## Deep Learning

Ian Goodfellow Yoshua Bengio Aaron Courville

## Contents

Website			vii
Acknowledgments			
Notation			xi
1	Intro 1.1 1.2	Oduction  Who Should Read This Book?	
Ι	Appl	lied Math and Machine Learning Basics	29
2	Line	ar Algebra	31
	2.1	Scalars, Vectors, Matrices and Tensors	31
	2.2	Multiplying Matrices and Vectors	
	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	
	2.9	The Moore-Penrose Pseudoinverse	45
	2.10	The Trace Operator	
	2.11	The Determinant	
	2.12	Example: Principal Components Analysis	48
3	Prob	pability and Information Theory	<b>53</b>
	3.1	Why Probability?	54

	3.2	Random Variables	56
	3.3	Probability Distributions	
	3.4	Marginal Probability	
	3.5	Conditional Probability	
	3.6	The Chain Rule of Conditional Probabilities	
	3.7	Independence and Conditional Independence	60
	3.8	Expectation, Variance and Covariance	
	3.9	Common Probability Distributions	
	3.10	Useful Properties of Common Functions	
	3.11	Bayes' Rule	70
	3.12	Technical Details of Continuous Variables	
	3.13	Information Theory	73
	3.14	Structured Probabilistic Models	
4	Num	nerical Computation	80
	4.1	Overflow and Underflow	80
	4.2	Poor Conditioning	
	4.3	Gradient-Based Optimization	82
	4.4	Constrained Optimization	
	4.5	Example: Linear Least Squares	
5	Mac	hine Learning Basics	98
	5.1	Learning Algorithms	99
	5.2	Capacity, Overfitting and Underfitting	
	5.3	Hyperparameters and Validation Sets	
	5.4	Estimators, Bias and Variance	
	5.5	Maximum Likelihood Estimation	
	5.6	Bayesian Statistics	135
	5.7	Supervised Learning Algorithms	
	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
	Б		400
II	Dee	p Networks: Modern Practices	166
6	Deep	p Feedforward Networks	168
	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	
7	Regi	ularization for Deep Learning	<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	
8	Optimization for Training Deep Models 274			
	8.1	How Learning Differs from Pure Optimization	275	
	8.2	Challenges in Neural Network Optimization	282	
	8.3		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	
9	Convolutional Networks 330			
	9.1	The Convolution Operation	331	
	9.2	Motivation	335	
	9.3	Pooling	339	
	9.4		345	
	9.5		347	
	9.6	Structured Outputs	358	
	9.7		360	
	9.8		362	
	9.9	_	363	

	9.10	The Neuroscientific Basis for Convolutional Networks	34
	9.11	Convolutional Networks and the History of Deep Learning 37	71
10	Sequ	ence Modeling: Recurrent and Recursive Nets 37	'3
	10.1	Unfolding Computational Graphs	<b>'</b> 5
	10.2	Recurrent Neural Networks	78
	10.3	Bidirectional RNNs	)4
	10.4	Encoder-Decoder Sequence-to-Sequence Architectures 39	96
	10.5	Deep Recurrent Networks	8(
	10.6	Recursive Neural Networks	0(
	10.7	The Challenge of Long-Term Dependencies	)1
	10.8	Echo State Networks	)4
	10.9	Leaky Units and Other Strategies for Multiple Time Scales 40	06
	10.10	The Long Short-Term Memory and Other Gated RNNs 40	
	10.11	Optimization for Long-Term Dependencies	.3
	10.12	Explicit Memory	6
11	Pract	tical Methodology 42	21
	11.1	Performance Metrics	2
	11.2	Default Baseline Models	
	11.3	Determining Whether to Gather More Data	
	11.4	Selecting Hyperparameters	
	11.5	Debugging Strategies	
	11.6	Example: Multi-Digit Number Recognition	-0
12	Appl	ications 44	3
	12.1	Large-Scale Deep Learning	13
	12.2	Computer Vision	
	12.3	Speech Recognition	8
	12.4	Natural Language Processing	
	12.5	Other Applications	<b>'</b> 8
	Ъ	T : D 1	
III	Dee	ep Learning Research 48	6
<b>13</b>		ar Factor Models 48	
	13.1	Probabilistic PCA and Factor Analysis	
	13.2	Independent Component Analysis (ICA)	
	13.3	Slow Feature Analysis	
	13.4	Sparse Coding 49	16

	13.5	Manifold Interpretation of PCA	499	
14	Autoencoders 502			
	14.1	Undercomplete Autoencoders	503	
	14.2	Regularized Autoencoders		
	14.3	Representational Power, Layer Size and Depth		
	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		
<b>15</b>	Rep	resentation Learning	<b>52</b> 6	
	15.1	Greedy Layer-Wise Unsupervised Pretraining	528	
	15.2	Transfer Learning and Domain Adaptation		
	15.3	Semi-Supervised Disentangling of Causal Factors	541	
	15.4	Distributed Representation	546	
	15.5	Exponential Gains from Depth		
	15.6	Providing Clues to Discover Underlying Causes	554	
<b>16</b>	Structured Probabilistic Models for Deep Learning 558			
	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>	Mon	te Carlo Methods	<b>59</b> 0	
	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	
18	Conf	fronting the Partition Function	605	
	18.1	The Log-Likelihood Gradient		
	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
Appr	roximate Inference	631
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
Deep	Generative Models	654
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	
20.10	Directed Generative Nets	692
	•	
20.15	Conclusion	720
oliogra	aphy	721
$\mathbf{lex}$		777
	18.4 18.5 18.6 18.7  Appr 19.1 19.2 19.3 19.4 19.5  Deep 20.1 20.2 20.3 20.4 20.5 20.6 20.7 20.8 20.9 20.10 20.11 20.12 20.13 20.14 20.15  bliogra	18.4 Score Matching and Ratio Matching 18.5 Denoising Score Matching 18.6 Noise-Contrastive Estimation 18.7 Estimating the Partition Function  Approximate Inference 19.1 Inference as Optimization 19.2 Expectation Maximization 19.3 MAP Inference and Sparse Coding 19.4 Variational Inference and Learning 19.5 Learned Approximate Inference  Deep Generative Models 20.1 Boltzmann Machines 20.2 Restricted Boltzmann Machines 20.3 Deep Belief Networks 20.4 Deep Boltzmann Machines 20.5 Boltzmann Machines 20.6 Convolutional Boltzmann Machines 20.7 Boltzmann Machines for Real-Valued Data 20.8 Other Boltzmann Machines 20.9 Back-Propagation through Random Operations 20.10 Directed Generative Nets 20.11 Drawing Samples from Autoencoders 20.12 Generative Stochastic Networks 20.13 Other Generation Schemes 20.14 Evaluating Generative Models 20.15 Conclusion

## Website

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.