inverse Matrices . . . . .	35
2.4 Linear Dependence and Span . . . . .	36
2.5 Norms . . . . .	38
2.6 Special Kinds of Matrices and Vectors . . . . .	39
2.7 Eigendecomposition . . . . .	41
2.8 Singular Value Decomposition . . . . .	43
2.9 The Moore-Penrose Pseudoinverse . . . . .	44
2.10 The Trace Operator . . . . .	45
2.11 The Determinant . . . . .	46
2.12 Example: Principal Components Analysis . . . . .	47
<b>3 Probability and Information Theory</b>	<b>52</b>
3.1 Why Probability? . . . . .	53

3.2	Random Variables . . . . .	55
3.3	Probability Distributions . . . . .	55
3.4	Marginal Probability . . . . .	57
3.5	Conditional Probability . . . . .	58
3.6	The Chain Rule of Conditional Probabilities . . . . .	58
3.7	Independence and Conditional Independence . . . . .	59
3.8	Expectation, Variance and Covariance . . . . .	59
3.9	Common Probability Distributions . . . . .	61
3.10	Useful Properties of Common Functions . . . . .	66
3.11	Bayes' Rule . . . . .	69
3.12	Technical Details of Continuous Variables . . . . .	70
3.13	Information Theory . . . . .	71
3.14	Structured Probabilistic Models . . . . .	74
<b>4</b>	<b>Numerical Computation</b>	<b>79</b>
4.1	Overflow and Underflow . . . . .	79
4.2	Poor Conditioning . . . . .	81
4.3	Gradient-Based Optimization . . . . .	81
4.4	Constrained Optimization . . . . .	92
4.5	Example: Linear Least Squares . . . . .	94
<b>5</b>	<b>Machine Learning Basics</b>	<b>97</b>
5.1	Learning Algorithms . . . . .	98
5.2	Capacity, Overfitting and Underfitting . . . . .	109
5.3	Hyperparameters and Validation Sets . . . . .	119
5.4	Estimators, Bias and Variance . . . . .	121
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 . . . . .	203
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 . . . . .	244
7.7	Multi-Task Learning . . . . .	245
7.8	Early Stopping . . . . .	246
7.9	Parameter Tying and Parameter Sharing . . . . .	251
7.10	Sparse Representations . . . . .	253
7.11	Bagging and Other Ensemble Methods . . . . .	255
7.12	Dropout . . . . .	257
7.13	Adversarial Training . . . . .	267
7.14	Tangent Distance, Tangent Prop, and Manifold Tangent Classifier	269
<b>8</b>	<b>Optimization for Training Deep Models</b>	<b>275</b>
8.1	How Learning Differs from Pure Optimization . . . . .	276
8.2	Challenges in Neural Network Optimization . . . . .	283
8.3	Basic Algorithms . . . . .	295
8.4	Parameter Initialization Strategies . . . . .	301
8.5	Algorithms with Adaptive Learning Rates . . . . .	307
8.6	Approximate Second-Order Methods . . . . .	311
8.7	Optimization Strategies and Meta-Algorithms . . . . .	319
<b>9</b>	<b>Convolutional Networks</b>	<b>332</b>
9.1	The Convolution Operation . . . . .	333
9.2	Motivation . . . . .	335
9.3	Pooling . . . . .	341
9.4	Convolution and Pooling as an Infinitely Strong Prior . . . . .	347
9.5	Variants of the Basic Convolution Function . . . . .	349
9.6	Structured Outputs . . . . .	360
9.7	Data Types . . . . .	362
9.8	Efficient Convolution Algorithms . . . . .	364
9.9	Random or Unsupervised Features . . . . .	364

9.10	The Neuroscientific Basis for Convolutional Networks . . . . .	366
9.11	Convolutional Networks and the History of Deep Learning . . . . .	372
<b>10</b>	<b>Sequence Modeling: Recurrent and Recursive Nets</b>	<b>375</b>
10.1	Unfolding Computational Graphs . . . . .	376
10.2	Recurrent Neural Networks . . . . .	380
10.3	Bidirectional RNNs . . . . .	396
10.4	Encoder-Decoder Sequence-to-Sequence Architectures . . . . .	398
10.5	Deep Recurrent Networks . . . . .	400
10.6	Recursive Neural Networks . . . . .	402
10.7	The Challenge of Long-Term Dependencies . . . . .	403
10.8	Echo State Networks . . . . .	406
10.9	Skip Connections through Time . . . . .	408
10.10	Leaky Units and a Spectrum of Different Time Scales . . . . .	409
10.11	The Long Short-Term Memory and Other Gated RNNs . . . . .	410
10.12	Optimization for Long-Term Dependencies . . . . .	414
10.13	Regularizing to Encourage Information Flow . . . . .	418
10.14	Organizing the State at Multiple Time Scales . . . . .	419
10.15	Explicit Memory . . . . .	420
<b>11</b>	<b>Practical methodology</b>	<b>423</b>
11.1	Performance Metrics . . . . .	424
11.2	Default Baseline Models . . . . .	427
11.3	Determining Whether to Gather More Data . . . . .	428
11.4	Selecting Hyperparameters . . . . .	429
11.5	Debugging Strategies . . . . .	438
11.6	Example: Multi-Digit Number Recognition . . . . .	442
<b>12</b>	<b>Applications</b>	<b>445</b>
12.1	Large Scale Deep Learning . . . . .	445
12.2	Computer Vision . . . . .	454
12.3	Speech Recognition . . . . .	460
12.4	Natural Language Processing . . . . .	462
12.5	Other Applications . . . . .	478
<b>III</b>	<b>Deep Learning Research</b>	<b>487</b>
<b>13</b>	<b>Linear Factor Models</b>	<b>490</b>
13.1	Probabilistic PCA and Factor Analysis . . . . .	491

13.2	Independent Component Analysis (ICA) . . . . .	492
13.3	Slow Feature Analysis . . . . .	495
13.4	Sparse Coding . . . . .	497
13.5	Manifold Interpretation of PCA . . . . .	500
<b>14</b>	<b>Autoencoders</b>	<b>503</b>
14.1	Undercomplete Autoencoders . . . . .	503
14.2	Regularized Autoencoders . . . . .	505
14.3	Representational Power, Layer Size and Depth . . . . .	509
14.4	Stochastic Encoders and Decoders . . . . .	510
14.5	Denoising Autoencoders . . . . .	511
14.6	Learning Manifolds with Autoencoders . . . . .	516
14.7	Contractive Autoencoders . . . . .	522
14.8	Predictive Sparse Decomposition . . . . .	524
14.9	Applications of Autoencoders . . . . .	525
<b>15</b>	<b>Representation Learning</b>	<b>527</b>
15.1	Greedy Layer-Wise Unsupervised Pretraining . . . . .	529
15.2	Transfer Learning and Domain Adaptation . . . . .	538
15.3	Semi-Supervised Disentangling of Causal Factors . . . . .	542
15.4	Distributed Representation . . . . .	547
15.5	Exponential Gains from Depth . . . . .	555
15.6	What is a Good Representation? . . . . .	557
<b>16</b>	<b>Structured Probabilistic Models for Deep Learning</b>	<b>560</b>
16.1	The Challenge of Unstructured Modeling . . . . .	561
16.2	Using Graphs to Describe Model Structure . . . . .	565
16.3	Sampling from Graphical Models . . . . .	581
16.4	Advantages of Structured Modeling . . . . .	583
16.5	Learning about Dependencies . . . . .	584
16.6	Inference and Approximate Inference . . . . .	585
16.7	The Deep Learning Approach to Structured Probabilistic Models	586
<b>17</b>	<b>Monte Carlo Methods</b>	<b>592</b>
17.1	Sampling and Monte Carlo Methods . . . . .	592
17.2	Markov Chain Monte Carlo Methods . . . . .	598
17.3	Gibbs Sampling . . . . .	601
17.4	The Challenge of Mixing between Separated Modes . . . . .	602

<b>18</b>	<b>Confronting the Partition Function</b>	<b>608</b>
18.1	The Log-Likelihood Gradient of Undirected Models . . . . .	609
18.2	Stochastic Maximum Likelihood and Contrastive Divergence . . .	611
18.3	Pseudolikelihood . . . . .	619
18.4	Score Matching and Ratio Matching . . . . .	621
18.5	Denoising Score Matching . . . . .	623
18.6	Noise-Contrastive Estimation . . . . .	624
18.7	Estimating the Partition Function . . . . .	627
<b>19</b>	<b>Approximate inference</b>	<b>635</b>
19.1	Inference as Optimization . . . . .	637
19.2	Expectation Maximization . . . . .	638
19.3	MAP Inference and Sparse Coding . . . . .	639
19.4	Variational Inference and Learning . . . . .	642
19.5	Learned Approximate Inference . . . . .	655
<b>20</b>	<b>Deep Generative Models</b>	<b>658</b>
20.1	Boltzmann Machines . . . . .	658
20.2	Restricted Boltzmann Machines . . . . .	662
20.3	Deep Belief Networks . . . . .	665
20.4	Deep Boltzmann Machines . . . . .	668
20.5	Boltzmann Machines for Real-Valued Data . . . . .	683
20.6	Convolutional Boltzmann Machines . . . . .	691
20.7	Boltzmann Machines for Structured or Sequential Outputs . . . .	693
20.8	Other Boltzmann Machines . . . . .	695
20.9	Back-Propagation through Random Operations . . . . .	696
20.10	Directed Generative Nets . . . . .	700
20.11	Auto-Regressive Networks . . . . .	714
20.12	Drawing Samples from Autoencoders . . . . .	720
20.13	Generative Stochastic Networks . . . . .	723
20.14	Other Generation Schemes . . . . .	725
20.15	Evaluating Generative Models . . . . .	726
20.16	Conclusion . . . . .	729
	<b>Bibliography</b>	<b>730</b>
	<b>Index</b>	<b>783</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.