

Introducing Machine Learning

Sample pages



Professional

Dino Esposito
Francesco Esposito

Contents

<i>Acknowledgments</i>	<i>xix</i>
<i>About the Authors</i>	<i>xxi</i>
<i>Introduction</i>	<i>xxiii</i>

PART I **LAYING THE GROUNDWORK OF MACHINE LEARNING**

Chapter 1	How Humans Learn	3
	The Journey Toward Thinking Machines	4
	The Dawn of Mechanical Reasoning	4
	Godel’s Incompleteness Theorems	4
	Formalization of Computing Machines	5
	Toward the Formalization of Human Thought	5
	The Birth of Artificial Intelligence as a Discipline	6
	The Biology of Learning	7
	What Is Intelligent Software, Anyway?	7
	How Neurons Work	8
	The Carrot-and-Stick Approach	14
	Adaptability to Changes	15
	Artificial Forms of Intelligence	16
	Primordial Intelligence	16
	Expert Systems	16
	Autonomous Systems	19
	Artificial Forms of Sentiment	20
	Summary	22
Chapter 2	Intelligent Software	23
	Applied Artificial Intelligence	23
	Evolution of Software Intelligence	24
	Expert Systems	25

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
--------------------	-----

Sample pages

Mapping Problems and Algorithms

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

—Muhammad 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.)

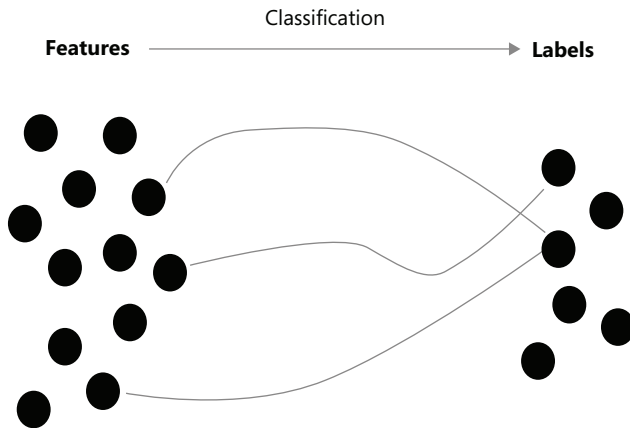


FIGURE 3-1 A graphical representation of a classification function

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.)

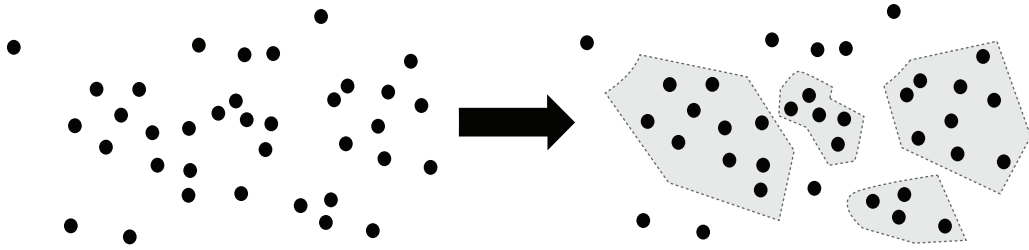


FIGURE 3-2 The final outcome of a clustering algorithm run on a dataset

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.

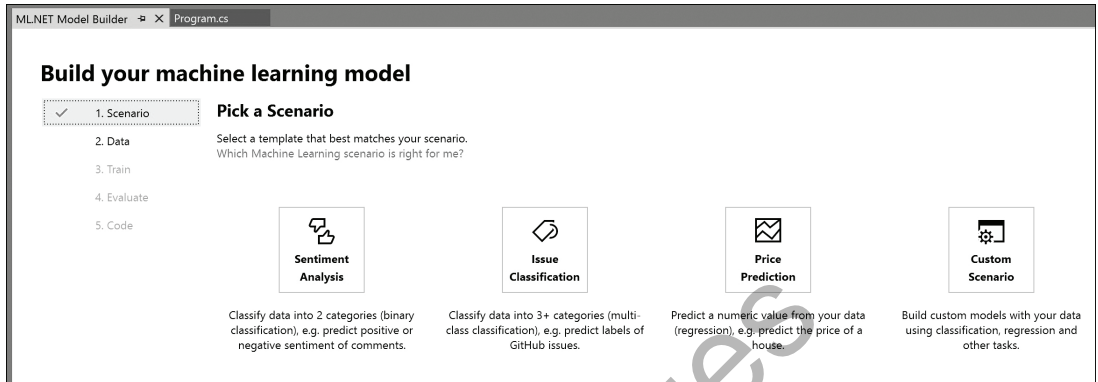


FIGURE 3-3 The main page of the Model Builder Visual Studio extension

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.

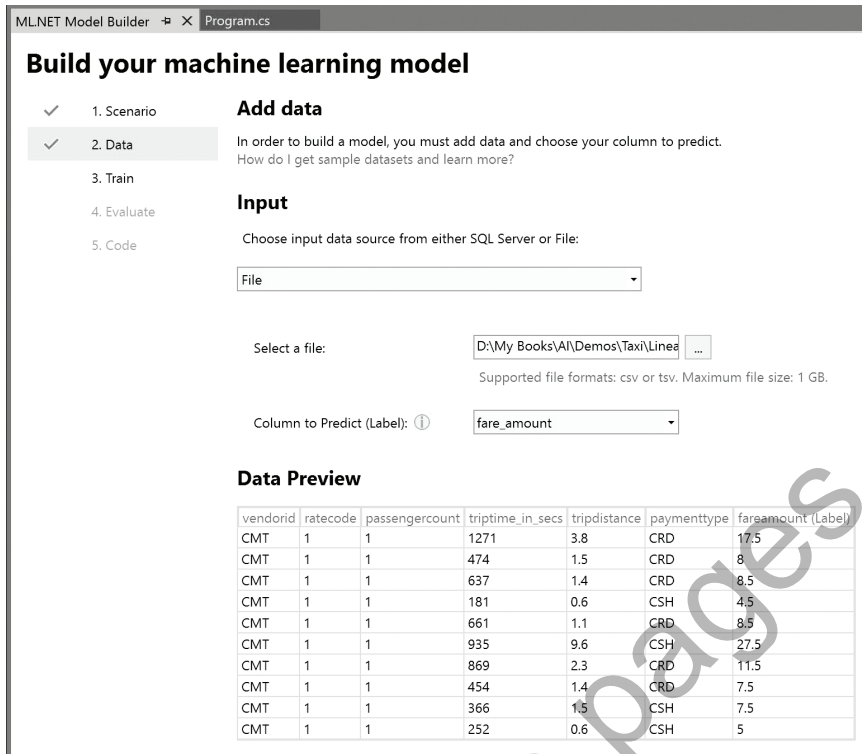


FIGURE 3-4 Loading data into the model

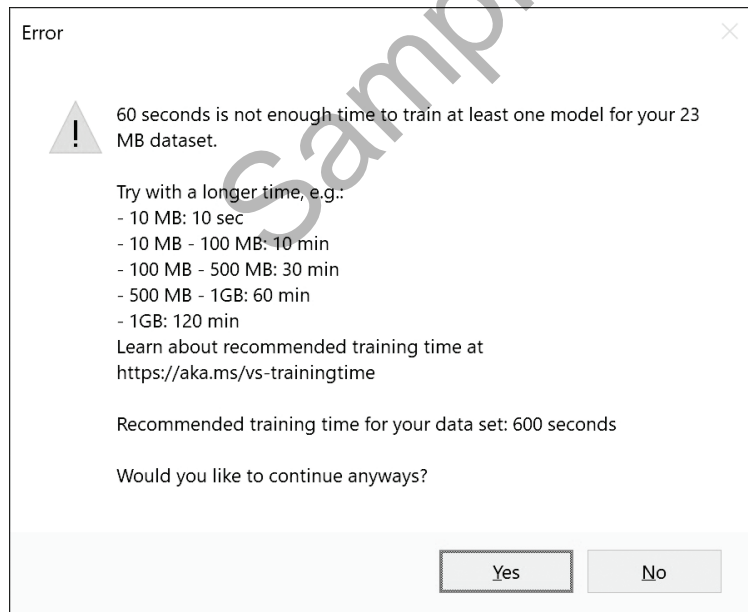


FIGURE 3-5 Estimating the training time

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

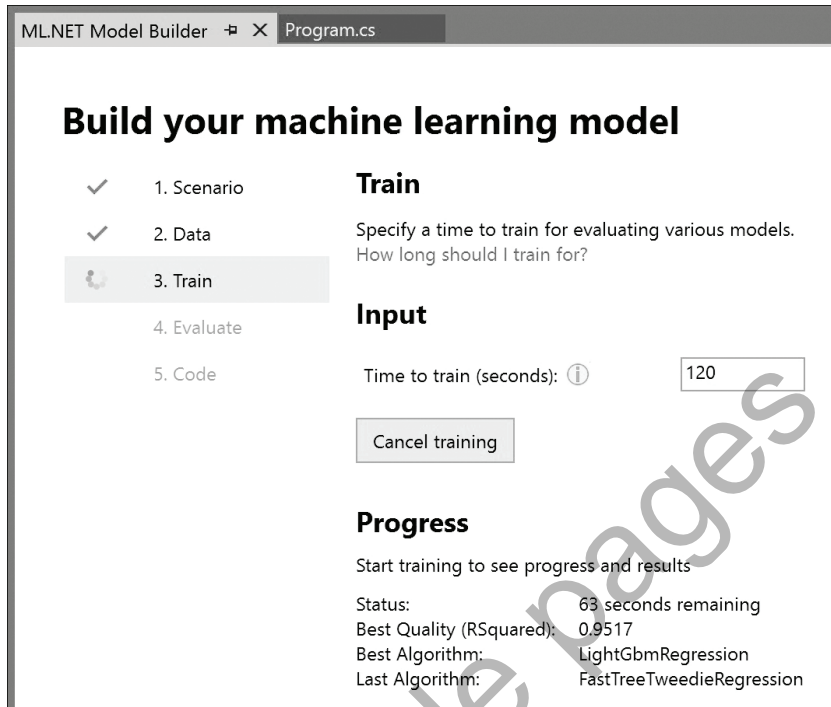


FIGURE 3-6 AutoML tries different algorithms and uses some metrics to evaluate the quality.

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.)

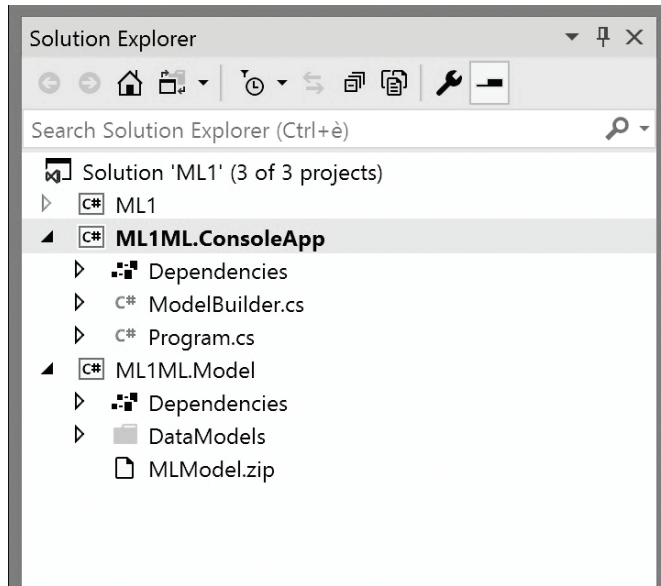


FIGURE 3-7 Autogenerated projects added by the Model Builder

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.