

## Machine Learning Fundamentals

This lucid, accessible introduction to supervised machine learning presents core concepts in a focused and logical way that is easy for beginners to follow. The author assumes basic calculus, linear algebra, probability and statistics but no prior exposure to machine learning. Coverage includes widely used traditional methods such as SVMs, boosted trees, HMMs, and LDAs, plus popular deep learning methods such as convolution neural nets, attention, transformers, and GANs. Organized in a coherent presentation framework that emphasizes the big picture, the text introduces each method clearly and concisely “from scratch” based on the fundamentals. All methods and algorithms are described by a clean and consistent style, with a minimum of unnecessary detail. Numerous case studies and concrete examples demonstrate how the methods can be applied in a variety of contexts.

Hui Jiang is a Professor of Electrical Engineering and Computer Science at York University, where he has been since 2002. His main research interests include machine learning, particularly deep learning, and its applications to speech and audio processing, natural language processing, and computer vision. Over the past 30 years, he has worked on a wide range of research problems from these areas and published hundreds of technical articles and papers in the mainstream journals and top-tier conferences. His works have won the prestigious IEEE Best Paper Award and the ACL Outstanding Paper honor.

Simplicity is the ultimate sophistication.

—Leonardo da Vinci

# Machine Learning Fundamentals

A Concise Introduction

Hui Jiang

*York University, Toronto*



**CAMBRIDGE**  
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre,  
New Delhi – 110025, India

103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

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

Information on this title: [www.cambridge.org/9781108837040](http://www.cambridge.org/9781108837040)

DOI: 10.1017/9781108938051

© Hui Jiang 2021

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2021

*A catalogue record for this publication is available from the British Library.*

ISBN 978-1-108-83704-0 Hardback

ISBN 978-1-108-94002-3 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

# Contents

<b>Preface</b>	<b>xi</b>
<b>Notation</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 <b>What Is Machine Learning?</b> . . . . .	1
1.2 <b>Basic Concepts in Machine Learning</b> . . . . .	4
1.2.1 Classification versus Regression . . . . .	4
1.2.2 Supervised versus Unsupervised Learning . . . . .	5
1.2.3 Simple versus Complex Models . . . . .	5
1.2.4 Parametric versus Nonparametric Models . . . . .	7
1.2.5 Overfitting versus Underfitting . . . . .	8
1.2.6 Bias–Variance Trade-Off . . . . .	10
1.3 <b>General Principles in Machine Learning</b> . . . . .	11
1.3.1 Occam’s Razor . . . . .	11
1.3.2 No-Free-Lunch Theorem . . . . .	11
1.3.3 Law of the Smooth World . . . . .	12
1.3.4 Curse of Dimensionality . . . . .	14
1.4 <b>Advanced Topics in Machine Learning</b> . . . . .	15
1.4.1 Reinforcement Learning . . . . .	15
1.4.2 Meta-Learning . . . . .	16
1.4.3 Causal Inference . . . . .	16
1.4.4 Other Advanced Topics . . . . .	16
Exercises . . . . .	18
<b>2 Mathematical Foundation</b>	<b>19</b>
2.1 <b>Linear Algebra</b> . . . . .	19
2.1.1 Vectors and Matrices . . . . .	19
2.1.2 Linear Transformation as Matrix Multiplication . . . . .	20
2.1.3 Basic Matrix Operations . . . . .	21

2.1.4	Eigenvalues and Eigenvectors . . . . .	23
2.1.5	Matrix Calculus . . . . .	25
2.2	<b>Probability and Statistics</b> . . . . .	27
2.2.1	Random Variables and Distributions . . . . .	27
2.2.2	Expectation: Mean, Variance, and Moments . . . . .	28
2.2.3	Joint, Marginal, and Conditional Distributions . . . . .	30
2.2.4	Common Probability Distributions . . . . .	33
2.2.5	Transformation of Random Variables . . . . .	40
2.3	<b>Information Theory</b> . . . . .	41
2.3.1	Information and Entropy . . . . .	41
2.3.2	Mutual Information . . . . .	43
2.3.3	KL Divergence . . . . .	46
2.4	<b>Mathematical Optimization</b> . . . . .	48
2.4.1	General Formulation . . . . .	49
2.4.2	Optimality Conditions . . . . .	50
2.4.3	Numerical Optimization Methods . . . . .	59
	Exercises . . . . .	64
3	<b>Supervised Machine Learning (in a Nutshell)</b>	67
3.1	Overview . . . . .	67
3.2	Case Studies . . . . .	72
4	<b>Feature Extraction</b>	77
4.1	<b>Feature Extraction: Concepts</b> . . . . .	77
4.1.1	Feature Engineering . . . . .	77
4.1.2	Feature Selection . . . . .	78
4.1.3	Dimensionality Reduction . . . . .	79
4.2	<b>Linear Dimension Reduction</b> . . . . .	79
4.2.1	Principal Component Analysis . . . . .	80
4.2.2	Linear Discriminant Analysis . . . . .	84
4.3	<b>Nonlinear Dimension Reduction (I): Manifold Learning</b> . . . . .	86
4.3.1	Locally Linear Embedding . . . . .	87
4.3.2	Multidimensional Scaling . . . . .	88
4.3.3	Stochastic Neighborhood Embedding . . . . .	89
4.4	<b>Nonlinear Dimension Reduction (II): Neural Networks</b> . . . . .	90
4.4.1	Autoencoder . . . . .	90
4.4.2	Bottleneck Features . . . . .	91
	Lab Project I . . . . .	92
	Exercises . . . . .	93

<b>DISCRIMINATIVE MODELS</b>	<b>95</b>
<b>5 Statistical Learning Theory</b>	<b>97</b>
5.1 Formulation of Discriminative Models . . . . .	97
5.2 Learnability . . . . .	99
5.3 Generalization Bounds . . . . .	100
5.3.1 Finite Model Space: $ \mathcal{H} $ . . . . .	100
5.3.2 Infinite Model Space: VC Dimension . . . . .	102
Exercises . . . . .	105
<b>6 Linear Models</b>	<b>107</b>
6.1 Perceptron . . . . .	108
6.2 Linear Regression . . . . .	112
6.3 Minimum Classification Error . . . . .	113
6.4 Logistic Regression . . . . .	114
6.5 Support Vector Machines . . . . .	116
6.5.1 Linear SVM . . . . .	116
6.5.2 Soft SVM . . . . .	121
6.5.3 Nonlinear SVM: The Kernel Trick . . . . .	123
6.5.4 Solving Quadratic Programming . . . . .	126
6.5.5 Multiclass SVM . . . . .	127
Lab Project II . . . . .	129
Exercises . . . . .	130
<b>7 Learning Discriminative Models in General</b>	<b>133</b>
7.1 A General Framework to Learn Discriminative Models . . . . .	133
7.1.1 Common Loss Functions in Machine Learning . . . . .	135
7.1.2 Regularization Based on $L_p$ Norm . . . . .	136
7.2 Ridge Regression and LASSO . . . . .	139
7.3 Matrix Factorization . . . . .	140
7.4 Dictionary Learning . . . . .	145
Lab Project III . . . . .	149
Exercises . . . . .	150
<b>8 Neural Networks</b>	<b>151</b>
8.1 Artificial Neural Networks . . . . .	152
8.1.1 Basic Formulation of Artificial Neural Networks . . . . .	152
8.1.2 Mathematical Justification: Universal Approximator . . . . .	154
8.2 Neural Network Structures . . . . .	156
8.2.1 Basic Building Blocks to Connect Layers . . . . .	156
8.2.2 Case Study I: Fully Connected Deep Neural Networks . . . . .	165
8.2.3 Case Study II: Convolutional Neural Networks . . . . .	166
8.2.4 Case Study III: Recurrent Neural Networks (RNNs) . . . . .	170

8.2.5	Case Study IV: Transformer . . . . .	172
8.3	<b>Learning Algorithms for Neural Networks</b> . . . . .	174
8.3.1	Loss Function . . . . .	175
8.3.2	Automatic Differentiation . . . . .	176
8.3.3	Optimization Using Stochastic Gradient Descent . . . . .	188
8.4	<b>Heuristics and Tricks for Optimization</b> . . . . .	189
8.4.1	Other SGD Variant Optimization Methods: ADAM . . . . .	192
8.4.2	Regularization . . . . .	194
8.4.3	Fine-Tuning Tricks . . . . .	196
8.5	<b>End-to-End Learning</b> . . . . .	197
8.5.1	Sequence-to-Sequence Learning . . . . .	198
	Lab Project IV . . . . .	200
	Exercises . . . . .	201
9	<b>Ensemble Learning</b> . . . . .	203
9.1	<b>Formulation of Ensemble Learning</b> . . . . .	203
9.1.1	Decision Trees . . . . .	205
9.2	<b>Bagging</b> . . . . .	208
9.2.1	Random Forests . . . . .	208
9.3	<b>Boosting</b> . . . . .	209
9.3.1	Gradient Boosting . . . . .	210
9.3.2	AdaBoost . . . . .	212
9.3.3	Gradient Tree Boosting . . . . .	214
	Lab Project V . . . . .	216
	Exercises . . . . .	217
	<b>GENERATIVE MODELS</b> . . . . .	219
10	<b>Overview of Generative Models</b> . . . . .	221
10.1	<b>Formulation of Generative Models</b> . . . . .	221
10.2	<b>Bayesian Decision Theory</b> . . . . .	222
10.2.1	Generative Models for Classification . . . . .	223
10.2.2	Generative Models for Regression . . . . .	227
10.3	<b>Statistical Data Modeling</b> . . . . .	228
10.3.1	Plug-In MAP Decision Rule . . . . .	229
10.4	<b>Density Estimation</b> . . . . .	231
10.4.1	Maximum-Likelihood Estimation . . . . .	231
10.4.2	Maximum-Likelihood Classifier . . . . .	234
10.5	<b>Generative Models (in a Nutshell)</b> . . . . .	234
10.5.1	Generative versus Discriminative Models . . . . .	236
	Exercises . . . . .	237



<b>11 Unimodal Models</b>	<b>239</b>
11.1 Gaussian Models . . . . .	240
11.2 Multinomial Models . . . . .	243
11.3 Markov Chain Models . . . . .	245
11.4 Generalized Linear Models . . . . .	250
11.4.1 Probit Regression . . . . .	252
11.4.2 Poisson Regression . . . . .	252
11.4.3 Log-Linear Models . . . . .	253
Exercises . . . . .	256
<b>12 Mixture Models</b>	<b>257</b>
12.1 Formulation of Mixture Models . . . . .	257
12.1.1 Exponential Family (e-Family) . . . . .	259
12.1.2 Formal Definition of Mixture Models . . . . .	261
12.2 Expectation-Maximization Method . . . . .	261
12.2.1 Auxiliary Function: Eliminating Log-Sum . . . . .	262
12.2.2 Expectation-Maximization Algorithm . . . . .	265
12.3 Gaussian Mixture Models . . . . .	268
12.3.1 K-Means Clustering for Initialization . . . . .	270
12.4 Hidden Markov Models . . . . .	271
12.4.1 HMMs: Mixture Models for Sequences . . . . .	272
12.4.2 Evaluation Problem: Forward–Backward Algorithm . . . . .	276
12.4.3 Decoding Problem: Viterbi Algorithm . . . . .	279
12.4.4 Training Problem: Baum–Welch Algorithm . . . . .	280
Lab Project VI . . . . .	287
Exercises . . . . .	288
<b>13 Entangled Models</b>	<b>291</b>
13.1 Formulation of Entangled Models . . . . .	291
13.1.1 Framework of Entangled Models . . . . .	292
13.1.2 Learning of Entangled Models in General . . . . .	294
13.2 Linear Gaussian Models . . . . .	296
13.2.1 Probabilistic PCA . . . . .	296
13.2.2 Factor Analysis . . . . .	298
13.3 Non-Gaussian Models . . . . .	300
13.3.1 Independent Component Analysis (ICA) . . . . .	300
13.3.2 Independent Factor Analysis (IFA) . . . . .	301
13.3.3 Hybrid Orthogonal Projection and Estimation (HOPE) . . . . .	302
13.4 Deep Generative Models . . . . .	303
13.4.1 Variational Autoencoders (VAE) . . . . .	304
13.4.2 Generative Adversarial Nets (GAN) . . . . .	307
Exercises . . . . .	309

<b>14 Bayesian Learning</b>	<b>311</b>
14.1 <b>Formulation of Bayesian Learning</b>	311
14.1.1 Bayesian Inference	313
14.1.2 Maximum a Posterior Estimation	314
14.1.3 Sequential Bayesian Learning	315
14.2 <b>Conjugate Priors</b>	318
14.2.1 Maximum-Marginal-Likelihood Estimation	323
14.3 <b>Approximate Inference</b>	324
14.3.1 Laplace's Method	324
14.3.2 Variational Bayesian (VB) Methods	326
14.4 <b>Gaussian Processes</b>	332
14.4.1 Gaussian Processes as Nonparametric Priors	333
14.4.2 Gaussian Processes for Regression	335
14.4.3 Gaussian Processes for Classification	338
Exercises	340
<b>15 Graphical Models</b>	<b>343</b>
15.1 <b>Concepts of Graphical Models</b>	343
15.2 <b>Bayesian Networks</b>	346
15.2.1 Conditional Independence	346
15.2.2 Representing Generative Models as Bayesian Networks	351
15.2.3 Learning Bayesian Networks	353
15.2.4 Inference Algorithms	355
15.2.5 Case Study I: Naive Bayes Classifier	361
15.2.6 Case Study II: Latent Dirichlet Allocation	362
15.3 <b>Markov Random Fields</b>	366
15.3.1 Formulation: Potential and Partition Functions	366
15.3.2 Case Study III: Conditional Random Fields	368
15.3.3 Case Study IV: Restricted Boltzmann Machines	370
Exercises	372
<b>APPENDIX</b>	<b>375</b>
<b>A Other Probability Distributions</b>	<b>377</b>
<b>Bibliography</b>	<b>381</b>
<b>Index</b>	<b>397</b>

## Preface

Machine learning used to be a niche area originating out of pattern classification in electrical engineering and artificial intelligence in computer science. Today, machine learning has grown into a very diverse discipline spanning a variety of topics in mathematics, science, and engineering. Because of the widespread use and increased power of computers, machine learning has found a plethora of relevant applications in almost all engineering domains and has made a huge impact on our society. In particular, with the boom of deep learning in recent years, thousands of new researchers and practitioners across academia and industry join forces every year to tackle machine learning and its applications. In many universities, machine learning has become one of the most popular advanced elective courses, highly demanded by senior undergraduates and graduates in almost all computer science and electrical engineering programs. The number of industrial job positions in machine learning, deep learning, and data science has dramatically increased in recent years, and this trend is expected to continue for at least the next 10 years due to the availability of a huge amount of data over the internet and personal devices.

### Why This Book?

There are already plenty of well-written textbooks for machine learning, most of which exhaustively cover a wide range of topics in machine learning. In teaching my machine learning courses, I found that they are too challenging for beginners because of the vast range of presented topics and the overwhelming technical details associated with them. Many beginners have trouble with the heavy mathematical notation and equations, whereas others drown in all the technical details and fail to grasp the essence of these machine learning methods.

In contrast, this book is intended to present the fundamental machine learning concepts, algorithms, and principles in a concise and lucid manner, without heavy mathematical machinery and excess detail. I have been selective in terms of the topics so that it can all be covered in an introductory course, rather than making it comprehensive enough to cover all machine learning topics. I chose to cover only relatively mature topics primarily related to supervised learning, which I believe are not only fundamental to the field of machine learning but also significant enough to have made an impact in both academia and industry. In other words, some satisfactory and feasible solutions have already been developed for these topics so that they are able to address not just toy problems

but many interesting problems arising in the real world. At the same time, I have tried to omit many minor issues surrounding the central topics so that beginners will not be distracted by these purely technical details.

Instead of covering the selected topics separately, one after another, I have tried to organize all machine learning topics into a coherent structure to give readers a big picture of the entire field. All topics are arranged into coherent groups, and the individual chapters are dedicated to covering all logically relevant methods in each group. After reading each chapter, readers can immediately understand the differences between them, grasp their relevance, and also know how these methods fit into the big picture of machine learning.

This book also aims to reflect the latest advancements in the field. I have included significant coverage on several important recent techniques, such as *transformers*, which have come to dominate many natural-language-processing tasks; *batch norm* and *ADAM optimization*, which are popular in learning large and deep neural networks; and recently popular deep generative models such as *variational autoencoders (VAEs)* and *generative adversarial nets (GANs)*.

For all topics in this book, I provide enough technical depth to explain the motivation, principles, and methodology in a professional manner. As much as possible, I derive the machine learning methods from scratch using rigorous mathematics to highlight the core ideas behind them. For critical theoretical results, I have included many important theorems and some light proofs. The important mathematical topics and methods that modern machine learning methods are built on are thoroughly reviewed in Chapter 2. However, readers do need a good background in *calculus*, *linear algebra*, and *probability and statistics* to be able to follow the descriptions and discussions in this book. Throughout the book, I have also done my best to present all technical content using clean and consistent mathematical notations and represent all algorithms in this book as concise linear algebra formulas, which can be translated almost line by line into efficient code using a programming language supporting vectorization, such as MATLAB or Python.

### Whom Is This Book For?

This book is primarily written as a textbook for an introductory course on machine learning for senior undergraduate students in computer science and computer/ software/electrical engineering programs or first-year graduate students in many science, engineering, and applied mathematics programs who are interested in basic machine learning methods for their own research problems. I also hope it will be useful as a self-study or reference book for researchers who wish to apply machine learning methods to solve their own problems, as well

as industrial practitioners who want to understand the concepts and principles behind the popular machine learning methods they implement. Given the large number of machine learning software programs and toolkits freely available today, it is often not hard to write code to run fairly complicated machine learning algorithms. However, in many cases, knowledge of the principles and mathematics behind these algorithms is required to tune these algorithms in order to deliver optimal results for the task at hand.

### Online Resources

This book is accompanied by the following GitHub repository:

<https://github.com/iNCML/MachineLearningBook>

This website provides a variety of supplementary materials to support this book, including the following:

- ▶ Lecture slides per chapter
- ▶ Code samples for some lab projects (MATLAB or Python)

Meanwhile, readers and instructors can also provide their feedback, suggestions, and comments on this book as *issues* through the GitHub repository. I will reply to these requests as much as possible.

### How to Use This Book

I have made much effort to keep this book succinct and only cover the most important issues for each selected topic. I encourage readers to read all chapters in order because I have tried my best to arrange a wide range of machine learning topics in a coherent structure. For each machine learning method, I have thoroughly covered the motivation, main ideas, concepts, methodology, and algorithms in the main text and sometimes have left extensive issues and extra technical details or extensions as chapter-end exercises. Readers may optionally follow these links to work on these exercises and practice the main ideas discussed in the text.

#### ▶ For a Semester-Long Course

Instructors may use this book as the primary or alternate textbook for a standard semester-long introductory course (about 10–12 weeks) on machine learning in the fourth year of a computer science, engineering, or applied mathematics program. I suggest covering the following topics in order:

- Chapter 1: Introduction (0.5 week)
- Chapter 2: Mathematical Foundation (1.5 weeks)
- Chapter 4: Feature Extraction (1 week)
- Chapter 5: Statistical Learning Theory (0.5 week)
  - §5.1 Formulation of Discriminative Models
  - §5.2 Learnability
- Chapter 6: Linear Models (1.5 weeks)
- Chapter 7: Learning Discriminative Models (1 week)
  - §7.1 General Framework
  - §7.2 Ridge and LASSO
  - §7.3 Matrix Factorization
- Chapter 8: Neural Networks (2 weeks)
- Chapter 9: Ensemble Learning (1 week)
- Chapter 10: Overview of Generative Models (1 week)
- Chapter 11: Unimodal Models (1 week)
  - §11.1 Gaussian Models
  - §11.2 Multinomial Models
  - §11.3 Markov Chain Models
- Chapter 12 Mixture Models (1 week)
  - §12.1 Formulation
  - §12.2 EM Method
  - §12.3 Gaussian Mixture Models

► **For a Year-Long Full Course**

Instructors may also use this book as the primary or alternate textbook for a year-long full course on machine learning (20–24 weeks) to give balanced coverage of both discriminative and generative models. The first half focuses on the mathematical preparation and discriminative models, whereas the second half gives full exposure to a variety of topics in generative models, including Chapter 13: Entangled Models, Chapter 14: Bayesian Learning, and Chapter 15: Graphical Models.

If time is tight, instructors may skip some optional topics, such as §4.3 Manifold Learning, §7.4 Dictionary Learning, §11.4 Generalized Linear Models, §12.4 Hidden Markov Models, or §14.4 Gaussian Processes.

► **For Self-Study**

All self-study readers are strongly recommended to go through the book in order. This will give a smooth transition from one topic to another, generally progressing gradually from easy topics to hard ones. Depending on one's own interests, readers may choose to skip any of the following advanced topics without affecting the understanding of other parts:

- §4.3 Manifold Learning
- §7.4 Dictionary Learning
- §11.4 Generalized Linear Models
- §12.4 Hidden Markov Models
- §14.4 Gaussian Processes

### Acknowledgments

Writing a textbook is a very challenging task. This book would not have been possible without help and supports from a large number of people.

Most content in this book evolved from the lecture notes I have used for many years to teach a machine learning course in the Department of Electrical Engineering and Computer Science at York University in Toronto, Canada. I am grateful to York University for the long-standing support of my teaching and research there.

I also thank Zoubin Ghahramani, David Blei, and Huy Vu for granting permission to use their materials in this book.

Many people have helped to significantly improve this book by proofreading the early draft and providing valuable comments and suggestions, including Dong Yu, Kelvin Jiang, Behnam Asadi, Jia Pan, Yong Ge, William Fu, Xiaodan Zhu, Chao Wang, Jiebo Luo, Hanjia Lyu, Joyce Luo, Qiang Huo, Chunxiao Zhou, Wei Zhang, Maria Koshkina, Zhuoran Li, Junfei Wang, and Parham Eftekhari. My special thanks to all of them!

Finally, I would like to thank my family, Iris and Kelvin, and my parents for their endless support and love throughout the time of writing this book as well as my career and life.

## Notation

This list describes some of the symbols that are used within this book.

$\mu$	The mean vector of a multivariate Gaussian
$\Sigma$	The covariance matrix of a multivariate Gaussian
$E[\cdot]$	The expectation or the mean
$E_X[\cdot]$	The expectation with respect to $X$
$\mathcal{H}$	Model space
$\mathbb{N}$	The set of natural numbers
$\mathbb{R}$	The set of real numbers
$\mathbb{R}^n$	The set of $n$ -dimensional real vectors
$\mathbb{R}^{m \times n}$	The set of $m \times n$ real matrices
$\mathcal{W}$	The set of all parameters in a neural network
$\mathbf{S}$	The sample covariance matrix
$\mathbf{w} * \mathbf{x}$	The convolution sum of $\mathbf{w}$ and $\mathbf{x}$
$\mathbf{w} \cdot \mathbf{x}$	The inner product of two vectors $\mathbf{w}$ and $\mathbf{x}$
$\mathbf{w} \odot \mathbf{x}$	The element-wise multiplication of $\mathbf{w}$ and $\mathbf{x}$
$\mathbf{W}$	A weight matrix
$\mathbf{w}$	A weight vector
$\mathbf{x}$	A feature vector
$\nabla f(\mathbf{x})$	The gradient of a function $f(\mathbf{x})$
$\Pr(A)$	The probability of an event $A$
$\ \mathbf{w}\ $	The norm (or $L_2$ norm) of a vector $\mathbf{w}$
$\ \mathbf{w}\ _p$	The $L_p$ norm of a vector $\mathbf{w}$
$f(\mathbf{x}; \theta)$	A function of $\mathbf{x}$ with the parameter $\theta$
$f_\theta(\mathbf{x})$	A function of $\mathbf{x}$ with the parameter $\theta$
$l(\theta)$	A log-likelihood function of the model parameter $\theta$
$m \ll n$	$m$ is much less than $n$
$p(\mathbf{x}, \mathbf{y})$	A joint distribution of $\mathbf{x}$ and $\mathbf{y}$
$p(\mathbf{y}   \mathbf{x})$	A conditional distribution of $\mathbf{y}$ given $\mathbf{x}$
$p_\theta(\mathbf{x})$	A probability distribution of $\mathbf{x}$ with the parameter $\theta$
$Q(\mathcal{W}; \mathbf{x})$	An objective function of the model parameters $\mathcal{W}$ given the data $\mathbf{x}$
$\theta$	Model parameter



## Summary of the General Notation Rules

Notation	Meaning	Examples
Lowercase letters	A scalar	$x, y, n, m, x_i, x_{ij}$
	A function	$f(\cdot), p(\cdot), g(\cdot), h(\cdot)$
Lowercase letters in bold	A column vector	$\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{a}, \mathbf{b}$ $\boldsymbol{\mu}, \boldsymbol{\nu}$
Uppercase letters	A random variable	$X, Y, X_i, X_j$
	A function	$Q(\cdot), \Phi(\cdot, \cdot)$
Uppercase letters in bold	A matrix	$\mathbf{A}, \mathbf{W}, \mathbf{S}$ $\boldsymbol{\Sigma}, \boldsymbol{\Phi}$
Uppercase letters in blackboard bold	A set of numbers	$\mathbb{N}, \mathbb{R}$
	A set of parameters	$\mathbb{B}, \mathbb{W}, \mathbb{V}$
Uppercase letters in calligraphy	A set of data	$\mathcal{D}, \mathcal{D}_N$