

LEARNING MADE EASY



2nd Edition

# Machine Learning

for  
**dummies**<sup>®</sup>  
A Wiley Brand



Perform in-depth analysis  
to extract interesting details

Learn how machine learning  
algorithms are invaluable

Implement algorithms in  
Python<sup>®</sup> and TensorFlow<sup>®</sup>

**John Paul Mueller**  
**Luca Massaron**

Bestselling authors of the first edition





# Machine Learning

2nd Edition

**by John Paul Mueller and Luca Massaron**

**for  
dummies®**  
A Wiley Brand

## Machine Learning For Dummies®, 2nd Edition

Published by: **John Wiley & Sons, Inc.**, 111 River Street, Hoboken, NJ 07030-5774, [www.wiley.com](http://www.wiley.com)

Copyright © 2021 by John Wiley & Sons, Inc., Hoboken, New Jersey

Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without the prior written permission of the Publisher. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

**Trademarks:** Wiley, For Dummies, the Dummies Man logo, Dummies.com, Making Everything Easier, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and may not be used without written permission. All other trademarks are the property of their respective owners. John Wiley & Sons, Inc. is not associated with any product or vendor mentioned in this book.

LIMIT OF LIABILITY/DISCLAIMER OF WARRANTY: THE PUBLISHER AND THE AUTHOR MAKE NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE ACCURACY OR COMPLETENESS OF THE CONTENTS OF THIS WORK AND SPECIFICALLY DISCLAIM ALL WARRANTIES, INCLUDING WITHOUT LIMITATION WARRANTIES OF FITNESS FOR A PARTICULAR PURPOSE. NO WARRANTY MAY BE CREATED OR EXTENDED BY SALES OR PROMOTIONAL MATERIALS. THE ADVICE AND STRATEGIES CONTAINED HEREIN MAY NOT BE SUITABLE FOR EVERY SITUATION. THIS WORK IS SOLD WITH THE UNDERSTANDING THAT THE PUBLISHER IS NOT ENGAGED IN RENDERING LEGAL, ACCOUNTING, OR OTHER PROFESSIONAL SERVICES. IF PROFESSIONAL ASSISTANCE IS REQUIRED, THE SERVICES OF A COMPETENT PROFESSIONAL PERSON SHOULD BE SOUGHT. NEITHER THE PUBLISHER NOR THE AUTHOR SHALL BE LIABLE FOR DAMAGES ARISING HEREFROM. THE FACT THAT AN ORGANIZATION OR WEBSITE IS REFERRED TO IN THIS WORK AS A CITATION AND/OR A POTENTIAL SOURCE OF FURTHER INFORMATION DOES NOT MEAN THAT THE AUTHOR OR THE PUBLISHER ENDORSES THE INFORMATION THE ORGANIZATION OR WEBSITE MAY PROVIDE OR RECOMMENDATIONS IT MAY MAKE. FURTHER, READERS SHOULD BE AWARE THAT INTERNET WEBSITES LISTED IN THIS WORK MAY HAVE CHANGED OR DISAPPEARED BETWEEN WHEN THIS WORK WAS WRITTEN AND WHEN IT IS READ.

For general information on our other products and services, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002. For technical support, please visit <https://hub.wiley.com/community/support/dummies>.

Wiley publishes in a variety of print and electronic formats and by print-on-demand. Some material included with standard print versions of this book may not be included in e-books or in print-on-demand. If this book refers to media such as a CD or DVD that is not included in the version you purchased, you may download this material at <http://booksupport.wiley.com>. For more information about Wiley products, visit [www.wiley.com](http://www.wiley.com).

Library of Congress Control Number: 2020952332

ISBN: 978-1-119-72401-8

ISBN 978-1-119-72406-3 (ebk); ISBN 978-1-119-72405-6 (ebk)

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1

# Contents at a Glance

<b>Introduction</b>	1
<b>Part 1: Introducing How Machines Learn</b>	5
CHAPTER 1: Getting the Real Story about AI	7
CHAPTER 2: Learning in the Age of Big Data	23
CHAPTER 3: Having a Glance at the Future	37
<b>Part 2: Preparing Your Learning Tools</b>	47
CHAPTER 4: Installing a Python Distribution	49
CHAPTER 5: Beyond Basic Coding in Python	67
CHAPTER 6: Working with Google Colab	87
<b>Part 3: Getting Started with the Math Basics</b>	115
CHAPTER 7: Demystifying the Math Behind Machine Learning	117
CHAPTER 8: Descending the Gradient	139
CHAPTER 9: Validating Machine Learning	153
CHAPTER 10: Starting with Simple Learners	175
<b>Part 4: Learning from Smart and Big Data</b>	197
CHAPTER 11: Preprocessing Data	199
CHAPTER 12: Leveraging Similarity	221
CHAPTER 13: Working with Linear Models the Easy Way	243
CHAPTER 14: Hitting Complexity with Neural Networks	271
CHAPTER 15: Going a Step Beyond Using Support Vector Machines	307
CHAPTER 16: Resorting to Ensembles of Learners	319
<b>Part 5: Applying Learning to Real Problems</b>	339
CHAPTER 17: Classifying Images	341
CHAPTER 18: Scoring Opinions and Sentiments	361
CHAPTER 19: Recommending Products and Movies	383
<b>Part 6: The Part of Tens</b>	405
CHAPTER 20: Ten Ways to Improve Your Machine Learning Models	407
CHAPTER 21: Ten Guidelines for Ethical Data Usage	415
CHAPTER 22: Ten Machine Learning Packages to Master	423
<b>Index</b>	431



# Table of Contents

<b>INTRODUCTION</b>	<b>1</b>
About This Book	1
Foolish Assumptions	2
Icons Used in This Book	3
Beyond the Book	3
Where to Go from Here	4
 <b>PART 1: INTRODUCING HOW MACHINES LEARN</b>	 <b>5</b>
<b>CHAPTER 1: Getting the Real Story about AI</b>	<b>7</b>
Moving beyond the Hype	8
Dreaming of Electric Sheep	9
Understanding the history of AI and machine learning	10
Exploring what machine learning can do for AI	11
Considering the goals of machine learning	12
Defining machine learning limits based on hardware	12
Overcoming AI Fantasies	13
Discovering the fad uses of AI and machine learning	14
Considering the true uses of AI and machine learning	15
Being useful; being mundane	16
Considering the Relationship between AI and Machine Learning	17
Considering AI and Machine Learning Specifications	18
Defining the Divide between Art and Engineering	19
Predicting the Next AI Winter	20
 <b>CHAPTER 2: Learning in the Age of Big Data</b>	 <b>23</b>
Considering the Machine Learning Essentials	24
Defining Big Data	25
Considering the Sources of Big Data	26
Building a new data source	26
Using existing data sources	29
Locating test data sources	29
Specifying the Role of Statistics in Machine Learning	30
Understanding the Role of Algorithms	31
Defining what algorithms do	32
Considering the five main techniques	32
Defining What Training Means	34

<b>CHAPTER 3: Having a Glance at the Future</b>	37
Creating Useful Technologies for the Future	38
Considering the role of machine learning in robots.	38
Using machine learning in health care.	39
Creating smart systems for various needs	40
Using machine learning in industrial settings.	40
Understanding the role of updated processors and other hardware	41
Discovering the New Work Opportunities with Machine Learning	42
Working for a machine	42
Working with machines	43
Repairing machines.	44
Creating new machine learning tasks.	44
Devising new machine learning environments.	45
Avoiding the Potential Pitfalls of Future Technologies	46
 <b>PART 2: PREPARING YOUR LEARNING TOOLS</b>	 47
<b>CHAPTER 4: Installing a Python Distribution</b>	49
Using Anaconda for Machine Learning	50
Getting Anaconda	50
Defining why Anaconda is used in this book.	51
Installing Anaconda on Linux.	52
Installing Anaconda on Mac OS X	53
Installing Anaconda on Windows	54
Downloading the Datasets and Example Code.	57
Using Jupyter Notebook	57
Defining the code repository.	59
Understanding the datasets used in this book.	64
 <b>CHAPTER 5: Beyond Basic Coding in Python</b>	 67
Defining the Basics You Should Know	68
Considering Python basics.	68
Working with functions.	72
Working with modules	76
Storing Data Using Sets, Lists, and Tuples.	78
Creating sets.	78
Performing operations on sets	78
Using lists	79
Creating and using tuples	82
Defining Useful Iterators	83
Working with ranges.	83
Iterating multiple lists using zip.	84
Working with generators using yield	84



Indexing Data Using Dictionaries .....	85
Creating dictionaries.....	85
Storing and retrieving data from dictionaries.....	85
<b>CHAPTER 6: Working with Google Colab.....</b>	<b>87</b>
Defining Google Colab .....	88
Understanding what Google Colab does.....	88
Considering the online coding difference .....	90
Using local runtime support .....	91
Working with Google Colab features .....	91
Getting a Google Account.....	94
Creating the account.....	94
Signing in .....	95
Working with Notebooks .....	96
Creating a new notebook.....	96
Opening existing notebooks .....	97
Uploading a notebook .....	99
Saving notebooks .....	100
Downloading notebooks .....	103
Performing Common Tasks .....	103
Creating code cells .....	104
Creating text cells .....	106
Creating special cells.....	107
Editing cells.....	108
Moving cells .....	108
Using Hardware Acceleration .....	108
Viewing Your Notebook .....	109
Displaying the table of contents .....	110
Getting notebook information.....	110
Checking code execution .....	110
Executing the Code .....	111
Sharing Your Notebook .....	112
Getting Help .....	113
<b>PART 3: GETTING STARTED WITH THE MATH BASICS.....</b>	<b>115</b>
<b>CHAPTER 7: Demystifying the Math Behind Machine Learning .....</b>	<b>117</b>
Working with Data.....	118
Learning the terminology.....	119
Understanding scalar and vector operations .....	120
Performing vector multiplication .....	123
Creating a matrix.....	123
Understanding basic operations.....	125

Performing matrix multiplication . . . . .	126
Glancing at advanced matrix operations . . . . .	128
Using vectorization effectively . . . . .	129
Exploring the World of Probabilities . . . . .	130
Getting an overview of probability . . . . .	130
Operating on probabilities . . . . .	131
Conditioning chance by Bayes' theorem . . . . .	132
Describing the Use of Statistics . . . . .	135
<b>CHAPTER 8: Descending the Gradient . . . . .</b>	<b>139</b>
Acknowledging Different Kinds of Learning . . . . .	140
Supervised learning . . . . .	140
Unsupervised learning . . . . .	141
Reinforcement learning . . . . .	141
The learning process . . . . .	142
Mapping an unknown function . . . . .	142
Exploring cost functions . . . . .	145
Descending the optimization curve . . . . .	147
Optimizing with big data . . . . .	148
Leveraging sampling . . . . .	149
Using parallelism . . . . .	150
Learning out-of-core . . . . .	151
<b>CHAPTER 9: Validating Machine Learning . . . . .</b>	<b>153</b>
Considering the Use of Example Data . . . . .	154
Checking Out-of-Sample Errors . . . . .	155
Understanding the concept of samples . . . . .	155
Looking for the holy grail of generalization . . . . .	156
Experimenting how bias and variance work . . . . .	158
Keeping model complexity in mind . . . . .	160
Keeping solutions balanced . . . . .	162
Depicting learning curves . . . . .	163
Training, Validating, and Testing . . . . .	165
Considering the split . . . . .	165
Resorting to cross-validation . . . . .	166
Looking for alternatives in validation . . . . .	167
Optimizing by Cross-Validation . . . . .	169
Sources of predictive performance . . . . .	169
Exploring the hyper-parameter space . . . . .	170
Selecting relevant features . . . . .	171
Avoiding Sample Bias and Leakage Traps . . . . .	173
<b>CHAPTER 10: Starting with Simple Learners . . . . .</b>	<b>175</b>
Discovering the Incredible Perceptron . . . . .	176
Falling short of a miracle . . . . .	176
Hitting the nonseparability limit . . . . .	179

Growing Greedy Classification Trees .....	180
Predicting outcomes by splitting data .....	181
Pruning overgrown trees .....	185
Taking a Probabilistic Turn .....	188
Understanding Naïve Bayes .....	189
Estimating response with Naïve Bayes .....	192
<b>PART 4: LEARNING FROM SMART AND BIG DATA .....</b>	<b>197</b>
<b>CHAPTER 11: Preprocessing Data .....</b>	<b>199</b>
Gathering and Cleaning Data .....	200
Repairing Missing Data .....	201
Identifying missing data .....	201
Choosing the right replacement strategy .....	203
Transforming Distributions .....	205
Creating Your Own Features .....	207
Understanding the need to create features .....	207
Creating features automatically .....	208
Explaining the basics of SVD .....	210
Reorganizing data .....	212
Delimiting Anomalous Data .....	215
Using a univariate strategy .....	215
Resorting to Multivariate Models .....	217
<b>CHAPTER 12: Leveraging Similarity .....</b>	<b>221</b>
Measuring Similarity between Vectors .....	222
Understanding similarity .....	222
Computing distances for learning .....	223
Using Distances to Locate Clusters .....	224
Checking assumptions and expectations .....	226
Inspecting the gears of the K-means algorithm .....	227
Tuning the K-Means Algorithm .....	229
Experimenting with K-means reliability .....	230
Experimenting with how centroids converge .....	233
Finding Similarity by K-Nearest Neighbors .....	238
Understanding the k parameter .....	238
Experimenting with a flexible algorithm .....	240
<b>CHAPTER 13: Working with Linear Models the Easy Way .....</b>	<b>243</b>
Starting to Combine Features .....	244
Getting an overview of regression .....	244
Solving problems with a machine learning approach .....	247
Understanding R squared .....	249
Mixing Features of Different Types .....	251
Switching to Probabilities .....	255

Specifying a binary response.....	255
Handling multiple classes.....	259
Guessing the Right Features .....	259
Defining the outcome of features that don't work together.....	259
Solving overfitting by using greedy selection .....	260
Addressing overfitting by regularization .....	262
Learning One Example at a Time .....	264
Using gradient descent.....	264
Understanding how SGD is different .....	265
<b>CHAPTER 14: Hitting Complexity with Neural Networks .....</b>	<b>271</b>
Revising the Perceptron .....	272
Pushing forth with feed-forward.....	274
Going even deeper down the rabbit hole .....	276
Pulling back with backpropagation.....	280
Representing the Way of Learning of a Network .....	283
Understanding the problem with overfitting .....	283
Choosing a framework .....	285
Getting your copy of TensorFlow and Keras .....	286
Opening the black box .....	289
Introducing Deep Learning .....	294
Understanding some deep learning essentials.....	295
Explaining the magic of convolutions.....	296
Understanding recurrent neural networks .....	300
<b>CHAPTER 15: Going a Step Beyond Using Support Vector Machines .....</b>	<b>307</b>
Revisiting the Separation Problem .....	308
Explaining the Algorithm .....	309
Avoiding the pitfalls of nonseparability .....	311
Applying nonlinearity .....	312
Explaining the kernel trick by example .....	314
Classifying and Estimating with SVM .....	316
<b>CHAPTER 16: Resorting to Ensembles of Learners .....</b>	<b>319</b>
Leveraging Decision Trees .....	320
Growing a forest of trees .....	322
Understanding the importance measures.....	327
Working with Almost Random Guesses.....	330
Bagging predictors with Adaboost .....	330
Boosting Smart Predictors.....	333
Meeting again with gradient descent.....	334
Considering the state of the art in tabular data .....	335

Averaging Different Predictors .....	336
Blending solutions.....	337
Stacking diverse solutions .....	337
<b>PART 5: APPLYING LEARNING TO REAL PROBLEMS.....</b>	<b>339</b>
<b>CHAPTER 17: Classifying Images.....</b>	<b>341</b>
Working with a Set of Images .....	342
Revising the State of the Art in Computer Vision .....	347
Extracting Visual Features .....	350
Recognizing Faces Using Eigenfaces.....	352
Classifying Images .....	356
<b>CHAPTER 18: Scoring Opinions and Sentiments .....</b>	<b>361</b>
Introducing Natural Language Processing.....	362
Revising the State of the Art in NLP .....	363
Understanding How Machines Read .....	364
Defining the input data. ....	364
Processing and enhancing text .....	366
Scraping textual datasets from the web .....	371
Handling problems with raw text .....	374
Using Scoring and Classification .....	375
Performing classification tasks .....	375
Analyzing reviews from e-commerce .....	378
<b>CHAPTER 19: Recommending Products and Movies .....</b>	<b>383</b>
Realizing the Revolution of E-Commerce.....	384
Downloading Rating Data.....	386
Trudging through the MovieLens dataset .....	386
Navigating through anonymous web data .....	390
Encountering the limits of rating data .....	392
Considering collaborative filtering .....	392
Catching the Limits of Behavioral Data .....	397
Integrating Text and Behaviors.....	399
Viewing the attributes.....	399
Obtaining statistics .....	400
Leveraging SVD .....	400
Understanding the SVD connection .....	400
Seeing SVD in action .....	401
<b>PART 6: THE PART OF TENS.....</b>	<b>405</b>
<b>CHAPTER 20: Ten Ways to Improve Your Machine Learning Models .....</b>	<b>407</b>
Studying Learning Curves.....	408
Using Cross-Validation Correctly.....	409

Choosing the Right Error or Score Metric .....	410
Searching for the Best Hyper-Parameters.....	410
Testing Multiple Models.....	411
Averaging Models .....	411
Stacking Models.....	412
Applying Feature Engineering .....	412
Selecting Features and Examples .....	413
Looking for More Data .....	414
<b>CHAPTER 21: Ten Guidelines for Ethical Data Usage.....</b>	<b>415</b>
Obtaining Permission .....	416
Using Sanitization Techniques.....	417
Avoiding Data Inference.....	418
Using Generalizations Correctly .....	418
Shunning Discriminatory Practices.....	419
Detecting Black Swans in Code .....	420
Understanding the Process .....	420
Considering the Consequences of an Action.....	421
Balancing Decision Making .....	421
Verifying a Data Source .....	422
<b>CHAPTER 22: Ten Machine Learning Packages to Master.....</b>	<b>423</b>
Gensim .....	423
imbalanced-learn.....	424
OpenCV.....	424
SciPy .....	425
SHAP .....	426
Statsmodels .....	427
Modin .....	427
PyTorch.....	428
Poetry .....	429
Snorkel .....	429
<b>INDEX.....</b>	<b>431</b>

# Introduction

---

The term *machine learning* has all sorts of meanings attached to it today, especially after Hollywood (and other movie studios) have gotten into the picture. Films such as *Ex Machina* have tantalized the imaginations of moviegoers the world over and made machine learning into all sorts of things that it really isn't. Of course, most of us have to live in the real world, where machine learning actually does perform an incredible array of tasks that have nothing to do with androids that can pass the Turing Test (fooling their makers into believing they're human). *Machine Learning For Dummies*, 2<sup>nd</sup> Edition gives you a view of machine learning in the real world and exposes you to the amazing feats you really can perform using this technology.

Even though the tasks that you perform using machine learning may seem a bit mundane when compared to the movie version, by the time you finish this book, you realize that these mundane tasks have the power to impact the lives of everyone on the planet in nearly every aspect of their daily lives. In short, machine learning is an incredible technology — just not in the way that some people have imagined.

This second edition of the book contains a significant number of changes, not the least of which is that it's using pure Python code for the examples now upon request from our readers. You can still download R versions of every example, which is actually better than before when only some of the examples were available in R. In addition, the book contains new topics, including an entire chapter that discusses machine learning ethics.

## About This Book

---

Machines and humans learn in entirely different ways, which is why the first part of this book is essential to your understanding of machine learning. Machines perform routine tasks at incredible speeds, but still require humans to do the actual thinking.

The second part of this book is about getting your system set up to use the various Python coding examples. The two setups work for desktop systems using Windows, Mac OS, or Linux, or mobile devices that have access to a Google Colab compatible browser.



TIP

If you're using R, you'll find a README file in the R download file that contains instructions for configuring your R Anaconda environment.

The third part of the book discusses math basics with regard to machine learning requirements. It prepares you to perform math tasks associated with algorithms used in machine learning to make either predictions or classifications from your data.

The fourth part of the book helps you discover what to do about data that isn't quite up to par. This part is also where you start learning about similarity and working with linear models. The most advanced chapter tells you how to work with ensembles of learners to perform tasks that might not otherwise be reasonable to complete.

The fifth part of the book is about practical application of machine learning techniques. You see how to do things like classify images, work with opinions and sentiments, and recommend products and movies.

The last part of the book contains helpful information to enhance your machine learning experience. This part of the book also contains a chapter specifically oriented toward ethical data use.

To make absorbing the concepts easy, this book uses the following conventions:

- » Text that you're meant to type just as it appears in the book is in **bold**. The exception is when you're working through a step list: Because each step is bold, the text to type is not bold.
- » Web addresses and programming code appear in monofont. If you're reading a digital version of this book on a device connected to the Internet, you can click or tap the web address to visit that website, like this: `https://www.dummies.com`.
- » When you need to type command sequences, you see them separated by a special arrow, like this: File ⇨ New File. In this example, you go to the File menu first and then select the New File entry on that menu.
- » When you see words in *italics* as part of a typing sequence, you need to replace that value with something that works for you. For example, if you see "Type **Your Name** and press Enter," you need to replace *Your Name* with your actual name.

## Foolish Assumptions

This book is designed for novice and professional alike. You can either read this book from cover to cover or look up topics and treat the book as a reference guide. However, we've made some assumptions about your level of knowledge when we put the



book together. You should already know how to use your device and work with the operating system that supports it. You also know how to perform tasks like downloading files and installing applications. You can interact with Internet well enough to locate the resources you need to work with the book. You know how to work with archives, such as the .zip file format. Finally, a basic knowledge of math is helpful.

## Icons Used in This Book

As you read this book, you see icons in the margins that indicate material of interest. This section briefly describes each icon.



TIP

The tips in this book are time-saving techniques or pointers to resources that you should try so that you can get the maximum benefit from machine learning.



WARNING

You should avoid doing anything that's marked with a Warning icon. Otherwise, you might find that your application fails to work as expected, you get incorrect answers from seemingly bulletproof code, or (in the worst-case scenario) you lose data.



TECHNICAL  
STUFF

Whenever you see this icon, think advanced tip or technique. Skip these bits of information whenever you like.



REMEMBER

This text usually contains an essential process or a bit of information that you must know to perform machine learning tasks successfully.

## Beyond the Book

If you want to email us, please do! Make sure you send your book-specific requests to: [John@JohnMuellerBooks.com](mailto:John@JohnMuellerBooks.com). We want to ensure that your book experience is the best one possible. The blog entries at <http://blog.johnmuellerbooks.com/> contain a wealth of additional information about this book. You can check out John's website at <http://www.johnmuellerbooks.com/>. You can also access other cool materials:

» **Cheat Sheet:** A cheat sheet provides you with some special notes on things you can do with machine learning that not every other scientist knows. You can find the Cheat Sheet for this book at [www.dummies.com](http://www.dummies.com). Type **Machine Learning For Dummies** in the Search box and click the Cheat Sheets option that appears.

» **Errata:** You can find errata by entering this book's title in the Search box at [www.dummies.com](http://www.dummies.com), which takes you to this book's page. In addition to errata, check out the blog posts with answers to reader questions and demonstrations of useful book-related techniques at <http://blog.johnmuelเลอร์books.com/>.

» **Companion files:** The source code is available for download. All the book examples tell you precisely which example project to use. You can find these files at this book's page at [www.dummies.com/go/machinelearningfd2e](http://www.dummies.com/go/machinelearningfd2e).

We've also had trouble with the datasets used in the previous edition of this book. Sometimes the datasets change or might become unavailable. Given that you likely don't want to download a large dataset unless you're interested in that example, we've made the non-toy datasets (those available with a package) available at <https://github.com/lmassaron/datasets>. You don't actually need to download them, though; the example code will perform that task for you automatically when you run it.

## Where to Go from Here

Most people will want to start this book from the beginning, because it contains a good deal of information about how the real world view of machine learning differs from what movies might tell you. However, if you already have a first grounding in the reality of machine learning, you can always skip to the next part of the book.

Chapter 4 is where you want to go if you want to use a desktop setup, while Chapter 6 is helpful when you want to use a mobile device. Your preexisting setup may not work with the book's examples because you might have different versions of the various products. It's essential that you use the correct product versions to ensure success. Even if you choose to go with your own setup, consider reviewing Chapter 5 unless you're an expert Python coder already.

If you're already an expert with Python and know how machine learning works, you could always skip to Chapter 7. Starting at Chapter 7 will help you get into the examples quickly so that you spend less time with basics and more time with intermediate machine learning tasks. You can always go back and review the previous materials as needed.

# 1

## **Introducing How Machines Learn**

#### IN THIS PART . . .

Discovering how AI really works and what it can do for you

Considering what the term *big data* means

Understanding the role of statistics in machine learning

Defining where machine learning will take society in the future

- » Seeing the dream; getting beyond the hype of artificial intelligence (AI)
- » Comparing AI to machine learning
- » Understanding the engineering portion of AI and machine learning
- » Delineating where engineering ends and art begins

## Chapter 1

# Getting the Real Story about AI

**A**rtificial Intelligence (AI), the appearance of intelligence in machines, is a huge topic today, and it's getting bigger all the time thanks to the success of new technologies (see some current examples at <https://thinkml.ai/top-5-ai-achievements-of-2019/>). However, most people are looking for everyday applications, such as talking to their smartphone. Talking to your smartphone is both fun and helpful to find out things like the location of the best sushi restaurant in town or to discover how to get to the concert hall. As you talk to your smartphone, it learns more about the way you talk and makes fewer mistakes in understanding your requests. The capability of your smartphone to learn and interpret your particular way of speaking is an example of an AI, and part of the technology used to make it happen is *machine learning*, the use of various techniques to allow algorithms to work better based on experience.

You likely make limited use of machine learning and AI all over the place today without really thinking about it. For example, the capability to speak to devices and have them actually do what you intend is an example of machine learning at work. Likewise, recommender systems, such as those found on Amazon, help you make purchases based on criteria such as previous product purchases or products that complement a current choice. The use of both AI and machine learning will only increase with time.

In this chapter, you delve into AI and discover what it means from several perspectives, including how it affects you as a consumer and as a scientist or engineer. You also discover that AI doesn't equal machine learning, even though the media often confuse the two. Machine learning is definitely different from AI, even though the two are related.

## Moving beyond the Hype

As any technology becomes bigger, so does the hype, and AI certainly has a lot of hype surrounding it. For one thing, some people have decided to engage in fear mongering rather than science. Killer robots, such as those found in the film *The Terminator*, really aren't going to be the next big thing. Your first real experience with an android AI is more likely to be in the form a health care assistant (<https://www.robotics.org/blog-article.cfm/The-Future-of-Elder-Care-is-Service-Robots/262>) or possibly as a coworker (<https://www.computerworld.com/article/2990849/meet-the-virtual-woman-who-may-take-your-job.html>). The reality is that you interact with AI and machine learning in far more mundane ways already. Part of the reason you need to read this chapter is to get past the hype and discover what AI can do for you today.



REMEMBER

You may also have heard machine learning and AI used interchangeably. AI includes machine learning, but machine learning doesn't fully define AI. This chapter helps you understand the relationship between machine learning and AI so that you can better understand how this book helps you move into a technology that used to appear only within the confines of science fiction novels.

Machine learning and AI both have strong engineering components. That is, you can quantify both technologies precisely based on *theory* (substantiated and tested explanations) rather than simply *hypothesis* (a suggested explanation for a phenomenon). In addition, both have strong science components, through which people test concepts and create new ideas of how expressing the thought process might be possible. Finally, machine learning also has an artistic component, and this is where a talented scientist can excel. In some cases, AI and machine learning both seemingly defy logic, and only the true artist can make them work as expected.

## YES, FULLY AUTONOMOUS WEAPONS EXIST

Before people send us their latest dissertations about fully autonomous weapons, yes, some benighted souls are working on such technologies. You'll find some discussions of the ethics of AI in this book, but for the most part, the book focuses on positive, helpful uses of AI to aid humans, rather than kill them, because most AI research reflects these uses. You can find articles on the pros and cons of AI online, such as the Towards Data Science article at <https://towardsdatascience.com/advantages-and-disadvantages-of-artificial-intelligence-182a5ef6588c> and the *Emerj* article at <https://emerj.com/ai-sector-overviews/autonomous-weapons-in-the-military/>.

If you really must scare yourself, you can find all sorts of sites, such as <https://www.reachingcriticalwill.org/resources/fact-sheets/critical-issues/7972-fully-autonomous-weapons>, that discuss the issue of fully autonomous weapons in some depth. Sites such as Campaign to Stop Killer Robots (<https://www.stopkillerrobots.org/>) can also fill in some details for you. We do encourage you to sign the letter banning autonomous weapons at <https://futureoflife.org/open-letter-autonomous-weapons/> — there truly is no need for them.

However, it's important to remember that bans against space-based, chemical, and certain laser weapons all exist. Countries recognize that these weapons don't solve anything. Countries will also likely ban fully autonomous weapons simply because the citizenry won't stand for killer robots. The bottom line is that the focus of this book is on helping you understand machine learning in a positive light.

## Dreaming of Electric Sheep

*Androids* (a specialized kind of robot that looks and acts like a human, such as Data in *Star Trek: The Next Generation*) and some types of *humanoid robots* (a kind of robot that has human characteristics but is easily distinguished from a human, such as C-3PO in *Star Wars*) have become the poster children for AI (see the dancing robots at <https://www.youtube.com/watch?v=1TckiTBaWkw>). They present computers in a form that people can *anthropomorphize* (give human characteristics to, even though they aren't human). In fact, it's entirely possible that one day you won't be able to distinguish between human and artificial life with ease. Science fiction authors, such as Philip K. Dick, have long predicted such an occurrence, and it seems all too possible today. The story "Do Androids Dream of Electric Sheep?" discusses the whole concept of more real than real. The idea appears as part of the plot in the movie *Blade Runner* (<https://www.warnerbros.com/movies/blade-runner>). However, some uses of robots today are just plain fun, as in the

Robot Restaurant show at <https://www.youtube.com/watch?v=l1vvTtz8hpg>. The sections that follow help you understand how close technology currently gets to the ideals presented by science fiction authors and the movies.



The current state of the art is lifelike, but you can easily tell that you're talking to an android. Viewing videos online can help you understand that androids that are indistinguishable from humans are nowhere near any sort of reality today. Check out the Japanese robots at <https://www.youtube.com/watch?v=LyyytwT-BMk> and <https://www.cnbc.com/2019/10/31/human-like-androids-have-entered-the-workplace-and-may-take-your-job.html>. One of the more lifelike examples is Erica (<https://www.youtube.com/watch?v=oR1wvLubFvg>), who is set to appear in a science fiction film. Her story appears on *HuffPost* at [https://www.huffpost.com/entry/erica-japanese-robot-science-fiction-film\\_n\\_5ef6523dc5b6acab284181c3](https://www.huffpost.com/entry/erica-japanese-robot-science-fiction-film_n_5ef6523dc5b6acab284181c3). The point is, technology is just starting to get to the point where people may eventually be able to create lifelike robots and androids, but they don't exist today.

## Understanding the history of AI and machine learning

There is a reason, other than anthropomorphization, that humans see the ultimate AI as one that is contained within some type of android. Ever since the ancient Greeks, humans have discussed the possibility of placing a mind inside a mechanical body. One such myth is that of a mechanical man called Talos (<http://www.ancient-wisdom.com/greekautomata.htm>). The fact that the ancient Greeks had complex mechanical devices, only one of which still exists (read about the Antikythera mechanism at <http://www.ancient-wisdom.com/antikythera.htm>), makes it quite likely that their dreams were built on more than just fantasy. Throughout the centuries, people have discussed mechanical persons capable of thought (such as Rabbi Judah Loew's Golem, <https://www.nytimes.com/2009/05/11/world/europe/11golem.html>).

AI is built on the hypothesis that mechanizing thought is possible. During the first millennium, Greek, Indian, and Chinese philosophers all worked on ways to perform this task. As early as the seventeenth century, Gottfried Leibniz, Thomas Hobbes, and René Descartes discussed the potential for rationalizing all thought as simply math symbols. Of course, the complexity of the problem eluded them (and still eludes us today, despite the advances you read about in Part 3 of this book). The point is that the vision for AI has been around for an incredibly long time, but the implementation of AI is relatively new.

The true birth of AI as we know it today began with Alan Turing's publication of "Computing Machinery and Intelligence" in 1950 (<https://www.csee.umbc.edu/~turing/>).



[edu/courses/471/papers/turing.pdf](http://edu/courses/471/papers/turing.pdf)). In this paper, Turing explored the idea of how to determine whether machines can think. Of course, this paper led to the Imitation Game involving three players. Player A is a computer and Player B is a human. Each must convince Player C (a human who can't see either Player A or Player B) that they are human. If Player C can't determine who is human and who isn't on a consistent basis, the computer wins.

A continuing problem with AI is too much optimism. The problem that scientists are trying to solve with AI is incredibly complex. However, the early optimism of the 1950s and 1960s led scientists to believe that the world would produce intelligent machines in as little as 20 years. After all, machines were doing all sorts of amazing things, such as playing complex games. AI currently has its greatest success in areas such as logistics, data mining, and medical diagnosis.

## Exploring what machine learning can do for AI

Machine learning relies on algorithms to analyze huge datasets. Currently, machine learning can't provide the sort of AI that the movies present. Even the best algorithms can't think, feel, present any form of self-awareness, or exercise free will. What machine learning can do is perform predictive analytics far faster than any human can. As a result, machine learning can help humans work more efficiently. The current state of AI, then, is one of performing analysis, but humans must still consider the implications of that analysis — making the required moral and ethical decisions. The “Considering the Relationship between AI and Machine Learning” section of this chapter delves more deeply into precisely how machine learning contributes to AI as a whole. The essence of the matter is that machine learning provides just the learning part of AI, and that part is nowhere near ready to create an AI of the sort you see in films.



REMEMBER

The main point of confusion between learning and intelligence is that people assume that simply because a machine gets better at its job (learning) it's also aware (intelligence). Nothing supports this view of machine learning. The same phenomenon occurs when people assume that a computer is purposely causing problems for them. The computer can't assign emotions and therefore acts only upon the input provided and the instruction contained within an application to process that input. A true AI will eventually occur when computers can finally emulate the clever combination used by nature:

- » **Genetics:** Slow learning from one generation to the next
- » **Teaching:** Fast learning from organized sources
- » **Exploration:** Spontaneous learning through media and interactions with others

# Considering the goals of machine learning

At present, AI is based on machine learning, and machine learning is essentially different from statistics. Yes, machine learning has a statistical basis, but it makes some different assumptions than statistics do because the goals are different. Table 1-1 lists some features to consider when comparing AI and machine learning to statistics.

**TABLE 1-1:** Comparing Machine Learning to Statistics

Technique	Machine Learning	Statistics
Data handling	Works with big data in the form of networks and graphs; raw data from sensors or the web text is split into training and test data.	Models are used to create predictive power on small samples.
Data input	The data is sampled, randomized, and transformed to maximize accuracy scoring in the prediction of out-of-sample (or completely new) examples.	Parameters interpret real-world phenomena and provide a stress on magnitude.
Result	Probability is taken into account for comparing what could be the best guess or decision.	The output captures the variability and uncertainty of parameters.
Assump-tions	The scientist learns from the data.	The scientist assumes a certain output and tries to prove it.
Distribution	The distribution is unknown or ignored before learning from data.	The scientist assumes a well-defined distribution.
Fitting	The scientist creates a best fit, but generalizable, model.	The result is fit to the present data distribution.

## Defining machine learning limits based on hardware

Huge datasets require huge amounts of memory. Unfortunately, the requirements don't end there. When you have huge amounts of data and memory, you must also have processors with multiple cores and high speeds. One of the problems that scientists are striving to solve is how to use existing hardware more efficiently. In some cases, waiting for days to obtain a result to a machine learning problem simply isn't possible. The scientists who want to know the answer need it quickly, even if the result isn't quite right. With this in mind, investments in better hardware also require investments in better science. This book considers some of the following issues as part of making your machine learning experience better:

- » **Obtaining a useful result:** As you work through the book, you discover that you need to obtain a useful result first, before you can refine it. In addition, sometimes tuning an algorithm goes too far and the result becomes quite fragile (and possibly useless outside a specific dataset).
- » **Asking the right question:** Many people get frustrated in trying to obtain an answer from machine learning because they keep tuning their algorithm without asking a different question. To use hardware efficiently, sometimes you must step back and review the question you're asking. The question might be wrong, which means that even the best hardware will never find the answer.
- » **Relying on intuition too heavily:** All machine learning questions begin as a hypothesis. A scientist uses intuition to create a starting point for discovering the answer to a question. Failure is more common than success when working through a machine learning experience. Your intuition adds the art to the machine learning experience, but sometimes intuition is wrong and you have to revisit your assumptions.



TECHNICAL  
STUFF

When you begin to realize the importance of environment to machine learning, you can also begin to understand the need for the right hardware and in the right balance to obtain a desired result. The current state-of-the-art systems actually rely on Graphical Processing Units (GPUs) to perform machine learning tasks. Relying on GPUs does speed the machine learning process considerably. A full discussion of using GPUs is outside the scope of this book, but you can read more about the topic at <https://devblogs.nvidia.com/parallelforall/bidmach-machine-learning-limit-gpus/> and <https://towardsdatascience.com/what-is-a-gpu-and-do-you-need-one-in-deep-learning-718b9597aa0d>.

## Overcoming AI Fantasies

As with many other technologies, AI and machine learning both have their fantasy or fad uses. For example, some people are using machine learning to create Picasso-like art from photos using products like NightCafé (<https://creator.nightcafe.studio/>), which supports people who really enjoy this art form. You can read all about using machine learning to create art at <https://www.washingtonpost.com/news/innovations/wp/2015/08/31/this-algorithm-can-create-a-new-van-gogh-or-picasso-in-just-an-hour/>. Of course, the problems with such use are many. For one thing, most people wouldn't really want a Picasso created in this manner except as a fad item (because no one had done it before). The point of art isn't in creating an interesting interpretation of a particular real-world representation, but rather in seeing how the artist interpreted it. The end of the article points out that the computer can only copy an

existing style at this stage — not create an entirely new style of its own. The following sections discuss AI and machine learning fantasies of various sorts.

## Discovering the fad uses of AI and machine learning

AI is entering an era of innovation that you used to read about only in science fiction. It can be hard to determine whether a particular AI use is real or simply the dream child of a determined scientist. For example, *The Six Million Dollar Man* ([https://en.wikipedia.org/wiki/The\\_Six\\_Million\\_Dollar\\_Man](https://en.wikipedia.org/wiki/The_Six_Million_Dollar_Man)) is a television series that looked fanciful at one time. When it was introduced, no one actually thought that we'd have real-world bionics at some point. However, Hugh Herr (<https://www.smithsonianmag.com/innovation/future-robotic-legs-180953040/>) and others (<https://www.fiercebiotech.com/medtech/using-onboard-ai-to-power-quicker-more-complex-prosthetic-hands>) have other ideas — bionic legs and arms really are possible now. Of course, they aren't available for everyone yet; the technology is only now becoming useful. Muddying the waters is *The Six Billion Dollar Man* movie, based partly on *The Six Million Dollar Man* television series (<https://www.cinemablend.com/new/Mark-Wahlberg-Six-Billion-Dollar-Man-Just-Made-Big-Change-91947.html>), which has suffered delays for various reasons (<https://screenrant.com/mark-wahlberg-six-billion-dollar-man-delays-updates/>). The fact is that AI and machine learning will both present opportunities to create some amazing technologies and that we're already at the stage of creating those technologies, but you still need to take what you hear with a huge grain of salt.

One of the more interesting uses of machine learning for entertainment purposes is the movie *B* (<https://www.cinemablend.com/news/2548939/one-sci-fi-movie-will-be-able-to-film-during-the-pandemic-thanks-to-casting-an-ai-robot-as-its-lead>), which stars an android named Erica. The inventors of Erica, Hiroshi Ishiguro and Kohei Ogawa, have spent a great deal of time trying to make her lifelike by trying to implement the human qualities of intent and desire (<https://www.yoichimatsuyama.com/conversation-with-evolving-robotic-species-interview-with-hiroshi-ishiguro/>). The result is something that encroaches on the uncanny valley (<https://www.scientificamerican.com/article/why-uncanny-valley-human-look-alikes-put-us-on-edge/>) in a new way. The plot of this movie will be on the same order as *Ex Machina* (<https://www.indiewire.com/2020/06/ex-machina-real-robot-erica-science-fiction-movie-1234569484/>).



REMEMBER

To make the future uses of AI and machine learning match the concepts that science fiction has presented over the years, real-world programmers, data scientists, and other stakeholders need to create tools. Nothing happens by magic, even though it may look like magic when you don't know what's happening behind the scenes. In order for the fad uses for AI and machine learning to become real-world uses, developers, data scientists, and others need to continue building real-world tools that may be hard to imagine at this point.

## Considering the true uses of AI and machine learning

You find AI and machine learning used in a great many applications today. The only problem is that the technology works so well that you don't know that it even exists. In fact, you might be surprised to find that many devices in your home already make use of both technologies. Both technologies definitely appear in your car and most especially in the workplace. In fact, the uses for both AI and machine learning number in the millions — all safely out of sight even when they're quite dramatic in nature. Here are just a few of the ways in which you might see AI used:

- » **Fraud detection:** You get a call from your credit card company asking whether you made a particular purchase. The credit card company isn't being nosy; it's simply alerting you to the fact that someone else could be making a purchase using your card. The AI embedded within the credit card company's code detected an unfamiliar spending pattern and alerted someone to it.
- » **Resource scheduling:** Many organizations need to schedule the use of resources efficiently. For example, a hospital may have to determine where to put a patient based on the patient's needs, availability of skilled experts, and the amount of time the doctor expects the patient to be in the hospital.
- » **Complex analysis:** Humans often need help with complex analysis because there are literally too many factors to consider. For example, the same set of symptoms could indicate more than one problem. A doctor or other expert might need help making a diagnosis in a timely manner to save a patient's life.
- » **Automation:** Any form of automation can benefit from the addition of AI to handle unexpected changes or events. A problem with some types of automation today is that an unexpected event, such as an object in the wrong place, can actually cause the automation to stop. Adding AI to the automation can allow the automation to handle unexpected events and continue as if nothing happened.
- » **Customer service:** The customer service line you call today may not even have a human behind it. The automation is good enough to follow scripts and use various resources to handle the vast majority of your questions. With

good voice inflection (provided by AI as well), you may not even be able to tell that you're talking with a computer.

- » **Safety systems:** Many of the safety systems found in machines of various sorts today rely on AI to take over the vehicle in a time of crisis. For example, many automatic braking systems rely on AI to stop the car based on all the inputs that a vehicle can provide, such as the direction of a skid.
- » **Machine efficiency:** AI can help control a machine in such a manner as to obtain maximum efficiency. The AI controls the use of resources so that the system doesn't overshoot speed or other goals. Every ounce of power is used precisely as needed to provide the desired services.

This list doesn't even begin to scratch the surface. You can find AI used in many other ways. However, it's also useful to view uses of machine learning outside the normal realm that many consider the domain of AI. Here are a few uses for machine learning that you might not associate with an AI:

- » **Access control:** In many cases, access control is a yes or no proposition. An employee smartcard grants access to a resource much in the same way that people have used keys for centuries. Some locks do offer the capability to set times and dates that access is allowed, but the coarse-grained control doesn't really answer every need. By using machine learning, you can determine whether an employee should gain access to a resource based on role and need. For example, an employee can gain access to a training room when the training reflects an employee role.
- » **Animal protection:** The ocean might seem large enough to allow animals and ships to cohabitate without problem. Unfortunately, many animals get hit by ships each year. A machine learning algorithm could allow ships to avoid animals by learning the sounds and characteristics of both the animal and the ship.
- » **Predicting wait times:** Most people don't like waiting when they have no idea of how long the wait will be. Machine learning allows an application to determine waiting times based on staffing levels, staffing load, complexity of the problems the staff is trying to solve, availability of resources, and so on.

## Being useful; being mundane

Even though the movies make it sound like AI is going to make a huge splash, and you do sometimes see some incredible uses for AI in real life, the fact of the matter is that most uses for AI are mundane, even boring. For example, a recent article