

Introducing Machine Learning



Professional



Dino Esposito
Francesco Esposito

FREE SAMPLE CHAPTER

SHARE WITH OTHERS



Introducing Machine Learning

Dino Esposito
Francesco Esposito

Introducing Machine Learning

**Published with the authorization of Microsoft Corporation by:
Pearson Education, Inc.**

Copyright © 2020 by Dino Esposito and Francesco Esposito

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms, and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit www.pearson.com/permissions.

No patent liability is assumed with respect to the use of the information contained herein. Although every precaution has been taken in the preparation of this book, the publisher and author assume no responsibility for errors or omissions. Nor is any liability assumed for damages resulting from the use of the information contained herein.

ISBN-13: 978-0-13-556566-7

ISBN-10: 0-13-556566-9

Library of Congress Control Number: 2019954810

ScoutAutomatedPrintCode

Trademarks

Microsoft and the trademarks listed at <http://www.microsoft.com> on the "Trademarks" webpage are trademarks of the Microsoft group of companies. All other marks are property of their respective owners.

Warning and Disclaimer

Every effort has been made to make this book as complete and as accurate as possible, but no warranty or fitness is implied. The information provided is on an "as is" basis. The author, the publisher, and Microsoft Corporation shall have neither liability nor responsibility to any person or entity with respect to any loss or damages arising from the information contained in this book or from the use of the programs accompanying it.

Special Sales

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact intlcs@pearson.com.

Editor-in-Chief

Brett Bartow

Executive Editor

Loretta Yates

Development Editor

Mark Renfrow

Assistant Sponsoring Editor

Charvi Arora

Managing Editor

Sandra Schroeder

Senior Project Editor

Tonya Simpson

Copy Editor

Chuck Hutchinson

Indexer

Cheryl Ann Lenser

Proofreader

Abigail Manheim

Technical Editor

Cesar De la Torre Llorente

Editorial Assistant

Cindy Teeters

Cover Designer

Twist Creative, Seattle

Cover Image

IMOGI graphics

Compositor

codeMantra

Dedications

To Michela and her dreams

To my loved ones, to whom I couldn't help but dedicate a book

This page intentionally left blank

Contents at a Glance

Introduction

xxiii

| | | |
|-----------------|--|-----|
| PART I | LAYING THE GROUNDWORK OF MACHINE LEARNING | |
| CHAPTER 1 | How Humans Learn | 3 |
| CHAPTER 2 | Intelligent Software | 23 |
| CHAPTER 3 | Mapping Problems and Algorithms | 33 |
| CHAPTER 4 | General Steps for a Machine Learning Solution | 49 |
| CHAPTER 5 | The Data Factor | 67 |
| PART II | MACHINE LEARNING IN .NET | |
| CHAPTER 6 | The .NET Way | 77 |
| CHAPTER 7 | Implementing the ML.NET Pipeline | 93 |
| CHAPTER 8 | ML.NET Tasks and Algorithms | 105 |
| PART III | FUNDAMENTALS OF SHALLOW LEARNING | |
| CHAPTER 9 | Math Foundations of Machine Learning | 135 |
| CHAPTER 10 | Metrics of Machine Learning | 151 |
| CHAPTER 11 | How to Make Simple Predictions: Linear Regression | 165 |
| CHAPTER 12 | How to Make Complex Predictions and Decisions: Trees | 181 |
| CHAPTER 13 | How to Make Better Decisions: Ensemble Methods | 197 |
| CHAPTER 14 | Probabilistic Methods: Naïve Bayes | 211 |
| CHAPTER 15 | How to Group Data: Classification and Clustering | 229 |
| PART IV | FUNDAMENTALS OF DEEP LEARNING | |
| CHAPTER 16 | Feed-Forward Neural Networks | 255 |
| CHAPTER 17 | Design of a Neural Network | 273 |
| CHAPTER 18 | Other Types of Neural Networks | 291 |
| CHAPTER 19 | Sentiment Analysis: An End-to-End Solution | 309 |

PART V **FINAL THOUGHTS**

| | | |
|------------|--------------------------------------|-----|
| CHAPTER 20 | AI Cloud Services for the Real World | 327 |
| CHAPTER 21 | The Business Perception of AI | 339 |
| | <i>Index</i> | 351 |

| | |
|--|-----------|
| General Artificial Intelligence..... | 27 |
| Unsupervised Learning | 27 |
| Supervised Learning | 29 |
| Summary | 32 |
| Chapter 3 Mapping Problems and Algorithms | 33 |
| Fundamental Problems | 33 |
| Classifying Objects | 34 |
| Predicting Results | 36 |
| Grouping Objects | 38 |
| More Complex Problems..... | 40 |
| Image Classification | 41 |
| Object Detection | 41 |
| Text Analytics | 42 |
| Automated Machine Learning..... | 42 |
| Aspects of an AutoML Platform | 42 |
| The AutoML Model Builder in Action | 45 |
| Summary | 48 |
| Chapter 4 General Steps for a Machine Learning Solution | 49 |
| Data Collection..... | 50 |
| Data-Driven Culture in the Organization..... | 50 |
| Storage Options | 51 |
| Data Preparation | 52 |
| Improving Data Quality..... | 53 |
| Cleaning Data | 53 |
| Feature Engineering | 54 |
| Finalizing the Training Dataset | 56 |
| Model Selection and Training | 58 |
| The Algorithm Cheat Sheet..... | 59 |
| The Case for Neural Networks..... | 61 |
| Evaluation of the Model Performance | 62 |

| | |
|---|----|
| Deployment of the Model | 64 |
| Choosing the Appropriate Hosting Platform | 64 |
| Exposing an API | 65 |
| Summary | 66 |

Chapter 5 The Data Factor 67

| | |
|---|----|
| Data Quality | 67 |
| Data Validity | 68 |
| Data Collection | 69 |
| Data Integrity | 70 |
| Completeness | 70 |
| Uniqueness | 70 |
| Timeliness | 71 |
| Accuracy | 71 |
| Consistency | 71 |
| What's a Data Scientist, Anyway? | 71 |
| The Data Scientist at Work | 72 |
| The Data Scientist Tool Chest | 73 |
| Data Scientists and Software Developers | 73 |
| Summary | 74 |

PART II MACHINE LEARNING IN .NET

Chapter 6 The .NET Way 77

| | |
|---|----|
| Why (Not) Python? | 78 |
| Why Is Python So Popular in Machine Learning? | 78 |
| Taxonomy of Python Machine Learning Libraries | 80 |
| End-to-End Solutions on Top of Python Models | 82 |
| Introducing ML.NET | 83 |
| Creating and Consuming Models in ML.NET | 84 |
| Elements of the Learning Context | 87 |
| Summary | 91 |

| | | |
|------------------|---|------------|
| Chapter 7 | Implementing the ML.NET Pipeline | 93 |
| | The Data to Start From | 93 |
| | Exploring the Dataset | 94 |
| | Applying Common Data Transformations | 94 |
| | Considerations on the Dataset | 95 |
| | The Training Step | 96 |
| | Picking an Algorithm | 96 |
| | Measuring the Actual Value of an Algorithm | 97 |
| | Planning the Testing Phase | 97 |
| | A Look at the Metrics | 98 |
| | Price Prediction from Within a Client Application | 99 |
| | Getting the Model File | 99 |
| | Setting Up the ASP.NET Application | 99 |
| | Making a Taxi Fare Prediction | 100 |
| | Devising an Adequate User Interface | 102 |
| | Questioning Data and Approach to the Problem | 103 |
| | Summary | 103 |
| | | |
| Chapter 8 | ML.NET Tasks and Algorithms | 105 |
| | The Overall ML.NET Architecture | 105 |
| | Involved Types and Interfaces | 105 |
| | Data Representation | 107 |
| | Supported Catalogs | 109 |
| | Classification Tasks | 111 |
| | Binary Classification | 111 |
| | Multiclass Classification | 116 |
| | Clustering Tasks | 122 |
| | Preparing Data for Work | 122 |
| | Training the Model | 123 |
| | Evaluating the Model | 124 |
| | Transfer Learning | 126 |
| | Steps for Building an Image Classifier | 127 |

| | |
|---|-----|
| Applying Necessary Data Transformations | 127 |
| Composing and Training the Model | 129 |
| Margin Notes on Transfer Learning | 131 |
| Summary | 132 |

PART III FUNDAMENTALS OF SHALLOW LEARNING

Chapter 9 Math Foundations of Machine Learning 135

| | |
|---|-----|
| Under the Umbrella of Statistics | 135 |
| The Mean in Statistics | 136 |
| The Mode in Statistics | 138 |
| The Median in Statistics | 139 |
| Bias and Variance | 141 |
| The Variance in Statistics | 142 |
| The Bias in Statistics | 144 |
| Data Representation | 145 |
| Five-number Summary | 145 |
| Histograms | 146 |
| Scatter Plots | 148 |
| Scatter Plot Matrices | 148 |
| Plotting at the Appropriate Scale | 149 |
| Summary | 150 |

Chapter 10 Metrics of Machine Learning 151

| | |
|--|-----|
| Statistics vs. Machine Learning | 151 |
| The Ultimate Goal of Machine Learning | 152 |
| From Statistical Models to Machine Learning Models | 153 |
| Evaluation of a Machine Learning Model | 155 |
| From Dataset to Predictions | 155 |
| Measuring the Precision of a Model | 157 |
| Preparing Data for Processing | 162 |
| Scaling | 162 |

| | |
|--|------------|
| Standardization | 163 |
| Normalization | 163 |
| Summary | 163 |
| Chapter 11 How to Make Simple Predictions: Linear Regression | 165 |
| The Problem | 165 |
| Guessing Results Guided by Data | 166 |
| Making Hypotheses About the Relationship | 167 |
| The Linear Algorithm | 169 |
| The General Idea | 169 |
| Identifying the Cost Function | 170 |
| The Ordinary Least Square Algorithm | 171 |
| The Gradient Descent Algorithm | 174 |
| How Good Is the Algorithm? | 178 |
| Improving the Solution | 178 |
| The Polynomial Route | 178 |
| Regularization | 179 |
| Summary | 180 |
| Chapter 12 How to Make Complex Predictions and Decisions: Trees | 181 |
| The Problem | 181 |
| What's a Tree, Anyway? | 182 |
| Trees in Machine Learning | 183 |
| A Sample Tree-Based Algorithm | 183 |
| Design Principles for Tree-Based Algorithms | 185 |
| Decision Trees versus Expert Systems | 185 |
| Flavors of Tree Algorithms | 186 |
| Classification Trees | 187 |
| How the CART Algorithm Works | 187 |
| How the ID3 Algorithm Works | 191 |

| | |
|--|------------|
| Regression Trees | 194 |
| How the Algorithm Works | 194 |
| Tree Pruning | 195 |
| Summary | 196 |
| Chapter 13 How to Make Better Decisions: Ensemble Methods | 197 |
| The Problem | 197 |
| The Bagging Technique | 198 |
| Random Forest Algorithms | 198 |
| Steps of the Algorithms | 200 |
| Pros and Cons | 202 |
| The Boosting Technique | 203 |
| The Power of Boosting | 203 |
| Gradient Boosting | 206 |
| Pros and Cons | 210 |
| Summary | 210 |
| Chapter 14 Probabilistic Methods: Naïve Bayes | 211 |
| Quick Introduction to Bayesian Statistics | 211 |
| Introducing Bayesian Probability | 212 |
| Some Preliminary Notation | 212 |
| Bayes' Theorem | 214 |
| A Practical Code Review Example | 215 |
| Applying Bayesian Statistics to Classification | 216 |
| Initial Formulation of the Problem | 217 |
| A Simplified (Yet Effective) Formulation | 217 |
| Practical Aspects of Bayesian Classifiers | 218 |
| Naïve Bayes Classifiers | 219 |
| The General Algorithm | 219 |
| Multinomial Naïve Bayes | 220 |
| Bernoulli Naïve Bayes | 223 |
| Gaussian Naïve Bayes | 224 |

| | |
|--|-----|
| Naïve Bayes Regression | 226 |
| Foundation of Bayesian Linear Regression | 226 |
| Applications of Bayesian Linear Regression | 228 |
| Summary | 228 |

Chapter 15 How to Group Data: Classification and Clustering 229

| | |
|---|-----|
| A Basic Approach to Supervised Classification | 230 |
| The K-Nearest Neighbors Algorithm | 230 |
| Steps of the Algorithm | 232 |
| Business Scenarios | 234 |
| Support Vector Machine | 235 |
| Overview of the Algorithm | 235 |
| A Quick Mathematical Refresher | 239 |
| Steps of the Algorithm | 240 |
| Unsupervised Clustering | 245 |
| A Business Case: Reducing the Dataset | 245 |
| The K-Means Algorithm | 246 |
| The K-Modes Algorithm | 247 |
| The DBSCAN Algorithm | 248 |
| Summary | 251 |

PART IV FUNDAMENTALS OF DEEP LEARNING

Chapter 16 Feed-Forward Neural Networks 255

| | |
|--|-----|
| A Brief History of Neural Networks | 255 |
| The McCulloch-Pitt Neuron | 255 |
| Feed-Forward Networks | 256 |
| More Sophisticated Networks | 256 |
| Types of Artificial Neurons | 257 |
| The Perceptron Neuron | 257 |
| The Logistic Neuron | 260 |

| | |
|---|------------|
| Training a Neural Network | 263 |
| The Overall Learning Strategy | 263 |
| The Backpropagation Algorithm | 264 |
| Summary | 270 |
| Chapter 17 Design of a Neural Network | 273 |
| Aspects of a Neural Network | 273 |
| Activation Functions | 274 |
| Hidden Layers | 277 |
| The Output Layer | 281 |
| Building a Neural Network | 282 |
| Available Frameworks | 282 |
| Your First Neural Network in Keras | 284 |
| Neural Networks versus Other Algorithms | 287 |
| Summary | 289 |
| Chapter 18 Other Types of Neural Networks | 291 |
| Common Issues of Feed-Forward Neural Networks | 291 |
| Recurrent Neural Networks | 292 |
| Anatomy of a Stateful Neural Network | 292 |
| LSTM Neural Networks | 295 |
| Convolutional Neural Networks | 298 |
| Image Classification and Recognition | 298 |
| The Convolutional Layer | 299 |
| The Pooling Layer | 301 |
| The Fully Connected Layer | 303 |
| Further Neural Network Developments | 304 |
| Generative Adversarial Neural Networks | 304 |
| Auto-Encoders | 305 |
| Summary | 307 |

Chapter 19 Sentiment Analysis: An End-to-End Solution 309

Preparing Data for Training 310

- Formalizing the Problem 310
- Getting the Data 311
- Manipulating the Data 311
- Considerations on the Intermediate Format 313

Training the Model 313

- Choosing the Ecosystem 314
- Building a Dictionary of Words 314
- Choosing the Trainer 315
- Other Aspects of the Network 319

The Client Application 321

- Getting Input for the Model 321
- Getting the Prediction from the Model 322
- Turning the Response into Usable Information 323

Summary 323

PART V FINAL THOUGHTS

Chapter 20 AI Cloud Services for the Real World 327

Azure Cognitive Services 327

Azure Machine Learning Studio 329

- Azure Machine Learning Service 331
- Data Science Virtual Machines 333

On-Premises Services 333

- SQL Server Machine Learning Services 333
- Machine Learning Server 334

Microsoft Data Processing Services 334

- Azure Data Lake 334
- Azure Databricks 334
- Azure HDInsight 335
- .NET for Apache Spark 335

| | |
|--------------------------|-----|
| Azure Data Share | 336 |
| Azure Data Factory | 336 |
| Summary | 336 |

Chapter 21 The Business Perception of AI 339

| | |
|--|-----|
| Perception of AI in the Industry | 339 |
| Realizing the Potential | 339 |
| What Artificial Intelligence Can Do for You | 340 |
| Challenges Around the Corner | 342 |
| End-to-End Solutions | 343 |
| Let's Just Call It Consulting | 344 |
| The Borderline Between Software and Data Science | 344 |
| Agile AI | 346 |
| Summary | 349 |

| | |
|--------------------|-----|
| <i>Index</i> | 351 |
|--------------------|-----|

This page intentionally left blank

Acknowledgments

Writing a book with your son is a special experience even when it's the umpteenth book you write. For this one, I just put down in (hopefully clear) words Francesco's thoughts, vision, and his deep, and largely unexplained, understanding of machine learning. I definitely learned a lot from writing as much as I hope you will learn from reading.

If I learned a lot, well, that was mostly because of two people.

It's not the first time I have done some writing under the technical supervision of Cesar De la Torre Llorente, and it's always been a heavenly experience. I love his pragmatism and accuracy in devising, before designing, software products. He's currently principal program manager on the .NET product group at Microsoft and is in charge of the development of ML.NET. This is not specifically a book on ML.NET, but if the parts of the book that illustrate the .NET way to machine learning are accurate, well, that's because of the great help we received from Cesar.

There's an aspect in renewable energy that is little known: you need intelligent software to make it happen. At least on a functional level, it is vital to make accurate production, outage, fault, and price forecasts. Now, I don't think there are many people on this planet with a decade's worth of experience in this area that only recently was appointed the label *artificial intelligence*. Tiago Santos has been our guide in the random forest of machine learning and real-world artificial intelligence. "AI is just software" is now our shared motto.

If I've been able to give my career yet another turn (from Windows to web development and from software architecture to machine learning), it's also because two other people keep my creativity constantly stimulated. Giorgio Garcia-Agreda of Crionet made real my dreams as a tennis fanatic come true, up to singing "Easy like Sunday morning" in front of the tennis bigwigs. Simone Massaro of BaxEnergy discovered a fascinating new space where my renewable energy as a thinker can be freely expressed, sometimes even in front of top managers.

Any book is the result of teamwork, and it is our pleasure to call out the names of those who ultimately made it possible: Loretta Yates, as the acquisition editor; Charvi Arora, the managing editor, and Tonya Simpson, production editor.

—Dino

I finished high school one year early, and all I wanted was some money to practice as a professional investor. I had the wrong parents, though, and that approach didn't work. So, I asked my dad how to make money. "That's your problem," he said. "I can only teach you all I know." So, he taught me how to do things right and forgot to teach me how to do it wrong. As a result, today, we make the same mistakes in software. At this point, with some money in my hands, I was blissfully neglecting college when, on a hot summer afternoon, my dad told me, "Be honest: if you don't want to train your brain further, resign from college." As a result, a month later, I was back in class with a radically different mindset. I love mathematics, and I can't live doing anything different from it.

Then I met Gianfranco—friend, business partner, father, grandfather. He's a real professional investor, and he too taught me how to do things right and forgot to teach me how to do it wrong. As a result, today, we make the same mistakes in finance.

At school, at work, in the stock market, I study and try things. Sometimes they work, sometimes they don't, and whenever they don't, I learn something. It's the stick and carrot principle: the essence of learning for humans and machines. This book stems from my obsession for mathematical rigor and my dad's obsession for clarity. We used the stick on ourselves during the writing to ensure that carrots would be available during the reading.

This book is for you, Mom, because you'd love me anyway, regardless of triumph, disaster, or other impostors. This is for you, Maicol, because you'd love me even more if I stopped making noise on Sunday mornings.

This is for you, Alessandro, because you remind me when it's time to stop, and for you, Antonino, because you remind me of when I was too much of a smartass to be nice. This is for you, Sara, because you always give me a place to go the day before Christmas. And for you, Giorgio, because I'll always be a junior in front of you. This is for you, Grandma Concetta and Grandpa Salvatore, for the sausages, and for you, Grandma Leda, for your being as lively as any of us youngsters.

This is for you, Tiago, because we met only once to date, but enough to learn how much I have to learn from you.

This is also for all those I couldn't mention, including any of my present and future loves, so very complicated that would deserve a book of its own! And this is for me too, to help understand what I want to be.

—Francesco

About the Authors

DINO ESPOSITO



If I look back, I count more 20 books authored and 1000+ articles in a 25-year-long career. I've been writing the "Cutting Edge" column for *MSDN Magazine* month after month for 22 consecutive years. It is commonly recognized that such books and articles have helped the professional growth of thousands of .NET and ASP.NET developers and software architects worldwide.

After I escaped a dreadful COBOL project, in 1992 I started as a C developer, and since then, I have witnessed MFC and ATL, COM and DCOM, the debut of .NET, the rise and fall of Silverlight, and the ups and downs of various architectural patterns. In 1995 I led a team of five dreamers who actually deployed things that today we would call Google Photos and Shutterstock—desktop applications capable of dealing with photos stored in a virtual place that nobody had called the cloud yet. Since 2003 I have written Microsoft Press books about ASP.NET and also authored the best-seller *Microsoft .NET: Architecting Applications for the Enterprise*. I have a few successful Pluralsight courses on .NET architecture, ASP.NET MVC UI, and, recently, ML.NET. As architect of most of the backoffice applications that keep the professional tennis world tour running, I've been focusing on renewable energy, IoT, and artificial intelligence for the past two years as the corporate digital strategist at BaxEnergy.

You can get in touch with me through <https://youbiquitous.net> or twitter.com/despos, or you can connect to my LinkedIn network.

FRANCESCO ESPOSITO



I was 12 or so in the early days of the Windows Phone launch, and I absolutely wanted one of those devices in my hands. I could have asked Dad or Mom to buy it, but I didn't know how they would react. As a normal teenager, I had exactly zero chance of having someone buy it for me. So, I found out I was quite good at making sense of programming languages and impressed some folks at Microsoft enough to have a device to test. A Windows Phone was only the beginning; then came my insane passion

for iOS and, later, the shortcuts of C#.

The current part of my life began when I graduated from high school, one year earlier than expected. By the way, only 0.006 percent of students do that in Italy. I felt as powerful as a semi-god and enrolled in mathematics. I failed my first exams, and the shock put me at work day and night on ASP.NET as a self punishment. I founded my small software company, Youbiquitous, and began living on my own money. In 2017, my innate love for mathematics was resurrected and put me back on track with studies and led me to take the plunge in financial investments and machine learning.

This book, then, is the natural consequence of the end of my childhood. I wanted to give something back to my dad and help him make sense of the deep mathematics behind neural networks and algorithms. By the way, I have a dream: developing a super-theory of intelligence that would mathematically explore why the artificial intelligence of today works and where we can go further.

You can get in touch with me at <https://youbiquitous.net>.

Introduction

We need men who can dream of things that never were, and ask why not.

—John F. Kennedy, Speech to the Irish Parliament, June 1963

There are two views of artificial intelligence that people face today, and they are nonexclusive. One is the view pushed and pursued by the vast majority of media; the other is the view pushed and pursued by the IT community. In both camps, there are some true experts and some true pundits.

The view pushed by media focuses on the impact that artificial intelligence as a whole, in known and yet-to-know forms, may possibly have on our lives in some unfathomable future. The view pushed by the IT community (where software and data science experts belong) presents machine learning as the foundation of a new generation of software services that are just more intelligent than current services.

In the middle ground between the mass of people that the media reach and the much smaller IT community sits the patrol of cloud giants. They're the ones who conduct research and move the state of the art one step further every day, releasing new services for everyone to potentially add intelligence to new and existing applications.

At the base of the artificial intelligence pyramid sit managers and executives. On one hand, they're eager to apply to business those stunning services they hear from the tech news to edge out their competitors. On the other hand, they face the staggering bills of the projects they embarked on with the best of hopes.

- Artificial intelligence is not a magic wand.
- Artificial intelligence is not a service to pay per use. Worse yet, it's neither a capital nor operating expenditure.
- Artificial intelligence is just software.

Any business decision about artificial intelligence is better if made through the lens of software development: set requirements, get a reliable partner, put a budget on the table, work, start again in full respect of agility.

Is it that easy, then?

While artificial intelligence is about software development, it's not exactly the same as building an e-commerce website or a booking platform.

- Don't embark on artificial intelligence projects if you don't have a clear idea of the problem to solve, the context of it, and the point(s) to make.
- Don't embark on ambitious and adventurous projects by following the sole example of your closest competitor.
- Don't embark on such projects if you're not ready to lose some good money.

Just address one pain point at a time, build a cross-functional team, and provide full access to data.

Who Should Read This Book?

In the preparation of this book, we received a lot of feedback about the structure and elaborated on it quite a few times. We radically changed the table of contents at least three times. The hard part is that we devised this book to be unique and innovative, pursuing an idea of machine learning and software development a bit far away from the reality we see. Hopefully, our vision is the vision of machine learning that comes from the near future!

We see machine learning bounded within the fences of data science, as an artifact to be delivered to developers to embed it into some web service or desktop application. This is waterfall—no more no less. Where is all the agile that companies and enterprises constantly talk about? Agile ML means that data scientists and developers work together, and business analysts and domain experts join the team. And data stakeholders—whether it's IT or DevOps or whatever else—also join to facilitate data access and manipulation. This is agile teamwork—no more, no less.

We see the (business) need of a convergence of skills—from data science to software development and from software development to data science. This entry-level book is good for both sides of the pipeline. It talks to developers and shows ML.NET in action (over Python and along with Python) before getting into the analysis of the mechanics of machine learning algorithms. It also talks to data scientists who need to learn more about software needs.

This book is ideal if you're a software developer willing to add data science and machine learning skills to your arsenal. It's also ideal if you're a data scientist willing to learn more about software. Both categories, though, need to learn more and more about the other.

This is the bet of this book. We've classified it as "introductory" because it expands in width instead of going deep. It provides .NET examples because we think that, while the Python ecosystem is rich and thriving, there's no reason not to look around for platforms that allow you to do some machine learning closer to the bare metal of software applications, software services, and microservices—where ultimately any learning pipeline (including TensorFlow, PyTorch, handcrafted Python code) ends up being used.

Who Should Not Read This Book?

This is an introductory-level book specifically devised to give a broad but clear and accurate overview of machine learning using the ML.NET platform for experimenting. If you're looking for tons of Python examples, this book is not ideal. If you're looking for how-to examples to copy and paste in your solutions, whether Python or ML.NET, we're not sure this book is ideal. If you're looking for the nitty-gritty details of the mathematics behind algorithms or for an annotated overview of some implementations of algorithms, again, this book may not be ideal. (We do include some mathematics, but we still only scratch the surface.)

Organization of This Book

This book is divided into five sections. Part I, "Laying the Groundwork of Machine Learning," provides a quick overview of the foundation of artificial intelligence, intelligent software, and the basic steps of any machine learning project within end-to-end solutions. Part II, "Machine Learning in .NET," focuses on the ML.NET library and outlines its core parts, such as tasks for data processing, training, and evaluation in the context of common problems such as regression and classification. Part III, "Fundamentals of Shallow Learning," touches on the mathematical details of families of algorithms commonly trained to solve real-life problems: regressors, decision trees, ensemble methods, Bayesian classifiers, support vector machines, K-means, online gradients. Part IV, "Fundamentals of Deep Learning," is dedicated to neural networks that may come into play when none of the previous algorithms are found suitable. Finally, Part V, "Final Thoughts," is about the business vision of artificial intelligence in general and machine learning in particular, and it provides a cursory review of the runtime services for data processing and computation made available by cloud platforms, specifically the Azure platform.

Code Samples

All the code illustrated in the book, including possible errata and extensions, can be found at MicrosoftPressStore.com/IntroMachineLearning/downloads.

PART I

Laying the Groundwork of Machine Learning

| | | |
|------------------|---|----|
| CHAPTER 1 | How Humans Learn | 3 |
| CHAPTER 2 | Intelligent Software | 23 |
| CHAPTER 3 | Mapping Problems and Algorithms | 33 |
| CHAPTER 4 | General Steps for a Machine Learning Solution ... | 49 |
| CHAPTER 5 | The Data Factor | 67 |

Mapping Problems and Algorithms

When I consider what people generally want in calculating, I found that it always is a number.

—*Muḥammad ibn Mūsā al-Khwārizmī*

Persian mathematician of eighth century whose name originated the word algorithm

More often than not, the user experience produced by machine learning looks like magic to users. At the end of the day, though, machine learning is only a new flavor of software—a new specialty much like web or database development—and a flavor of software that today is a real breakthrough.

A breakthrough technology is any technology that enables people to do things that weren't possible before. Behind the apparent magic of final effects, however, there is a series of cumbersome tasks and, more than everything else, there's a series of sequential and interconnected decisions along the way that are hard to make and time consuming. In a nutshell, they are critical decisions for the success of the solution.

This chapter has two purposes. First, it identifies the classes of problems that machine learning can realistically address and the algorithms known to be appropriate for each class. Second, it introduces a relatively new approach—automated machine learning or AutoML for short—that can automate the selection of the best machine learning pipeline for a given problem and a given dataset.

In this chapter, we'll describe classes of problems and classes of algorithms. We'll focus on the building blocks of a learning pipeline in the next chapter.

Fundamental Problems

As you saw in Chapter 2, "Intelligent Software," the whole area of machine-based learning can be split into supervised and unsupervised learning. It's an abstract partition of the space of algorithms, and the main discriminant for being supervised or unsupervised is whether or not the initial dataset includes valid answers. Put another way, we can reduce automated learning into the union of two learning approaches—*learning by example* (supervised) and *learning by discovery* (unsupervised).

Under these two forms of learning, we can identify a number of general problems and for each a number of general algorithms. This layout is reflected in the organization of any machine learning software development library you can find out there and use—whether it’s based on Python, Java, or .NET.



Note Not coincidentally, most of the topics covered in the following chapters match, to a large extent, the tasks of the newest Microsoft’s ML.NET framework (covered in Part II, “Machine Learning in .NET”) and algorithm cheat-sheet of scikit-learn—an extremely popular machine learning Python library. (See <https://scikit-learn.org>.)

Classifying Objects

The classification problem is about identifying the category an object belongs to. In this context, an object is a data item and is fully represented by an array of values (known as *features*). Each value refers to a measurable property that makes sense to consider in the scenario under analysis. It is key to note that classification can predict values only in a discrete, categorical set.

Variations of the Problem

The actual rules that govern the object-to-category mapping process lead to slightly different variations of the classification problem and subsequently different implementation tasks.

Binary Classification. The algorithm has to assign the processed object to one of only two possible categories. An example is deciding whether, based on a battery of tests for a particular disease, a patient should be placed in the “disease” or “no-disease” group.

Multiclass Classification. The algorithm has to assign the processed object to one of many possible categories. Each object can be assigned to one and only one category. For example, classifying the competency of a candidate, it can be any of poor/sufficient/good/great but not any two at the same time.

Multilabel Classification. The algorithm is expected to provide an array of categories (or labels) that the object belongs to. An example is how to classify a blog post. It can be about sports, technology, and perhaps politics at the same time.

Anomaly Detection. The algorithm aims to spot objects in the dataset whose property values are significantly different from the values of the majority of other objects. Those anomalies are also often referred to as *outliers*.

Commonly Used Algorithms

At the highest level of abstraction, classification is the process of predicting the group to which a given data item belongs. In stricter math terms, a classification algorithm is a function that maps input variables to discrete output variables. (See Figure 3-1.)

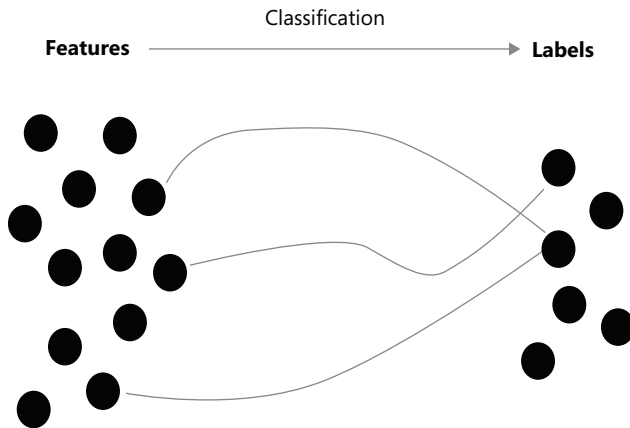


FIGURE 3-1 A graphical representation of a classification function

The classes of algorithms most commonly used for classification problems are as follows:

Decision Tree. A decision tree is a tailor-made binary tree that implements a sequence of rules to be progressively applied to each input object. Each leaf of the tree represents one of the possible output categories. Along the way, the input object is routed downward through the levels of the tree based on rules set at each node. Each rule is based on a possible value of one of the features. In other words, at each step, the key feature value of the input object (say, Age) is checked against the set value (say, 40), and the visit proceeds in the subtree that applies (say, less than or greater than or equal to 40). The number of nodes and the feature/value rules implemented are determined during the training of the algorithm.

Random Forest. This is a more specialized version of the decision tree algorithm. Instead of a single tree, the algorithm uses a forest of simpler trees trained differently and then provides a response that is some average of all the responses obtained.

Support Vector Machine. Conceptually, this algorithm represents the input values as points in an n -dimensional space and looks for a sufficiently wide gap between points. In two dimensions, you can imagine the algorithm looking for a curve that cuts the plane in two, leaving as much space as possible along the margin. In three dimensions, you can think of a plane that cuts the space in two.

Naïve Bayes. This algorithm works by computing the probability that a given object, given its values, may fall in one of the predefined categories. The algorithm is based on Bayes' theorem, which describes the likelihood of an event given some related conditions.

Logistic Regression. This algorithm calculates the probability of an object falling in a given category given its properties. The algorithm uses a sigmoid (logistic) function that, for its mathematical nature, lends itself well to be optimized to calculate a probability very close to 1 (or very close to 0). For this reason, the algorithm works well in either/or scenarios, and so it is mostly used in binary classification.

The preceding list is not exhaustive but includes the most-used classes of algorithms battle-tested for classification problems.



Important In the everyday jargon of machine learning, the term *algorithm* commonly refers to an entire family of algorithms that share the same general approach to the solution but may differ on a number of minor and not-so-minor implementation details. If you want to refer to a specific implementation of an algorithm, the term *trainer* (or even the term *estimator*) is more common. The term *pipeline*, instead, refers to the overall combination of data transformations, trainers, and evaluators that form the ultimately deployed machine learning *model*.

Common Problems Addressed by Classification

A number of real-life problems can be modeled as classification problems, whether binary, multiclass, or multilabel. Again, the following list can't and won't be exhaustive, but it is enough to give a clue about where to look when a concrete business issue surfaces:

- Spam and customer churn detection
- Data ranking and sentiment analysis
- Early diagnosis of a disease from medical images
- A recommender system built for customers
- News tagging
- Fraud or fault detection

Spam detection can be seen as a binary classification problem: an email is spam or is not. The same can be said for early diagnosis solutions although in this case the nature of the input data—images instead of records of data—requires a more sophisticated pipeline and probably would be solved using a neural network rather than any of the algorithms described earlier. Customer churn detection and sentiment analysis are multiclass problems, whereas news tagging and recommenders are multilabel problems. Finally, fraud or fault detection can be catalogued as an anomaly detection problem.

Predicting Results

Many would associate artificial intelligence with the ability to make smart predictions about future events. In spite of appearances, prediction is not magic but the result of a few statistical techniques, the most relevant of which is regression analysis. Regression measures the strength of a relationship set between one output variable and a series of input variables.

Regression is a supervised technique and is used to predict a continuous value (as opposed to discrete categorical values of classification).

Variations of the Problem

Regression is about finding a mathematical function that captures the relationship between input and output values. What kind of function? Different formulations of the regression function lead to different variations of the regression problem. Here are some macro areas:

Linear Regression. The algorithm seeks a linear, straight-line function so that all values, present and future, plot around it. The linear regression algorithm is fairly simple and, to a large extent, even unrealistic because, in practice, it means that a single value guides the prediction. Any realistic predictive scenarios, instead, bring in several different input data flows.

Multilinear Regression. In this case, the regression function responsible for the actual prediction is based on a larger number of input parameters. This fits in a much smoother way into the real world because to predict the price of a house, for example, you would use not only square footage but also historical trends, neighborhood, rooms, age, and maybe more factors.

Polynomial Regression. The relationship between the input values and the predicted value is modeled as an n th degree polynomial in one of the input values. In this regard, polynomial regression is a special case of multilinear regression and is useful when data scientists have reasons to hypothesize a curvilinear relationship.

Nonlinear Regression. Any techniques that need a nonlinear curve to describe the trend of the output value given a set of input data fall under the umbrella of nonlinear regression.

Commonly Used Algorithms

The solution to a regression problem is finding the curve that best follows the trend of input data. Needless to say, the training phase of the algorithm works on training data, but the deployed model, instead, needs to perform well on similar live data. The curve that predicts the output value based on the input is the curve that minimizes a given error function. The various algorithms define the error function in different ways and measure the error in different ways.

The classes of algorithms most commonly used for regression problems are as follows:

Gradient Descent. The gradient descent algorithm is expected to return the coefficients that minimize an error function. It works iteratively by first assigning default values to the coefficient and then measuring the error. If the error is large, it then looks at the gradient of the function and moves ahead in that direction, determining new values for the coefficients. It repeats the step until some stop condition is met.

Stochastic Dual Coordinate Ascent. This algorithm takes a different approach and essentially solves a dual problem—maximizing the value calculated by the function rather than minimizing the error. It doesn't use the gradient but proceeds along each axis until it finds a maximum and then moves to the next axis.

Regression Decision Tree. This algorithm builds a decision tree, as discussed previously, for classification problems. The main differences are the type of the error function used to decide

if the tree is deep enough and the way in which the feature value in each node is chosen (in this case, it is the mean of all values).

Gradient Boosting Machine. This algorithm combines multiple weaker algorithms (e.g., most commonly, a basic decision tree) and builds a unified, stronger learner. Typically, the prediction results from the weighed combination of the output of all the chained weak learners. Extremely popular algorithms in this class are XGBoost and LightGBM.



Important Both regression and classification cover very large areas of real-life problems. And often the actual problems faced can't be solved with any of these algorithms. Instead, they require a deeper learning approach via some neural network.

Common Problems Addressed by Regression

Regression is the task of predicting a continuous value, whether a quantity, a price, or a temperature.

- Price prediction (houses, stocks, taxi fares, energy)
- Production prediction (food, goods, energy, availability of water)
- Income prediction
- Time series forecasting

Time series regression is interesting because it can help understand and, better yet, predict the behavior of sophisticated dynamic systems that periodically report their status. This is fairly common in industrial plants where, even thanks to Internet of Things (IoT) devices, there's plenty of observational data. Time series regression is also commonly used in the forecasts of financial, industrial, and medical systems.

Grouping Objects

In machine learning, clustering refers to the grouping of objects represented as a set of input values. A clustering algorithm will place each object point into a specific group based on the assumption that objects in the same group have similar properties and objects in different groups have quite dissimilar properties.

At first, clustering may look like classification, and in fact, both problems are about deciding the category that a given data item belongs to. There's one key difference between the two, however. A clustering algorithm receives no guidance from the training dataset about the possible target groups. In other words, clustering is a form of unsupervised learning, and the algorithm is left alone to figure out how many groups the available dataset can be split on.

A clustering algorithm processes a dataset and returns an array of subsets. Those subsets receive no labels and no clues about the content from the algorithm itself. Any further analysis is left to the data science team. (See Figure 3-2.)

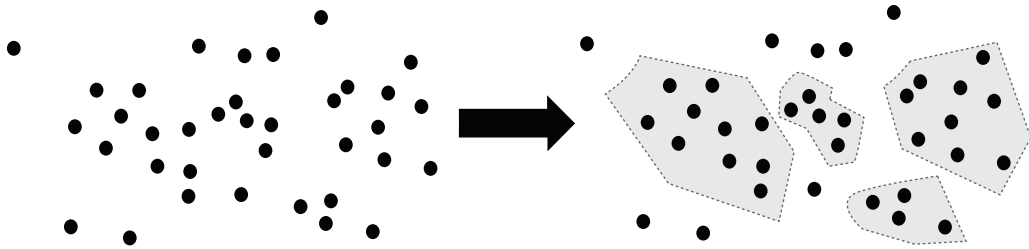


FIGURE 3-2 The final outcome of a clustering algorithm run on a dataset

Commonly Used Algorithms

The essence of clustering is analyzing data and identifying as many relevant clusters of data as it can find. While the idea of a cluster is fairly intuitive—a group of correlated data items—it still needs some formal definition of the concept of correlation to be concretely applied. In the end, a clustering algorithm looks for disjoint areas (not necessarily partitions) of the data space that contain data items with some sort of similarity.

This fact leads straight to another noticeable difference between clustering and regression or classification. You'll never deploy a clustering model in production and never run it on live data to get a label or a prediction. Instead, you may use the clustering step to make sense of the available data and plan some further supervised learning pipeline.

Clustering algorithms adopt one of the following approaches: partition-based, density-based, or hierarchy-based. Here are the most popular algorithms:

K-Means. This partition-based algorithm sets a fixed number of clusters (according to some preliminary data analysis) and randomly defines their data center. Next, it goes through the entire dataset and calculates the distance between each point and each of the data centers. The point finds its place in the cluster whose center is the nearest. The algorithm proceeds iteratively and recalculates the data center at each step.

Mean-Shift. This partition-based algorithm defines a circular sliding window (with arbitrary radius) and initially centers it at a random point. At each step, the algorithm shifts the center point of the window to the mean of the points within the radius. The method converges when no better center point is found. The process is repeated until all points fall in a window and overlapping windows are resolved, keeping only the window with the most points.

DBSCAN. This density-based algorithm starts from the first unvisited point in the dataset and includes all points located within a given range in a new cluster. If too few points are found, the point is marked as an outlier for the current iteration. Otherwise, all points within a given range of each point currently in the cluster are recursively added to the cluster. Iterations continue until there's at least one point not included in any cluster or their number is so small that it's OK to ignore them.

Agglomerative Hierarchical Clustering. This hierarchy-based algorithm initially treats each point as a cluster and proceeds iteratively, combining clusters that are close enough to a given distance metric. Technically, the algorithm would end when all the points fit in a single cluster, which would be the same as the original dataset. Needless to say, you can set a maximum number of iterations or use any other logic to decide when to stop merging clusters.

K-Means is by far the simplest and fastest algorithm, but, in some way, it violates the core principle of clustering because it sets a fixed number of groups. So, in the end, it's halfway between classification and clustering. In general, clustering algorithms have a linear complexity, with the notable exception of hierarchy-based methods. Not all algorithms, however, produce the same quality of clustering regardless of the distribution of the dataset. DBSCAN, for example, doesn't perform as well as others when the clusters are of varying density, but it's more efficient than, say, partition-based methods in the detection of outliers.

Common Problems Addressed by Clustering

Clustering is the method of many crucial business tasks in a number of different fields, including marketing, biology, insurance, and in general wherever screening of population, habits, numbers, media content, or text is relevant.

- Tagging digital content (videos, music, images, blog posts)
- Regrouping books and news based on author, topics, and other valuable information
- Discovering customer segments for marketing purposes
- Identifying suspicious fraudulent finance or insurance operations
- Performing geographical analysis for city planning or energy power plant planning

It is remarkable to consider that clustering solutions are often used in combination with a classification system. Clustering may be first used to find a reasonable number of categories for the data expected in production, and then a classification method could be employed on the identified clusters. In this case, categories will be manually labeled, looking at the content of identified clusters. In addition, the clustering method might be periodically rerun on a larger and updated dataset to see whether a better categorization of the content is possible.

More Complex Problems

Classification, regression, and clustering algorithms are sometimes referred to as *shallow learning*, in contrast to *deep learning*. Admittedly, the distinction between shallow learning and deep learning is a bit sketchy and cursory; yet, it marks the point of separating problems that can be solved with a relatively straight algorithm from those that require the introduction of some flavor of neural networks (more or less deep in terms of constituent layers) or the pipelining of multiple straight algorithms. Typically, these problems revolve around the area of cognition such as computer vision, creative work, and speech synthesis.

Image Classification

Image processing began in the late 1960s when a group of NASA scientists had the problem of converting analogic signals to digital images. The core of image processing is the simple application of mathematical functions to a matrix of pixels. A much more enhanced form of image processing is computer vision.

Computer vision isn't limited to processing data points but attempts to recognize patterns of pixels and how they match to forms (objects, animals, persons) in the real world. Computer vision is the branch of machine learning devoted to the emulation of the human eye, capable of capturing images and recognizing and classifying them based on properties such as size, color, and luminosity.

In the realm of computer vision, image classification is one of the most interesting sectors, especially for its applications to sensitive fields such as health care and security. Image classification is the process of taking a picture (or a video frame), analyzing it, and producing a response in the form of a categorical value (it's a dog) or a set of probabilistic values (70 percent, it's a dog; 20 percent, it's a wolf; 10 percent, it's a fox). In much the same way, an image classifier can guess mood, attitude, or even pain.

Even though many existing cloud services can recognize and classify images (even video frames), the problem of image classification can hardly be tackled outside a specific business context. In other words, you can hardly take a generic public cloud cognitive service and use it to process medical images (of a certain type) or monitor the live stream of a public camera. You need specific training for the algorithm tailor-made for the scenario you're facing.

An image classifier is typically a convolutional multilayer neural network. In such a software environment, each processing node receives input from the previous layers and passes processed data to the next. Depending on the number (and type) of layers, the resulting algorithm proves able (or not so able) to do certain things.

Object Detection

A side aspect of computer vision, tightly related to image classification, is object detection. With image classification, you can rely on a class of algorithms capable of looking at live streams of pictures and recognize elements in it. In other words, image classification can tell you what is in the processed picture. Object detection goes one step further and operates a sort of multiclass classification of the picture, telling about all the forms recognized and also about their relative position.

Object detection is very hot in technologies like self-driving cars and robotics. Advanced forms of object detection can also identify bounding boxes for the form to find and even draw precise boundaries around it. Object detection algorithms typically belong to either of two classes—classification-based or regression-based.

In this context, classification and regression don't refer to the straight shallow learning algorithms covered earlier in the chapter but relate to the learning approach taken by the neural network to come to a conclusion.

Text Analytics

Text analytics consists of parsing and tokenizing text, looking for patterns and trends. It is about learning relationships between named entities, performing lexical analysis, calculating and evaluating the frequency of words, and identifying sentence boundaries and lemmas. In a way, it's a statistical exercise of data mining and predictive analysis applied to text with the ultimate goal of taking software to interact with humans using the same natural language.

A typical application of text analytics is summarizing, indexing, and tagging the content of large digital free text databases and documents such as the comments (and complaints) left by customers of a public service. Text analytics often goes under the more expressive name of *natural language processing* (NLP) and is currently explored in more ambitious scenarios such as processing a live stream, performing speech recognition, and using recognized text for further parsing and information retrieval. Natural language processing applications are commonly built on top of neural networks in which the input text passes through multiple layers to be progressively parsed and tokenized until the networks produce a set of probabilistic intents.

There are quite a few applications of NLP available in the industry, buried in the folds of enterprise frameworks used in answering machine applications and call centers. However, if you just want to explore the power of the raw NLP, research a few of the existing test platforms, such as <https://knowledge-studio-demo.ng.bluemix.net>. The tool parses text, an excerpt of a police car accident report, and automatically extracts relevant facts, such as age of the involved people, characteristics of involved vehicles, location, and time.

Automated Machine Learning

Machine learning is a large field and is growing larger every day. As you'll see in much more detail in the next chapter, building an intelligent solution for a real-life business problem requires a *workflow* that essentially consists of a combination of different steps: data transformations, training algorithms, evaluation metrics, and, last but not least, domain knowledge, knowledge base, trial-and-error attitude, and imagination.

In this context, while the human ability to sort things out probably remains unparalleled, the community is seriously investigating the possibility of using automated, wizard-style tools to prepare a sketchy plan that could possibly represent the foundation of a true solution in a matter of minutes instead of days.

This is just the essence of the automated machine learning (AutoML) approach and consists of a framework that looks at your data and declared intent and intelligently suggests the steps to take that it determines most appropriate.

Aspects of an AutoML Platform

The typical end-to-end pipeline of any machine learning solution applied to a real-world problem most likely includes a number of steps, as outlined here:

- Preliminary analysis and cleaning of available data
- Identification of the properties (features) of the data that look most promising and relevant to solve the actual problem

- Selection of the algorithm
- Configuration of the parameters of the algorithm
- Definition of an appropriate validation model to measure the performance of the algorithm and indirectly the quality of the data it is set to use

Machine learning may not be for the faint-hearted, and even when one has a strong domain knowledge, the risk of feeling like a nonexpert newbie is fairly high.

Hence, AutoML is emerging as a solution to get people started quickly on machine learning projects and sometimes even effectively. AutoML offers the clear advantage of being fast and producing working solutions. The debatable point is not how objectively good the solution is that you can get out of an AutoML wizard, but the trade-off between what you get from AutoML and what you might be able to design by hand, especially if your team is not made up of domain and machine learning super-experts.



Note To some extent, the debate about the alleged superficiality of AutoML solutions recalls past debates about the use of high-level programming languages over Assembly and the use of system-managed memory over memory cells directly allocated by the programmer. Our frank opinion is that AutoML frameworks are excellent at doing their job on simple problems. They can't do much for complex problems, however. But unfortunately, as of today, most real-world problems are quite complex.

Common Features

An AutoML framework is made of two distinct parts: a public list of supported learning scenarios and an invisible runtime service that returns a deliverable model based on some input parameters. A learning scenario is essentially an expert subsystem designed to solve specific classes of problems using data in one of a few predefined formats. The runtime is a learning pipeline in which a set of predefined data transformations are performed on selected input given the learning objective; target features are selected; and the trainer is selected, configured, trained, and tested.

An AutoML framework will perform any of the following tasks in an automated way after the user has indicated the physical source of data (tabular files, relational databases, cloud-based data warehouses) and the learning objective:

- Preprocessing and loading of data from different formats including detection of missing and skewed values
- Understanding of the type of each dataset column to figure out whether the column is, say, a Boolean, a discrete number, a categorical value, or free text
- Application of built-in forms of feature engineering and selection, namely the addition or transformation of data columns in a way that makes particular sense for the learning objective
- Detection of the type of work required by the learning objective (binary classification, regression, anomaly detection) and selection of a range of most appropriate training algorithms

- Configuration of the hyperparameters of the selected training algorithms
- Training of the model, application of appropriate evaluation metrics, and testing of the model

In addition, an AutoML framework is also often capable of visualizing data and results in a fancy way that is also helpful to better understand the underpinnings of the problem at hand.

There are a couple of popular AutoML frameworks: one is from Google and one, the newest, from Microsoft. Let's first briefly examine the Google Cloud AutoML platform, and then we'll go for a deeper live demonstration of the Microsoft AutoML framework as integrated in Visual Studio 2019.

Google Cloud AutoML

The Google Cloud AutoML platform is located at <https://cloud.google.com/automl>. It comes as a suite of machine learning systems specifically designed to simplify as much as possible the building of models tailor-made for specific needs. The platform works much like a UI wizard and guides the user through the steps of selecting the scenario, data, and parameters and then does the apparent magic of returning a deployable artifact out of nowhere. Internally, the Google Cloud AutoML platform relies on Google's transfer learning technology, which allows the building of neural networks as the composition of predefined existing networks.

Google Cloud AutoML supports a few learning scenarios such as computer vision, object detection in videos and still images, and natural language processing and translation. As you can see, it's a group of pretty advanced and sophisticated scenarios. It also supports a simpler one, called AutoML Tables, that works on tabular datasets and tests multiple model types at the same time (regression, feedforward neural network, decision tree, ensemble methods).

Microsoft AutoML Model Builder

An AutoML framework is also integrated in Visual Studio 2019 and comes packaged with ML.NET—the newest Microsoft .NET-based library for machine learning. The AutoML Model Builder framework has both a visual, wizard-style interface in Visual Studio (more on this in a moment) and a command-line interface (CLI) for use from within command-based environments such as PowerShell. A quick but effective summary of AutoML CLI can be found at <https://bit.ly/2FaK7SP>.

In Microsoft's AutoML framework, developers choose a task, provide the data source, and indicate a maximum training duration. Needless to say, the selected maximum duration is a discriminant for the quality of the final model. The shorter time you choose, the less reliable the final model can be.



Note Compared to Google Cloud AutoML, the Microsoft AutoML solution currently focuses on simpler tasks and is available also on premise and then for shorter training cycles. The Google platform, instead, is cloud-based and suitable for longer and more realistic training cycles available through a paid subscription.

The AutoML Model Builder in Action

In Visual Studio 2019, after you install the latest version of the ML.NET Model Builder extension, you gain the ability to add a machine learning item to an existing project. When you do that, you're sent to a wizard like the one shown in Figure 3-3.

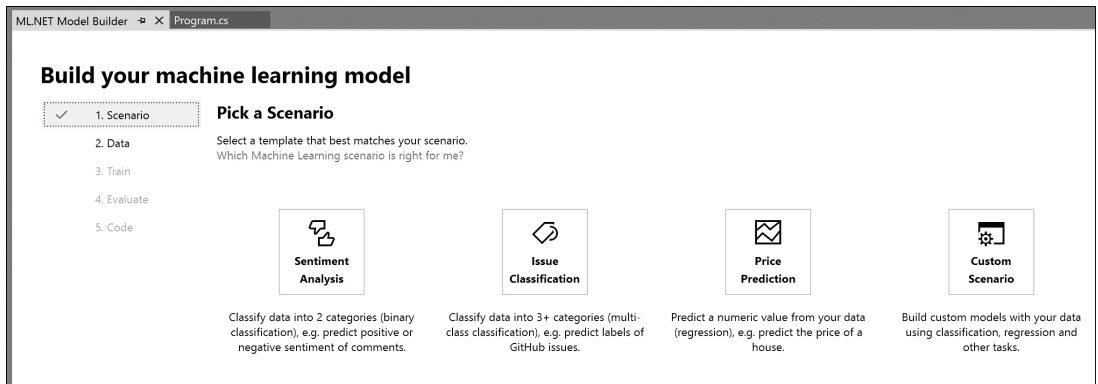


FIGURE 3-3 The main page of the Model Builder Visual Studio extension

As you can see, the wizard is articulated in five steps that broadly match the main steps of any machine learning pipeline. The first step of the builder is choosing the learning scenario—namely, the broad category of the problem for which you'd like to build a machine learning solution. In the version of the builder used for the test, the choice is not very large: Sentiment Analysis, Issue Classification, Price Prediction, and Custom Scenario. As an example, let's go for Price Prediction.

Exploring the Price Prediction Scenario

After you pick the scenario, the wizard asks you to load some data into the system. For the price prediction scenario, you can choose from a plain file or a SQL Server table. In the example shown in Figure 3-4, the loaded file is a CSV file. One key input to provide is the name of the column you want the final model to predict. In this case, the CSV file contains about one million rows, representing a taxi ride that really took place. The column to predict is the fare amount.

Training the Model

The third step is about the selection of the ideal trainer—the algorithm that is the most appropriate for the learning scenario and the data. This is where the power (and from a certain angle also the weakness) of the automated machine learning framework emerges. Some hard-coded logic, specific to the chosen scenario, tries a few training algorithms based on the allotted training time. Figure 3-5 shows an estimation of the training time necessary for a certain amount of data.

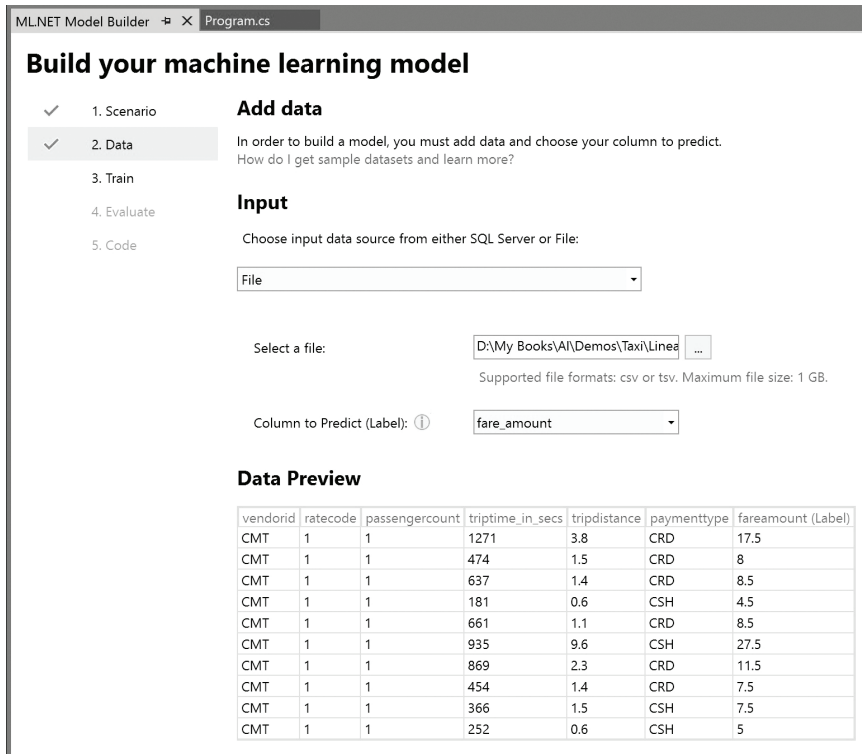


FIGURE 3-4 Loading data into the model

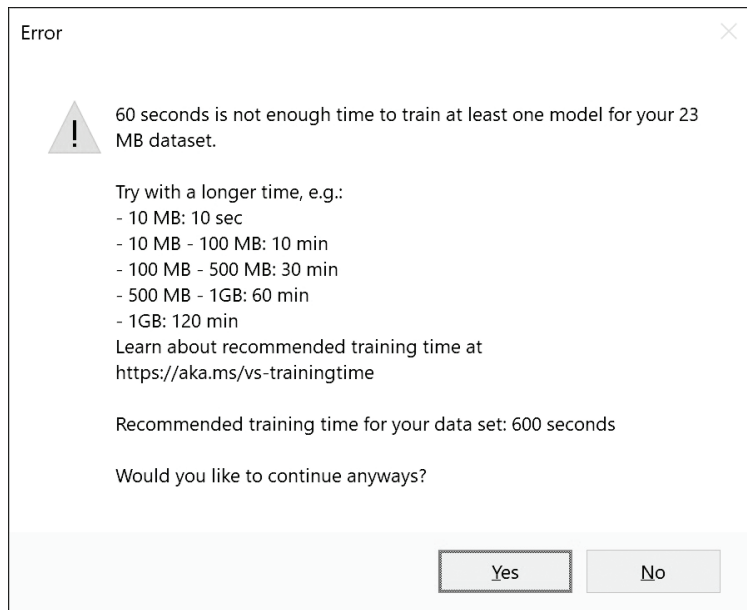


FIGURE 3-5 Estimating the training time

During the training phase, the system tries several different algorithms and uses an apt metric to evaluate its performance. (See Figure 3-6.)

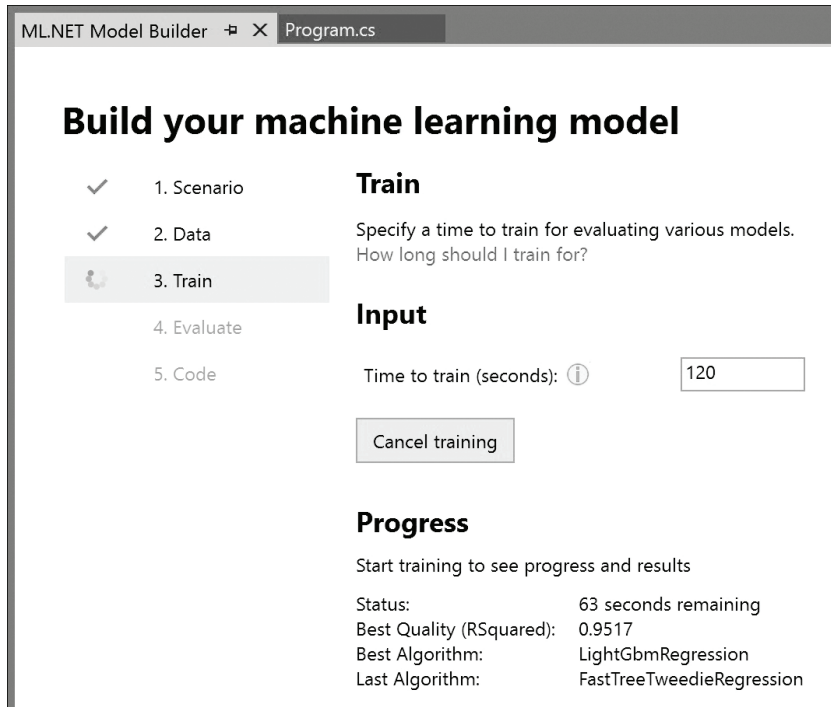


FIGURE 3-6 AutoML tries different algorithms and uses some metrics to evaluate the quality.

Evaluating the Results

At the end of the training, the AutoML system has data about a few algorithms it has tried with different hyperparameters. The metrics for evaluating the performance depend on the tasks and the algorithm. Price prediction is essentially a regression task for which the R-squared measure is the most commonly used. (We'll cover the math behind regression and R-squared in Chapter 11, "How to Make Simple Predictions: Linear Regression.") The theoretic ideal value of the R-squared metrics is 1; therefore, any value close enough to 1 is more than acceptable. Consider that in training, a resulting metric with a value of 1 (or very close to 1) is often the sign of overfitting—the model fits too much to the training data and potentially might not work effectively on live data once in production.

The AutoML process then suggests the use of the *LightGbmRegression* algorithm. If you want, you can just take the ZIP file with the final model ready for deployment. But what about looking into the actual set of data transformation and the actual code to possibly modify for further improvements?

The AutoML also offers the option to add the C# files to the current project for you to further edit them and retrain the model on a different dataset, for example. (See Figure 3-7.)

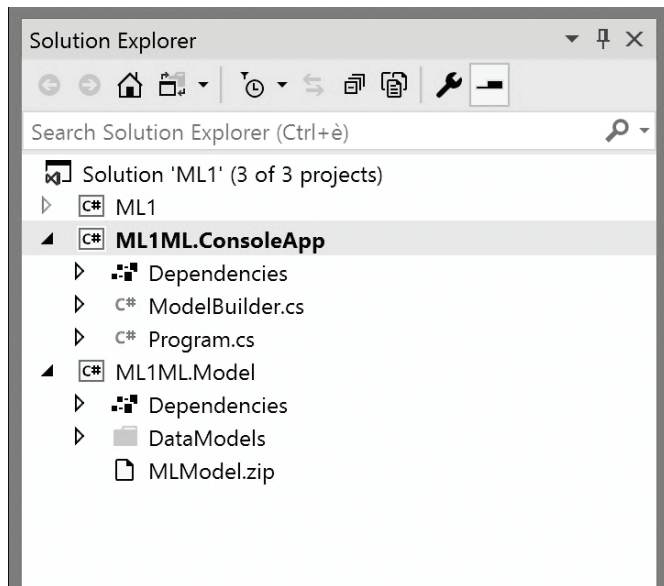


FIGURE 3-7 Autogenerated projects added by the Model Builder

As you can see, the figure contains two projects. One is a console application that contains a *ModelBuilder.cs* file packed with the code used to build the model. The other project is a class library and contains a sample client application seen as the foundation for using the model. This project also contains the actual model as a ZIP file.

Summary

Machine learning is ultimately intelligent software, but it is not the magic wand that movies and literature (and recently also sales/marketing departments) love to depict. More importantly, machine learning is not a physical black box you can pick from the shelves of a drugstore, bring home, mount, and use.

In the real world, you can't just "load data into the machine" and have the machine, in some way, just use it. In the real world, there are a few classes of approaches (mostly derived from statistics) such as regression, classification, and clustering and a bunch of concrete training algorithms. However, when to use which?

Determining which to use is a matter of experience and know-how, but it is also a matter of knowing data and how things actually work in the actual business domain. Does that mean that only experts can do machine learning? Yes, for the most part, that is just the point. However, nobody is born an expert, and everyone needs to get started in some way. This is the reason why automated tools for machine learning are emerging. In this chapter, we briefly looked at the Google Cloud AutoML and Visual Studio ML.NET Model Builder.

With the next chapter, we complete the preliminary path of machine learning, discussing the concept of a pipeline—namely, the sequence of steps that ultimately lead to the production of a deliverable model.

Index

A

- acceptance testing, 58
- accuracy of data, 71
- activation functions, 257–258, 274–277
 - linear, 274
 - ReLU, 276–277, 302
 - sigmoid, 261–262, 274–275
 - softmax, 275
 - step function as, 260
 - TanH, 275–276
- adaptability to change, 15–16
- adaptive boosting, 204–206
- agglomerative hierarchical clustering, 40
- agility in artificial intelligence (AI), 348
- AI. *See* artificial intelligence (AI)
- algorithms. *See also* models
 - for classification, 34–36
 - KNN algorithm, 230–234
 - SVM algorithm, 235–245
 - for clustering, 39–40
 - DBSCAN, 248–251
 - K-Means, 246–247
 - K-Modes, 247–248
 - decision trees
 - for classification, 186–193
 - design principles, 185
 - expert systems versus, 185
 - for regression, 194–195
 - ensemble methods, 198
 - bagging technique, 198–203
 - boosting technique, 203–210
 - for linear regression, 169–178
 - cost function identification, 170–171
 - evaluating, 178–180
 - gradient descent, 174–178
 - ordinary least square, 171–174
 - models versus, 58–59
 - for naïve Bayes classifiers, 219–220
 - for prediction, 37–38
 - training algorithms
 - backpropagation, 264–270
 - price prediction example (ML.NET pipeline), 96–97
 - selecting, 45–47, 59–61, 96–97
 - value of, measuring, 97
- al-Khwārizmī, Muḥammad ibn Mūsā, 33
- anomaly detection, 34
- Apache Spark, 335
- APIs, exposing, 65
- applied AI, 18, 23. *See also* expert systems
- aprioristic knowledge, 227
- architecture
 - of expert systems, 17–18
 - of human brain, 12–13
 - of ML.NET, 105
 - catalogs, 109–111
 - data representation, 107–109
 - types and interfaces, 105–107
 - of neural networks, 273–274
 - activation functions, 274–277
 - hidden layers, 277–280
 - output layer, 281
 - RNN (recurrent neural networks), 293–294
- arithmetic mean, 136–137
- artificial intelligence (AI). *See also* machine learning
 - agility in, 348
 - autonomous systems, 19–20
 - categories of, 20
 - examples of, 20
 - learning dimension, 19
 - supervised learning, 29–31
 - unsupervised learning, 27–29
 - challenges of, 342–343
 - cloud computing and, 344
 - end-to-end solutions, 343

artificial intelligence (AI)

- evolution of, 24–25
 - expert systems, 16–19, 25–26
 - examples of, 18
 - expertise versus intelligence, 25
 - history of, 17
 - internal architecture, 17–18
 - limitations of, 19
 - “Miracle on the Hudson” example, 25–26
 - updating, 26
 - history of, 6–7
 - hype surrounding, 25
 - perception of, 339
 - potential of, 339–340
 - primordial intelligence, 16–17
 - purpose of, 340
 - classification, 341
 - cognition, 341
 - content creation, 341
 - human emulation, 341
 - prediction, 340–341
 - rigidity in, 347–348
 - sentiment and, 20–21
 - types of, 16
 - waterfall methodology, 346–347
 - artificial neurons, 255–256
 - logistic neurons, 260–262
 - perceptrons, 257–260
 - activation function, 257–258
 - enabling learning, 260
 - feed-forward layers, 259–260
 - NAND gates, 258–259
 - Asimov, Isaac, 327
 - auto-encoders, 305–307
 - AutoML (automated machine learning), 42–48
 - features of, 43–44
 - Google Cloud AutoML, 44
 - Microsoft AutoML Model Builder, 44–48
 - in sentiment analysis, 315–316
 - autonomous systems, 19–20
 - categories of, 20
 - examples of, 20
 - learning dimension, 19
 - supervised learning, 29–31
 - inferred function, 31
 - labeled data, 30–31
 - prediction and classification, 29–30
 - unsupervised learning, 27–29
 - discovering data clusters, 27–28
 - evaluating data clusters, 29
 - availability of data, 69
 - average absolute deviation, 142
 - average pooling, 302
 - axons, 10–11
 - Azure Cognitive Services, 327–329
 - Azure Data Factory, 336
 - Azure Data Lake, 334
 - Azure Data Share, 336
 - Azure Databricks, 334–335
 - Azure HDInsight, 335
 - Azure Machine Learning Service (ML Service), 331–332
 - Azure Machine Learning Studio (ML Studio), 329–331
- ## B
- Babbage, Charles, 5
 - backpropagation algorithm, 264–270
 - bagging technique, 198–203
 - definition of, 198
 - random forests, 198–203
 - pros and cons, 202–203
 - steps in, 200–202
 - Ball Tree, 234
 - batch gradient, 177
 - Bayes, Thomas, 212
 - Bayes’ theorem, 214–215
 - Bayesian statistics, 211–216. *See also* naïve Bayes classifiers
 - Bayes’ theorem, 214–215
 - chain rule, 213
 - classification and, 216–218
 - conditional probability, 213
 - independent events, 213–214
 - intersection of events, 213
 - partitions of events, 214
 - sample scenario, 215–216
 - behavior learning, 14–15
 - bell curve, 224–225
 - Benedict XVI (pope), 105
 - Bernoulli Naïve Bayes (BNB), 223–224
 - bias (in statistics), 144–145
 - variance versus, 157–158
 - biased datasets, avoiding, 53–54
 - bidimensional case for cost function, 170–171
 - minimizing, 171–173
 - bidirectional LSTM, 316
 - bimodal datasets, 138
 - binary classification, 34

- in ML.NET, 111–116
 - data transformations on, 112–113
 - evaluating model, 115–116
 - sentiment analysis, 111–112
 - training model, 114–115
- neural networks versus, 287–289
- SVM algorithm as, 235
- binomial distribution, 223
- BNB (Bernoulli Naïve Bayes), 223–224
- Bohr, Niels, 93
- Boole, George, 4
- boosting technique, 203–210
 - adaptive boosting, 204–206
 - definition of, 198
 - gradient boosting, 204, 206–210
 - with imbalanced data, 203
- bootstrap technique, 200
- box plot, 146
- brain functionality. *See* human intelligence
- breakthrough technology, definition of, 33
- Brynjolfsson, Erik, 339

C

- Caffe, 284
- carrot-and-stick approach, 14–15
- CART (Classification and Regression Trees) algorithm, 187–191
- catalogs (ML.NET), 109–111
 - cross-cutting operation, 110–111
 - task-specific, 110
- categorical features, 136
- CDF (cumulative distribution function), 139–140
- centroids, 123
- cerebral cortex structure, 9–10
- chain rule, 213
- change, adaptability to, 15–16
- choosing. *See* selecting
- Church, Alonzo, 5
- Church-Turing thesis, 5
- classification, 34–36, 341
 - algorithms used, 34–36
 - Bayesian statistics and, 216–218. *See also* naïve
 - Bayes classifiers
 - confusion matrix and, 159–160
 - with decision trees, 181, 186–193
 - CART algorithm, 187–191
 - error function, 187
 - homogeneity, 186–187
 - ID3 algorithm, 191–193
 - sample scenario, 186
- definition of, 229
- of images, 41
 - data transformations on, 127–129
 - steps in, 127
 - training model, 129–131
- KNN algorithm, 230–234
 - Ball Tree, 234
 - brute-force implementation, 233
 - business scenario, 234
 - categorical data in, 232
 - distance calculations, 231–232
 - K-D tree, 233–234
 - number of neighbors, 230–231
 - training, 234
- in ML.NET, 111–121
 - binary classification, 111–116
 - multiclass classification, 116–121
- problems addressed by, 36
- regression versus, 165
- in supervised learning, 29–30
- SVM algorithm, 235–245
 - coefficients for prediction, 243
 - hyperplanes, 235–236
 - Lagrange multipliers, 240
 - linearly separable datasets, 238–239
 - mechanics of prediction, 240–241, 244
 - multiclass classification with, 244–245
 - nonlinearly separable datasets, 237–238
 - scalar product of vectors, 239–240
 - support vectors, 236–237
 - training, 242–243
 - vector operations, 239
- variations of, 34
- cleaning data, 53–54
- client application
 - for price prediction example (ML.NET), 99–103
 - designing user interface, 102–103
 - getting model file, 99
 - making predictions, 100–102
 - questioning data and problem approach, 103
 - setting up, 99–100
 - for sentiment analysis, 321
 - data collection, 321–322
 - output formatting, 323
 - prediction from, 322–323
- clock speed of human brain, 11

cloud computing, artificial intelligence (AI)

- cloud computing, artificial intelligence (AI) and, 344
 - clustering, 38–40. *See also* data clusters
 - algorithms used, 39–40
 - business scenario, 245–246
 - DBSCAN, 248–251
 - definition of, 229
 - K-Means, 246–247
 - ideal number of clusters, 247
 - steps in, 246–247
 - K-Modes, 247–248
 - in ML.NET, 122–126
 - data preparation, 122–123
 - evaluating model, 124–126
 - training model, 123–124
 - problems addressed by, 40
 - as unsupervised learning, 245
 - CNN (convolutional neural networks), 298–304
 - convolutional layer, 299–301
 - fully connected layer, 303–304
 - image classification in, 298
 - pooling layer, 301–303
 - code-breaking machines, 5
 - cognition, 341
 - cold data, 51
 - completeness of data, 70
 - complex relationships, 154–155
 - computed features, 56
 - The Computer and the Brain* (von Neumann), 11
 - computers, memory, 13
 - computing machines, formalizing, 5
 - Conda, 332
 - conditional probability, 213
 - conditional statements as primordial intelligence, 16–17
 - confusion matrix, 159–160
 - in multiclass classification, 121
 - consistency of data, 71
 - constrained optimization, 243
 - content creation, 341
 - continuous features, 136, 219
 - continuous training model, 53
 - convolutional layer in CNN (convolutional neural network), 299–301
 - convolutional neural networks (CNN). *See* CNN (convolutional neural networks)
 - correlation, 148
 - correlation analysis, 55
 - cost function
 - identifying, 170–171
 - minimizing, 171–174
 - covariance, 148
 - creative work in neural networks, 62
 - cross validation, 161
 - cross-cutting operation catalogs (ML.NET), 110–111
 - Cumberbatch, Benedict, 23
 - cumulative distribution function (CDF), 139–140
- ## D
- data accessibility, 51
 - data accuracy, 71
 - data availability, 69
 - data cleaning, 53–54
 - data clusters. *See also* clustering
 - discovering, 27–28
 - evaluating, 29
 - data collection, 50–52
 - data quality and, 69–70
 - data-driven culture, 50–51
 - for sentiment analysis, 311, 321–322
 - storage options, 51–52
 - data completeness, 70
 - data consistency, 71
 - data harmonization, 54
 - data integrity, 70–71
 - data lakes, 51–52, 69, 334
 - data ownership, 51
 - data preparation, 52–58, 162
 - cleaning data, 53–54
 - in clustering, 122–123
 - data quality improvement, 53
 - feature engineering, 54–56
 - in ML.NET, 88
 - normalization, 163
 - with Pandas library, 80
 - scaling, 162
 - for sentiment analysis, 310–313
 - data collection, 311
 - intermediate format for data transformations, 311–313
 - problem formalization, 310
 - standardization, 163
 - training dataset finalization, 56–58
 - data quality, 67–70
 - data collection and, 69–70
 - data validity, 68
 - improving, 53
 - data representation
 - in ML.NET, 107–109
 - for statistics, 145–150
 - five-number summary, 145–146

- histograms, 146–147
 - scale of plots, 149–150
 - scatter plot matrices, 148–149
 - scatter plots, 148
- data sampling, 69–70
- data science virtual machines (DSVMs), 333
- data scientists, 71–74
 - daily tasks, 72–73
 - definition of, 27
 - job description, 72
 - Python as language for, 79
 - software developers versus, 73–74, 344–346
 - tool chests, 73
- data timeliness, 71
- data transformations
 - in binary classification, 112–113
 - in clustering, 123
 - in multiclass classification, 117–118
 - price prediction example (ML.NET pipeline), 94–95
 - for sentiment analysis, 311–313
 - in transfer learning, 127–129
- data uniqueness, 70–71
- data validity, 68
- data views
 - row navigation, 108–109
 - schema of, 108
 - shuffling data, 109
- data visualization, 80–81
- data warehouses, 52, 69
- Databricks, 334–335
- data-driven culture, 50–51
- DBSCAN, 39, 248–251
- decision trees, 35
 - for classification, 186–193
 - CART algorithm, 187–191
 - error function, 187
 - homogeneity, 186–187
 - ID3 algorithm, 191–193
 - sample scenario, 186
 - definition of, 182
 - design principles, 185
 - examples of, 183–184
 - expert systems versus, 185
 - in machine learning, 183
 - for regression, 194–195
 - usages for, 181–182
- deduping, 54

- deep learning, 40
 - in ML.NET, 86–87
 - shallow learning versus, 289
- deep LSTM neural networks, 297–298
- deep RNN (recurrent neural networks), 295
- delta rule, 266
- DENDRAL, 17
- dendrites, 10
- density function, 225–226
- density-based clustering. *See* DBSCAN
- deploying models, 64–65
- designing user interfaces, 102–103
- detection of objects, 41
- dictionary of words, building, 314–315
- dimensionality reduction, 29, 56
- discovering data clusters, 27–28
- discrete features, 136, 219
- dissimilarity measure, 248
- distance calculations, 231–232
- documents, data within, 69
- dropout, 319–320
- DSVMs (data science virtual machines), 333
- dummy variables, 56

E

- ecosystems, selecting for sentiment analysis, 314
- edges, definition of, 182
- Einstein, Albert, 309
- Eisenhower, Dwight, 49
- elbow method, 247
- embedding layer, 306, 319
- encoding, 306
- end-to-end solutions
 - in artificial intelligence, 343
 - Python and, 82–83
- ensemble methods, 198
 - bagging technique, 198–203
 - boosting technique, 203–210
- entropy
 - definition of, 191
 - information gain and, 192
- epochs of training, 320
- error function for classification decision trees, 187
- estimator
 - bias and, 144
 - definition of, 152
- Euclid, 4

Euclidean distance

- Euclidean distance, 231–232
- evaluating
 - AutoML Model Builder results, 47–48
 - data clusters, 29
 - linear regression algorithms, 178–180
 - models, 62–64, 155–161
 - bias versus variance, 157–158
 - in binary classification, 115–116
 - in clustering, 124–126
 - confusion matrix, 159–160
 - cross validation, 161
 - linear versus nonlinear, 156
 - in multiclass classification, 119–121
 - noise in dataset, 156
 - regularization, 161
 - in transfer learning, 131
 - underfitting and overfitting, 158–159
- evaluators in ML.NET, 90
- evolution of software intelligence, 24–25
- expected value, variance and, 144
- expert systems, 16–19, 25–26
 - decision trees versus, 185
 - examples of, 18
 - expertise versus intelligence, 25
 - history of, 17
 - internal architecture, 17–18
 - limitations of, 19
 - “Miracle on the Hudson” example, 25–26
 - updating, 26
- experts, definition of, 25
- exposing APIs, 65
- external assessment, 124–126
- extracting dataset features, 55–56

F

- FaceApp, 305
- feature engineering, 54–56
- features
 - definition of, 135
 - extracting, 55–56
 - generating, 54–55
 - selecting, 55
 - type system of, 136
- feed-forward neural networks
 - history of neural networks, 256
 - limitations of, 256–257, 291–292
 - logistic neurons, 260–262

- perceptrons, 257–260
 - activation function, 257–258
 - enabling learning, 260
 - feed-forward layers, 259–260
 - NAND gates, 258–259
- training, 263–270
 - backpropagation algorithm, 264–270
 - gradient descent, 263
 - minibatch gradient, 264

- Feigenbaum, Edward, 17
- five-number summary, 145–146
- forward chaining, 18
- frameworks for neural networks, 282
 - Caffe, 284
 - Keras, 283, 284–287
 - MXNet, 284
 - PyTorch, 283
 - TensorFlow, 282–283
 - Theano, 284
- fully connected layer in CNN (convolutional neural network), 303–304
- functional completeness, 258–259

G

- GAN (generative adversarial neural networks), 304–305
- Gauss, Carl Friedrich, 171
- Gaussian distribution, 224–225
- Gaussian Naïve Bayes, 224–226
- general AI, 19, 27. *See also* autonomous systems
- general recursive functions, 5
- generative adversarial neural networks (GAN), 304–305
- geometric mean, 137–138
- GIL (Global Interpreter Lock), 78
- Godel, Kurt, 4–5, 291
- Google Cloud AutoML, 44
- gradient boosting, 38, 204, 206–210
 - hyperparameters, 209
 - implementations of, 209
 - pros and cons, 210
 - steps in, 206–209
- gradient descent, 37, 174–178
 - in feed-forward neural networks, 263
- graphs, definition of, 182
- grouping. *See* classification; clustering
- gRPC services, API exposure in, 65

H

- harmonic mean, 138
- harmonization of data, 54
- Hawking, Stephen, 3, 20–21
- heatmaps, 55
- heteroscedasticity, 169
- hidden layers in neural networks, 277–280
- Hilbert, David, 4, 211
- histograms, 138–139, 146–147
- history
 - of artificial intelligence, 6–7
 - of expert systems, 17
 - of machine learning, 4–7
 - artificial intelligence, 6–7
 - computing machine formalization, 5
 - Godel's theorems, 4–5
 - human thought formalization, 5–6
 - mechanical reasoning, 4
 - of neural networks, 255
 - feed-forward neural networks, 256
 - limitations of feed-forward neural networks, 256–257
 - McCulloch-Pitts neurons, 255–256
 - of Python, 78
- Hoare, Tony, 97
- holdout, 57, 98, 161
- homogeneity, 186–187
- hosting platforms, selecting, 64–65
- hosting scenarios in ML.NET, 91
- hot data, 51
- human emulation, 341
- human intelligence
 - behavior learning, 14–15
 - memory, 13
 - neurons, 8–13
 - brain architecture, 12–13
 - cerebral cortex structure, 9–10
 - computing power of, 11–12
 - number of, 8–9
 - physiology of, 10–11
- human thought, formalizing, 5–6
- hyperparameter tuning, 63
- hyperplanes, 235–236

I

- ID3 (Iterative Dichotomiser) algorithm, 191–193
- IDataLoader interface, 106

- IDataView interface, 106, 107–108
- IEstimator interface, 106
- image classification, 41
 - in convolutional neural networks, 298
 - data transformations on, 127–129
 - steps in, 127
 - training model, 129–131
- improving data quality, 53
- impurity, information gain and, 187–191
- Inception Model (IM), 127
- incompleteness theory (Godel), 4–5
- incremental learning, 227
- independent events, 213–214
- inferred function in supervised learning, 31
- information gain
 - definition of, 187
 - entropy and, 192
 - impurity and, 187–191
- integrity of data, 70–71
- intelligence. *See also* artificial intelligence (AI)
 - definition of, 7
 - expertise versus, 25
 - human intelligence
 - behavior learning, 14–15
 - memory, 13
 - neurons, 8–13
 - software intelligence
 - adaptability to change, 15–16
 - evolution of, 24–25
 - examples of, 7–8
- interfaces (ML.NET), 105–107
- intermediate format for sentiment analysis, 311–313
- interpretation of probability, 212
- interquartile range, 141
- intersection of events, 213
- irreducible error, 154
- ITransformer interface, 106

J

- JPEG compression, 307

K

- Kaku, Michio, 273
- K-D tree, 233–234
- Keras, 81–82, 283–287
 - environment preparation, 284–285

- models
 - creating, 285–286
 - training, 286
- for sentiment analysis, 317
- kernel functions, 238, 239
- k-fold, 57, 161
- K-Means, 39, 246–247
 - DBSCAN versus, 250–251
 - ideal number of clusters, 247
 - steps in, 246–247
 - training, 123–124
- K-Modes, 247–248
- KNN (K-Nearest Neighbors) algorithm, 230–234
 - Ball Tree, 234
 - brute-force implementation, 233
 - business scenario, 234
 - categorical data in, 232
 - distance calculations, 231–232
 - K-D tree, 233–234
 - number of neighbors, 230–231
 - training, 234

L

- labeled data in supervised learning, 30–31
- Lagrange multipliers, 240
- lambda calculus, 5
- Lao Tzu, 165
- Laplace, Pierre Simon de, 21, 212
- Laplace smoothing, 222
- Laplace's demon, 21
- leaf, definition of, 182
- learning
 - in humans, carrot-and-stick approach, 14–15
 - in neural networks, enabling, 260
- learning by discovery, 33
- learning by example, 33
- learning rate, 209
- leave-p-out technique, 161
- LeCun, Yann, 197
- Lederberg, Joshua, 17
- left tail, 139
- Leibniz, Gottfried, 4
- libraries (Python), 80–82
 - Keras, 81–82
 - Matplotlib, 80–81
 - NumPy, 81
 - Pandas, 80
 - PyTorch, 81–82
 - scikit-learn, 81
 - SciPy, 81
 - TensorFlow, 81–83
 - Theano, 81
- LightGBM, 209
- linear activation function, 274
- linear models, 156
 - fit of, 168–169
- linear regression, 37
 - algorithm for, 169–178
 - cost function identification, 170–171
 - evaluating, 178–180
 - gradient descent, 174–178
 - ordinary least square, 171–174
 - Bayes' theorem and, 227
 - example problem, 165–169
 - in supervised learning, 31
- linear relationships, 167
- linearly separable datasets, 238–239
- logistic neurons, 260–262
- logistic regression, 35, 169
- LSTM (Long Short-Term Memory) neural networks, 295–298
 - bidirectional, 316
 - deep LSTM, 297–298
 - memory context in, 296–297
- Lull, Raymond, 4

M

- machine learning, 19. *See also* artificial intelligence (AI); AutoML (automated machine learning); autonomous systems
 - decision trees in, 183
 - as general AI, 27
 - history of, 4–7
 - artificial intelligence, 6–7
 - computing machine formalization, 5
 - Godel's theorems, 4–5
 - human thought formalization, 5–6
 - mechanical reasoning, 4
 - models. *See* models
 - problems in, 342–343
 - classification, 34–36
 - clustering, 38–40
 - image classification, 41
 - object detection, 41
 - prediction, 36–38
 - text analytics, 42

- statistics in, 24. *See also* statistics
- statistics versus, 151
 - goals of, 152
 - models, 153–155
- steps in, 42
 - data collection, 50–52
 - data preparation, 52–58
 - model deployment, 64–65
 - model evaluation, 62–64
 - model selection, 58–62
- Machine Learning Server (ML Server), 334
- Manhattan distance, 231–232
- Matplotlib, 80–81
- max pooling, 302
- McCarthy, John, 6
- McCulloch, Warren, 255–256
- McCulloch-Pitts neurons, 255–256
- mean (in statistics), 136–138
 - arithmetic mean, 136–137
 - geometric mean, 137–138
 - harmonic mean, 138
- mean squared error (MSE), 145, 157
- mean-shift, 39
- measure of central tendency, 138
- mechanical reasoning, 4
- median (in statistics), 139–141
 - cumulative distribution function (CDF), 139–140
 - properties of, 140–141
 - quartiles, 141
- memory, humans versus computers, 13
- memory context
 - in LSTM neural networks, 296–297
 - in RNN (recurrent neural networks), 293
- Michelangelo Buonarroti, 67
- Microsoft AutoML Model Builder, 44–48
 - evaluation, 47–48
 - price prediction scenario example, 45–46
 - training algorithm selection, 45–47
- minibatch gradient, 177
 - in feed-forward neural networks, 264
- minimizing cost function, 171–174
- MinMax scaler, 162
- “Miracle on the Hudson” example, 25–26
- ML Server (Machine Learning Server), 334
- ML Service (Azure Machine Learning Service), 331–332
- ML Studio (Azure Machine Learning Studio), 329–331
- ML.NET, 83–84, 345
 - architecture, 105
 - catalogs, 109–111
 - data representation, 107–109
 - types and interfaces, 105–107
- classification tasks, 111–121
 - binary classification, 111–116
 - multiclass classification, 116–121
- clustering tasks, 122–126
 - data preparation, 122–123
 - evaluating model, 124–126
 - training model, 123–124
- learning context, 87
 - data preparation, 88
 - evaluators, 90
 - hosting scenarios, 91
 - root object, 87–88
 - trainers, 89
- models
 - creating, 84–85
 - deep learning, 86–87
 - shallow learning, 85–86
- price prediction example
 - client application, 99–103
 - dataset, 93–96
 - testing phase, 97–98
 - training algorithms, 96–97
- transfer learning, 126–131
 - data transformations on, 127–129
 - image classification, steps in, 127
 - training model, 129–131
- MNB (Multinomial Naive Bayes), 220–222
- mode (in statistics), 138–139
- Model Builder, 85–86
- models. *See also* algorithms
 - algorithms versus, 58–59
 - creating in Keras, 285–286
 - data preparation, 162
 - normalization, 163
 - scaling, 162
 - standardization, 163
 - data requirements, 64
 - deploying, 64–65
 - evaluating, 62–64, 155–161
 - bias versus variance, 157–158
 - in binary classification, 115–116
 - in clustering, 124–126
 - confusion matrix, 159–160
 - cross validation, 161
 - linear versus nonlinear, 156
 - in multiclass classification, 119–121

models

- noise in dataset, 156
 - regularization, 161
 - in transfer learning, 131
 - underfitting and overfitting, 158–159
 - ML.NET. *See* ML.NET
 - for price prediction example (ML.NET), getting, 99
 - Python. *See* Python
 - saving for sentiment analysis, 318
 - selecting, 58–62
 - algorithm selection, 59–61
 - neural networks, 61–62
 - statistical versus machine learning, 153–155
 - training, 81
 - in binary classification, 114–115
 - in clustering, 123–124
 - in Keras, 286
 - in multiclass classification, 118–119
 - in sentiment analysis, 313–320
 - in transfer learning, 129–131
 - training time, 63
 - MSE (mean squared error), 145, 157
 - multiclass classification, 34
 - in ML.NET, 116–121
 - confusion matrix, 121
 - data transformations on, 117–118
 - evaluating model, 119–121
 - training model, 118–119
 - with SVM algorithm, 244–245
 - multilabel classification, 34
 - multilinear case for cost function, 171
 - minimizing, 173–174
 - multilinear regression, 37
 - multilinear relationships, 167
 - multimodal datasets, 138
 - multinomial distribution, 221
 - Multnomial Naïve Bayes (MNB), 220–222
 - MXNet, 284
- ## N
- naïve Bayes classifiers, 35
 - algorithm for, 219–220
 - Bernoulli Naïve Bayes (BNB), 223–224
 - components of, 218–219
 - definition of, 219
 - formulation, 217–218
 - Gaussian Naïve Bayes, 224–226
 - Multnomial Naïve Bayes (MNB), 220–222
 - naïve Bayes regression, 226–228
 - NAND gates, 258–259
 - natural language processing (NLP), 42
 - navigating rows in data views, 108–109
 - .NET for Apache Spark, 335–336
 - neural networks, 61–62
 - architecture, 273–274
 - activation functions, 274–277
 - hidden layers, 277–280
 - output layer, 281
 - auto-encoders, 305–307
 - binary classification versus, 287–289
 - convolutional, 298–304
 - convolutional layer, 299–301
 - fully connected layer, 303–304
 - image classification in, 298
 - pooling layer, 301–303
 - feed-forward
 - limitations of, 291–292
 - logistic neurons, 260–262
 - perceptrons, 257–260
 - training, 263–270
 - frameworks, 282
 - Caffe, 284
 - Keras, 283, 284–287
 - MXNet, 284
 - PyTorch, 283
 - TensorFlow, 282–283
 - Theano, 284
 - generative adversarial, 304–305
 - history of, 255
 - feed-forward neural networks, 256
 - limitations of feed-forward neural networks, 256–257
 - McCulloch-Pitts neurons, 255–256
 - Inception Model (IM), 127
 - LSTM, 295–298
 - deep LSTM, 297–298
 - memory context in, 296–297
 - Python libraries, 81–82
 - recurrent, 292–295
 - architecture of, 293–294
 - deep RNN, 295
 - memory context in, 293
 - state management, 294–295
 - for sentiment analysis
 - dropout, 319–320
 - embedding layer, 319
 - epochs of training, 320
 - selecting, 315–318

- shallow versus deep learning, 289
- neurons, 8–13
 - brain architecture, 12–13
 - cerebral cortex structure, 9–10
 - computing power of, 11–12
 - number in neural networks, 280
 - number of, 8–9
 - physiology of, 10–11
 - role in neural networks, 277–280
- Ng, Andrew, 255
- n-grams, 222
- NLP (natural language processing), 42
- noise in dataset, 156
- nonlinear models, 156
- nonlinear regression, 37
- nonlinear relationships, 169
- nonlinearly separable datasets, 237–238
- normal distribution, 224–225
- normalization, 163
- notebooks, 80, 332
- numeric computing, 81
- NumPy, 81

O

- object detection, 41
- observations, definition of, 135
- One vs One, 244
- one-hot encoding, 56
- ONNX format, 323
- ordinary least square algorithm, 171–174
- outlier removal, 54
- outliers, impact of, 192–193
- output formatting for sentiment analysis, 323
- output layer in neural networks, 281
- overfitting, 158–159

P

- Page, Larry, 77
- Pandas, 80
- partial derivatives, definition of, 174
- partitions of events, 214
- PassGAN, 305
- perceptrons, 256, 257–260
 - activation function, 257–258
 - enabling learning, 260
 - feed-forward layers, 259–260
 - NAND gates, 258–259

- A Philosophical Essay on Probabilities* (Laplace), 21
- pipeline, 49
- Pitts, Walter, 255–256
- PoC (proof-of-concept) in data collection, 50
- polynomial regression, 37, 178–179
- pooling layer in CNN (convolutional neural network), 301–303
- precision
 - of human brain, 11
 - of models, 155–161
 - bias versus variance, 157–158
 - confusion matrix, 159–160
 - cross validation, 161
 - linear versus nonlinear, 156
 - noise in dataset, 156
 - regularization, 161
 - underfitting and overfitting, 158–159
- prediction, 36–38, 340–341. *See also* price prediction example (ML.NET pipeline)
 - algorithms used, 37–38
 - with ensemble methods, 198
 - bagging technique, 198–203
 - boosting technique, 203–210
 - as goal of machine learning, 152
 - with linear regression
 - algorithm for, 169–178
 - evaluating algorithm, 178–180
 - example problem, 165–169
 - problems addressed by, 38
 - in sentiment analysis, 322–323
 - in supervised learning, 29–30
 - in SVM algorithm
 - coefficients for, 243
 - mechanics of, 240–241, 244
 - variations of, 37
 - predictive maintenance, 20, 68
- Price, Richard, 212
- price prediction example (ML.NET pipeline)
 - client application, 99–103
 - designing user interface, 102–103
 - getting model file, 99
 - making predictions, 100–102
 - questioning data and problem approach, 103
 - setting up, 99–100
 - dataset, 93–96
 - conversion to C# class, 94
 - data transformations on, 94–95
 - limitations of, 95–96
 - testing phase, 97–98
 - training algorithms, 96–97

price prediction scenario (AutoML Model Builder)

- price prediction scenario (AutoML Model Builder), 45–46
- primordial intelligence, 16–17
- probability
 - Bayesian statistics, 211–216. *See also* naïve Bayes classifiers
 - Bayes' theorem, 214–215
 - chain rule, 213
 - classification and, 216–218
 - conditional probability, 213
 - independent events, 213–214
 - intersection of events, 213
 - partitions of events, 214
 - sample scenario, 215–216
 - density function, 225–226
 - interpretation of, 212
 - zero probability problem, 221–222
- problems in machine learning, 342–343
 - classification, 34–36
 - algorithms used, 34–36
 - problems addressed by, 36
 - variations of, 34
 - clustering, 38–40
 - algorithms used, 39–40
 - problems addressed by, 40
 - formalizing for sentiment analysis, 310
 - image classification, 41
 - object detection, 41
 - prediction, 36–38
 - algorithms used, 37–38
 - problems addressed by, 38
 - variations of, 37
 - text analytics, 42
- proof-of-concept (PoC) in data collection, 50
- pruning, 195
- Python, 78, 332
 - end-to-end solutions and, 82–83
 - environment preparation for Keras, 284–285
 - history of, 78
 - as language for scientists, 79
 - libraries in, 80–82
 - Keras, 81–82
 - Matplotlib, 80–81
 - NumPy, 81
 - Pandas, 80
 - PyTorch, 81–82
 - scikit-learn, 81
 - SciPy, 81
 - TensorFlow, 81–83
 - Theano, 81
 - simplicity of, 79
 - PyTorch, 81–82, 283

Q

- quality of data, 67–70
 - data collection and, 69–70
 - data validity, 68
 - improving, 53
- quartiles, 141
- Quicksort algorithm, 97

R

- random forests, 35, 198–203
 - pros and cons, 202–203
 - steps in, 200–202
- range (in statistics), 142
- range normalization, 54
- recurrent neural networks (RNN). *See* RNN (recurrent neural networks)
- reducing dataset size, 245–246
- regression, 36–38. *See also* linear regression; nonlinear regression
 - algorithms used, 37–38
 - classification versus, 165
 - with decision trees, 182, 194–195
 - naïve Bayes, 226–228
 - problems addressed by, 38
 - in supervised learning, 29–30
 - variations of, 37
- regression decision tree, 37
- regularization, 161, 179–180
- relationships
 - complex, 154–155
 - hypotheses about, 167–169
 - simple, 153
- ReLU activation function, 276–277, 302
- ridge regression, 179
- right tail, 139
- rigidity in artificial intelligence (AI), 347–348
- RNN (recurrent neural networks), 292–295
 - architecture of, 293–294
 - deep RNN, 295

- memory context in, 293
- state management, 294–295
- robust scaler, 162
- root object in ML.NET, 87–88
- Rosenblatt, Frank, 256
- row navigation in data views, 108–109
- Russell, Bertrand, 4
- Rutherford, Ernest, 135

S

- same padding (SP), 301
- saving models for sentiment analysis, 318
- scalar product of vectors, 239–240
- scale of plots (in statistics), 149–150
- scaling, 162
- scatter plot matrices, 148–149
- scatter plots, 148
- schema information in ML.NET, 106
- schema of data views, 108
- scikit-learn, 81
- SciPy, 81
- Searle, John, 6
- selecting
 - dataset features, 55
 - ecosystems for sentiment analysis, 314
 - hosting platforms, 64–65
 - models, 58–62
 - algorithm selection, 59–61
 - neural networks, 61–62
 - trainers for sentiment analysis, 315–318
 - training algorithms, 45–47, 59–61, 96–97
- semi-hot data, 51
- sentiment, artificial intelligence (AI) and, 20–21
- sentiment analysis, 309–310
 - client application, 321
 - data collection, 321–322
 - output formatting, 323
 - prediction from, 322–323
 - data preparation, 310–313
 - data collection, 311
 - intermediate format for data transformations, 311–313
 - problem formalization, 310
 - in ML.NET, 111–112
 - training model, 313–320
 - dictionary of words construction, 314–315
 - dropout, 319–320
 - ecosystem selection, 314
 - embedding layer, 319
 - epochs of training, 320
 - trainer selection, 315–318
- shallow learning, 40
 - deep learning versus, 289
 - in ML.NET, 85–86
- shuffling data in data views, 109
- sigmoid activation function, 261–262, 274–275
- silhouette method, 247
- simple relationships, 153
- simplicity of Python, 79
- softmax activation function, 275
- software developers, data scientists versus, 73–74, 344–346
- software intelligence. *See also* artificial intelligence (AI)
 - adaptability to change, 15–16
 - evolution of, 24–25
 - examples of, 7–8
- SP (same padding), 301
- sparse data, grouping, 55
- splitting test and training datasets, 57, 98
- SSMLS (SQL Server Machine Learning Services), 333
- standard deviation, variance and, 142–144
- standardization, 163
- stateful neural networks. *See* RNN (recurrent neural networks)
- statistics
 - Bayesian statistics, 211–216. *See also* naïve Bayes classifiers
 - Bayes' theorem, 214–215
 - chain rule, 213
 - classification and, 216–218
 - conditional probability, 213
 - independent events, 213–214
 - intersection of events, 213
 - partitions of events, 214
 - sample scenario, 215–216
 - bias, 144–145
 - data representation, 145–150
 - five-number summary, 145–146
 - histograms, 146–147
 - scale of plots, 149–150
 - scatter plot matrices, 148–149
 - scatter plots, 148
 - data sampling, 69–70
 - in machine learning, 24
 - machine learning versus, 151
 - goals of, 152
 - models, 153–155

statistics

- mean, 136–138
 - arithmetic mean, 136–137
 - geometric mean, 137–138
 - harmonic mean, 138
 - median, 139–141
 - cumulative distribution function (CDF), 139–140
 - properties of, 140–141
 - quartiles, 141
 - mode, 138–139
 - variance, 142–144
 - expected value and, 144
 - standard deviation and, 142–144
 - stochastic dual coordinate ascent, 37
 - stochastic gradient, 177
 - stop-words, 222
 - store-and-train model, 53
 - strong AI, 21
 - Sun Tzu, 343
 - supervised learning, 20, 29–31
 - inferred function, 31
 - labeled data, 30–31
 - as learning by example, 33
 - prediction and classification, 29–30
 - support vectors, 236–237
 - SVM (Support Vector Machine) algorithm, 35, 235–245
 - coefficients for prediction, 243
 - hyperplanes, 235–236
 - Lagrange multipliers, 240
 - linearly separable datasets, 238–239
 - mechanics of prediction, 240–241, 244
 - multiclass classification with, 244–245
 - nonlinearly separable datasets, 237–238
 - scalar product of vectors, 239–240
 - support vectors, 236–237
 - training, 242–243
 - vector operations, 239
 - synapses, 10–11
- ## T
- TanH activation function, 275–276
 - tasks (ML.NET) for training, 89
 - task-specific catalogs (ML.NET), 110
 - TensorFlow, 81–83, 282–283
 - transfer learning with, 126–131
 - data transformations on, 127–129
 - image classification, steps in, 127
 - training model, 129–131
 - test datasets, splitting from training datasets, 57, 98
 - testing phase, price prediction example (ML.NET pipeline), 97–98
 - text analytics, 42
 - text-based features, 136
 - Theano, 81, 284
 - thinking machines, 5–7. *See also* artificial intelligence (AI)
 - threads, 78
 - time-based data
 - collecting, 69
 - in neural networks, 61–62
 - timeline series, 30
 - timeliness of data, 71
 - timestamp features, 136
 - trainers
 - in ML.NET, 89
 - selecting for sentiment analysis, 315–318
 - training
 - auto-encoders, 306
 - feed-forward neural networks, 263–270
 - backpropagation algorithm, 264–270
 - gradient descent, 263
 - minibatch gradient, 264
 - KNN algorithm, 234
 - models
 - in binary classification, 114–115
 - in clustering, 123–124
 - in Keras, 286
 - in multiclass classification, 118–119
 - in sentiment analysis, 313–320
 - in transfer learning, 129–131
 - SVM algorithm, 242–243
 - training algorithms
 - backpropagation, 264–270
 - price prediction example (ML.NET pipeline), 96–97
 - selecting, 45–47, 59–61, 96–97
 - training datasets, 152
 - finalizing, 56–58
 - splitting from test datasets, 57, 98
 - transfer learning, 126–131
 - data transformations on, 127–129
 - image classification, steps in, 127
 - training model, 129–131
 - transformations. *See* data transformations
 - trees. *See* decision trees

trimodal datasets, 138
 Turing, Alan, 5, 6, 23, 229
 Turing machine, 5
 Turing test, 6
 type system of features, 136
 types (ML.NET), 105–107

U

underfitting, 158–159
 underflow, 222
 uniform representation of data, 54
 uniqueness of data, 70–71
 unsupervised learning, 20, 27–29
 with clustering, 245
 business scenario, 245–246
 DBSCAN, 248–251
 K-Means, 246–247
 K-Modes, 247–248
 discovering data clusters, 27–28
 evaluating data clusters, 29
 as learning by discovery, 33
 updating expert systems, 26
 user interfaces, designing, 102–103

V

valid padding, 301
 validity of data, 68
 value of algorithms, measuring, 97

variance, 142–144
 bias versus, 157–158
 expected value and, 144
 standard deviation and, 142–144
 variance threshold, 55
 vector of errors in backpropagation algorithm, 269–270
 vectors
 basic operations on, 239
 scalar product of, 239–240
 von Neumann, John, 6, 11, 21, 151

W

waterfall methodology, 346–347
 weak AI, 21
 weak learners, 197
 weather forecasting, 212
 web applications, API exposure in, 65
 Wirth, Niklaus, 181
 workflows, 42

X

XGBoost, 209

Z

Zen of Python, 79
 zero probability problem, 221–222