

Advanced Data Analysis

DATA 71200

Class 1

Course Description

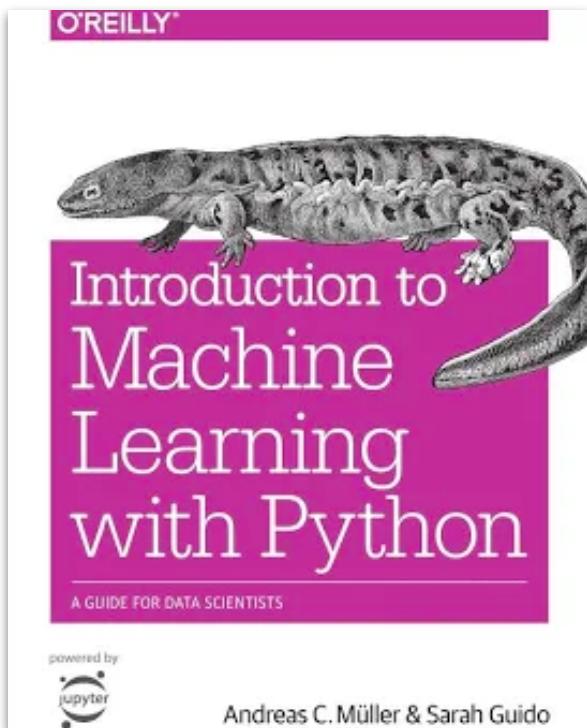
- ▶ This course will provide you with skills necessary to **apply machine learning techniques to data, and interpret and communicate their results.**
- ▶ You will also begin to develop **intuitions** about when machine learning is an appropriate tool versus other statistical methods.
- ▶ This course will cover both **supervised methods** (e.g., k-nearest neighbors, naïve Bayes classifiers, decision trees, and support vector machines) and **unsupervised methods** (e.g., principal component analysis and k-means clustering).
 - The supervised methods will focus primarily on “**classic” machine learning techniques** where features are designed rather than learned, although we will briefly look at recent deep learning models with neural networks.
- ▶ This is an **applied machine learning class** that emphasizes the intuitions and know-how needed to get learning algorithms to work in practice, rather than mathematical derivations.
- ▶ The course will be taught in **Python**, primarily using the **scikit-learn** library.

Course Objectives

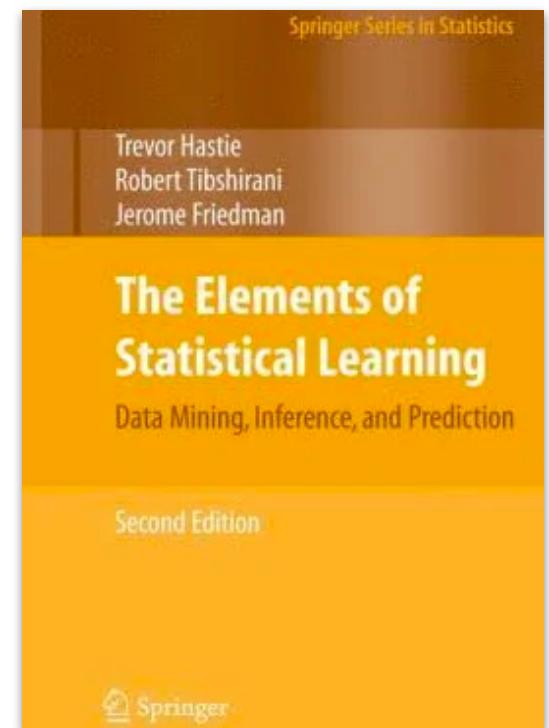
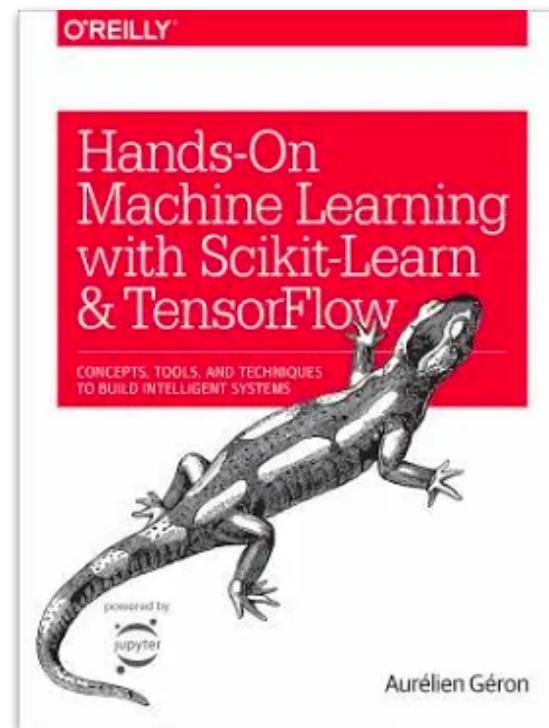
- ▶ By the end of the course, you will be able to
 - articulate the main assumptions underlying machine learning approaches
 - demonstrate the basic principles of dataset creation
 - articulate the importance of data representations
 - evaluate machine learning algorithms
 - articulate the difference between supervised and unsupervised learning
 - apply a range of supervised and unsupervised learning techniques

Textbooks

Required



Recommended



Grade Breakdown

Class Participation 10%

Datacamp Assignments 25%

Project 1: Dataset creation 15%

Project 2: Supervised learning 15%

Project 3: Unsupervised learning 15%

Final Paper 20%

Grade Breakdown Details

- ▶ **Class Participation: 10%**
 - The participation grade is a combination of attendance (including arriving on time); attentiveness, engagement, and participation during class; and general preparedness for class discussions.
- ▶ **Datacamp Assignments: 25%**
 - These projects are hands-on activities designed to both provide coding background and reinforce the concepts covered in class.

Grade Breakdown Details

- ▶ **Project 1 (Dataset creation): 15%**
 - Curation and cleaning of a labeled data set that you will use for the supervised and unsupervised learning tasks in project 2 and 3. The dataset can built from existing data and should be stored in your GitHub repositiory.
- ▶ **Project 2 (Supervised learning): 15%**
 - Application of two supervised learning techniques on the dataset you created in Project 1. This assignment should be completed as a Jupyter notebook your GitHub repository.

Grade Breakdown Details

- ▶ **Project 3 (Unsupervised learning): 15%**
 - Application of two unsupervised learning techniques on the dataset you created in Project 1. This assignment should be completed as a Jupyter notebook your GitHub repository.
- ▶ **Final Paper: 20%**
 - A 5--8 page paper describing the work you did in projects 1--3 (your dataset and your supervised and unsupervised experiments). The paper should describe both what you did technically and what you learned from the relative performance of the machine learning approaches you applied to your dataset. This assignment should be posted as a PDF in your GitHub repository.

Course Schedule

1-Jun Introduction/What is Machine Learning?

2-Jun Getting Started with Machine Learning

3-Jun Async

7-Jun Inspecting Data

8-Jun Representing Data

9-Jun Evaluation Methods

Course Schedule

10-Jun	Async
14-Jun	Async <i>Project 1 Due</i>
15-Jun	Supervised Learning (k-Nearest Neighbors, Linear Models, and Naive Bayes Classifiers)
16-Jun	Async
17-Jun	Supervised Learning (Decision Trees, Support Vector Machines and Uncertainty estimates from Classifiers)

Course Schedule

21-Jun	Async <i>Project 2 Due</i>
22-Jun	Unsupervised Learning (Dimensionality Reduction & Feature Extraction, and Manifold Learning)
23-Jun	Async
24-Jun	Unsupervised Learning (Clustering)
28-Jun	<i>Project 3 Due</i> <i>DataCamp Assignments Due</i>
1-Jul	<i>Final Project Due</i>

Coding Environment

▶ Python 3

- matplotlib, NumPy, Pandas, SciPy, scikit learn

▶ Jupyter notebooks

Anaconda Distribution

The World's Most Popular Python/R Data Science Platform [Download](#)

The open-source **Anaconda Distribution** is the easiest way to perform Python/R data science and machine learning on Linux, Windows, and Mac OS X. With over 19 million users worldwide, it is the industry standard for developing, testing, and training on a single machine, enabling *individual data scientists* to:

- Quickly download 7,500+ Python/R data science packages
- Manage libraries, dependencies, and environments with **Conda**
- Develop and train machine learning and deep learning models with **scikit-learn**, **TensorFlow**, and **Theano**
- Analyze data with scalability and performance with **Dask**, **NumPy**, **pandas**, and **Numba**
- Visualize results with **Matplotlib**, **Bokeh**, **Datashader**, and **Holoviews**



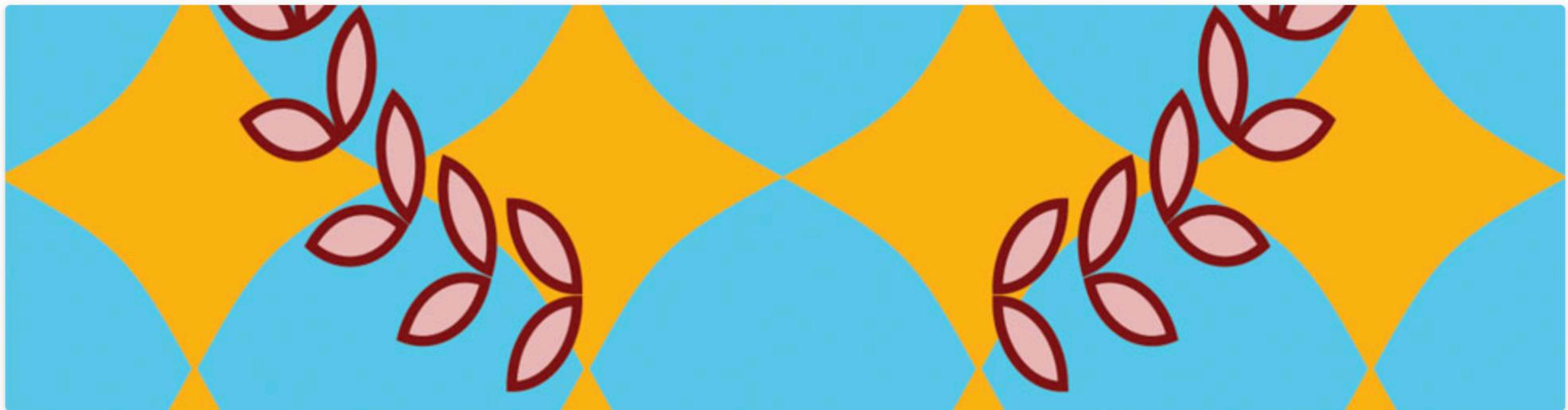
Tutorial: <https://machinelearningmastery.com/setup-python-environment-machine-learning-deep-learning-anaconda/>

Class Website

DATA 71200 Advanced Data Analysis Methods

An introduction to supervised and unsupervised machine learning methods

[HOME](#) [SYLLABUS](#) [COURSE SCHEDULE](#) [RESOURCES](#) [POSTS](#)



<https://data71200su21.commons.gc.cuny.edu/>

Data Camp



Search

Learn ▾

Assessment

Pricing

For Business

Sign in

THE SMARTEST WAY TO

Learn Data Science Online

The skills people and businesses need to succeed are changing. No matter where you are in your career or what field you work in, you will need to understand the language of data. With DataCamp, you learn data science today and apply it tomorrow.

Start Learning For Free



git Shell SPREADSHEETS

Create Your Free Account

LinkedIn

Facebook

Google

or



Email address



Password

Create Free Account

By continuing you accept the Terms of Use and Privacy Policy. You also accept that you are aware that your data will be stored outside of the EU and that you are above the age of 16.

Question Set 1

Géron (p. 4–9)

- ▶ **How would you define Machine Learning?**
- ▶ **Can you name four types of problems where it shines?**
- ▶ **What is a labeled training set?**

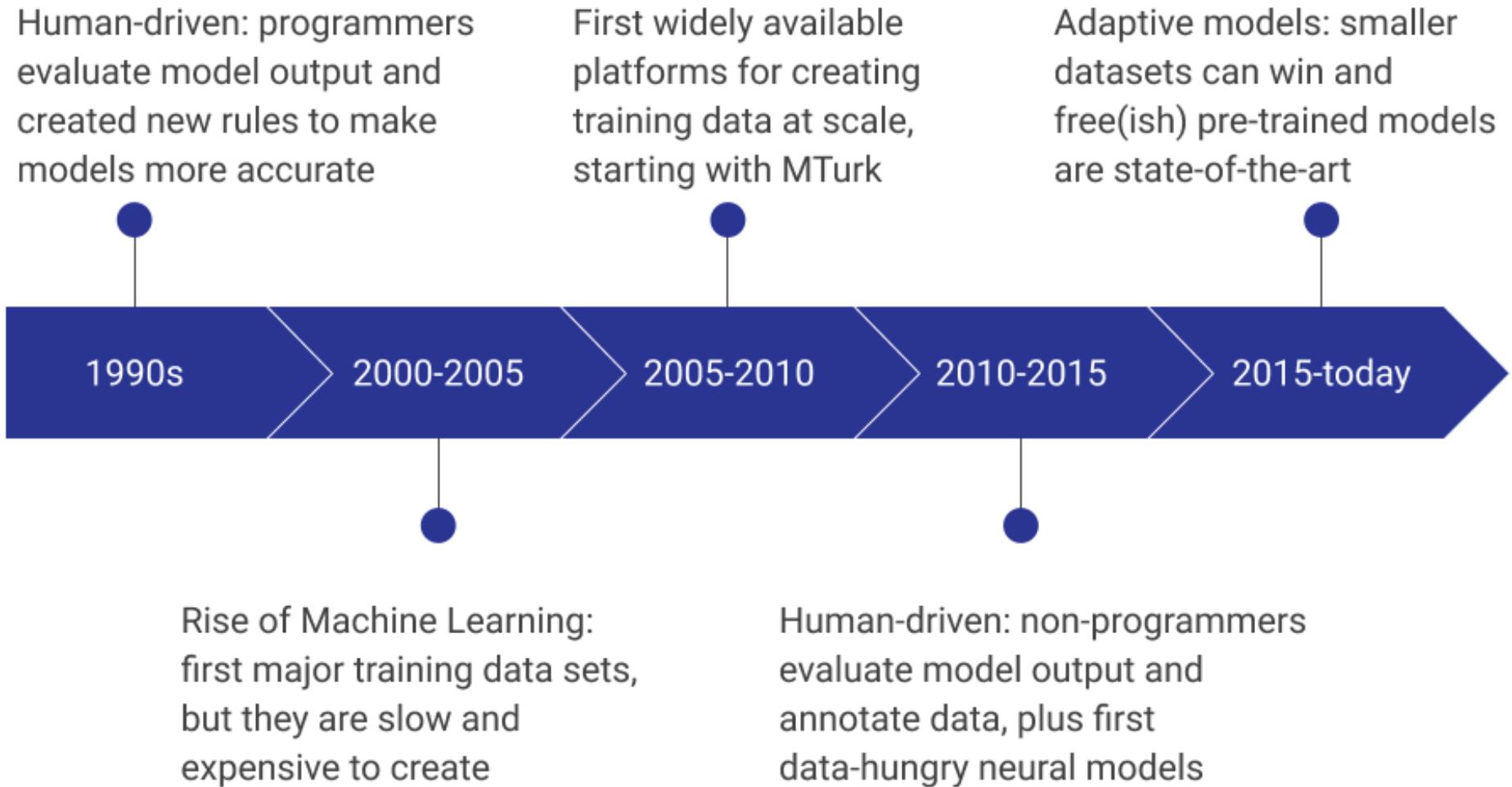
Question Set 1

- ▶ **How would you define Machine Learning?**
 - “Machine Learning is about building systems that can learn from data. Learning means getting better at some task, given some performance measure.”
- ▶ **Can you name four types of problems where it shines?**
 - “Machine Learning is great for complex problems for which we have no algorithmic solution, to replace long lists of hand-tuned rules, to build systems that adapt to fluctuating environments, and finally to help humans learn (e.g., data mining).”

Question Set 1

- ▶ **What is a labeled training set?**
 - “A labeled training set is a training set that contains the desired solution (a.k.a. a label) for each instance.”

Machine Learning



Key Questions

- ▶ “How can one construct computer systems that automatically improve through experience?”
- ▶ “What are the fundamental statistical-computational-information-theoretic laws that govern all learning systems, including computers, humans, and organizations?”
- ▶ “How accurately can the algorithm learn from a particular type and volume of training data?”
- ▶ “How robust is the algorithm to errors in its modeling assumptions or to errors in the training data”

Machine Learning vs Traditional Programming

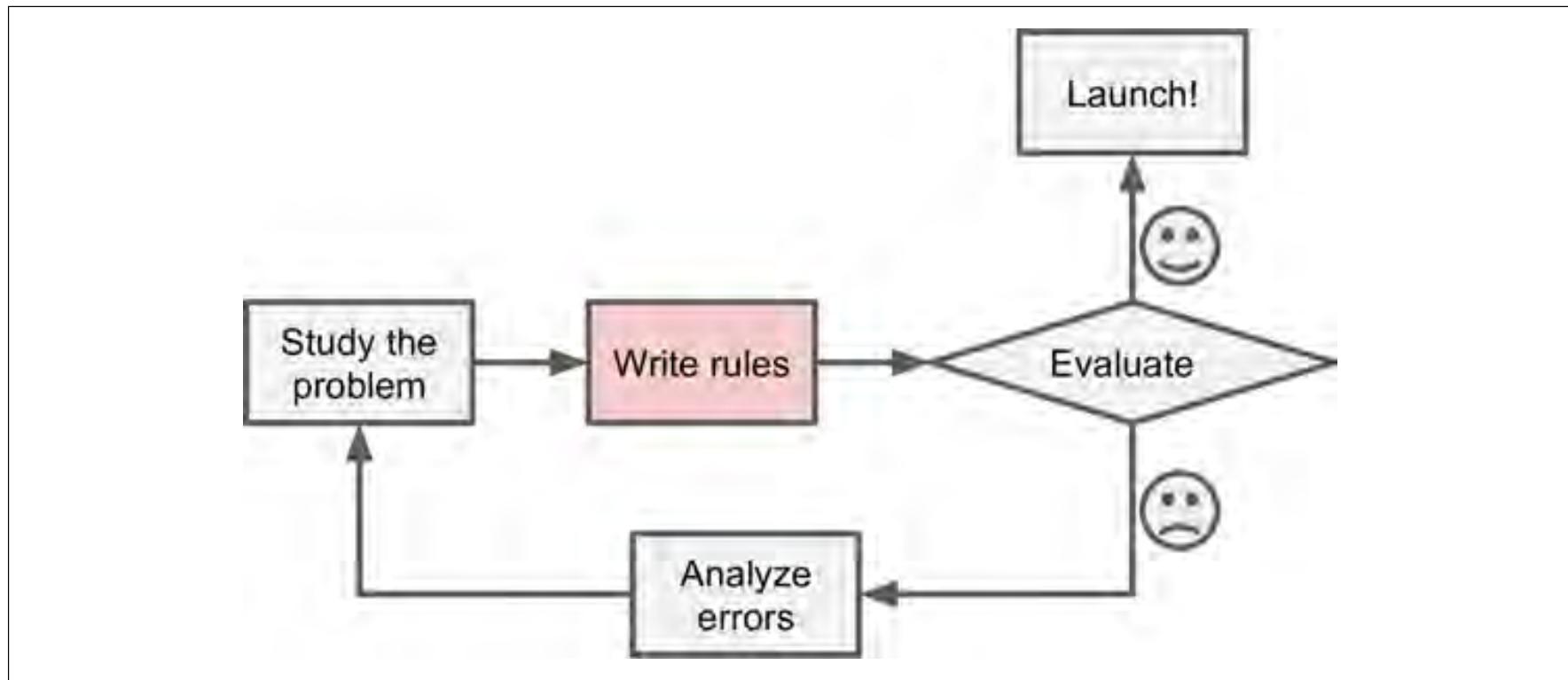


Figure 1-1. The traditional approach

Machine Learning vs Traditional Programming

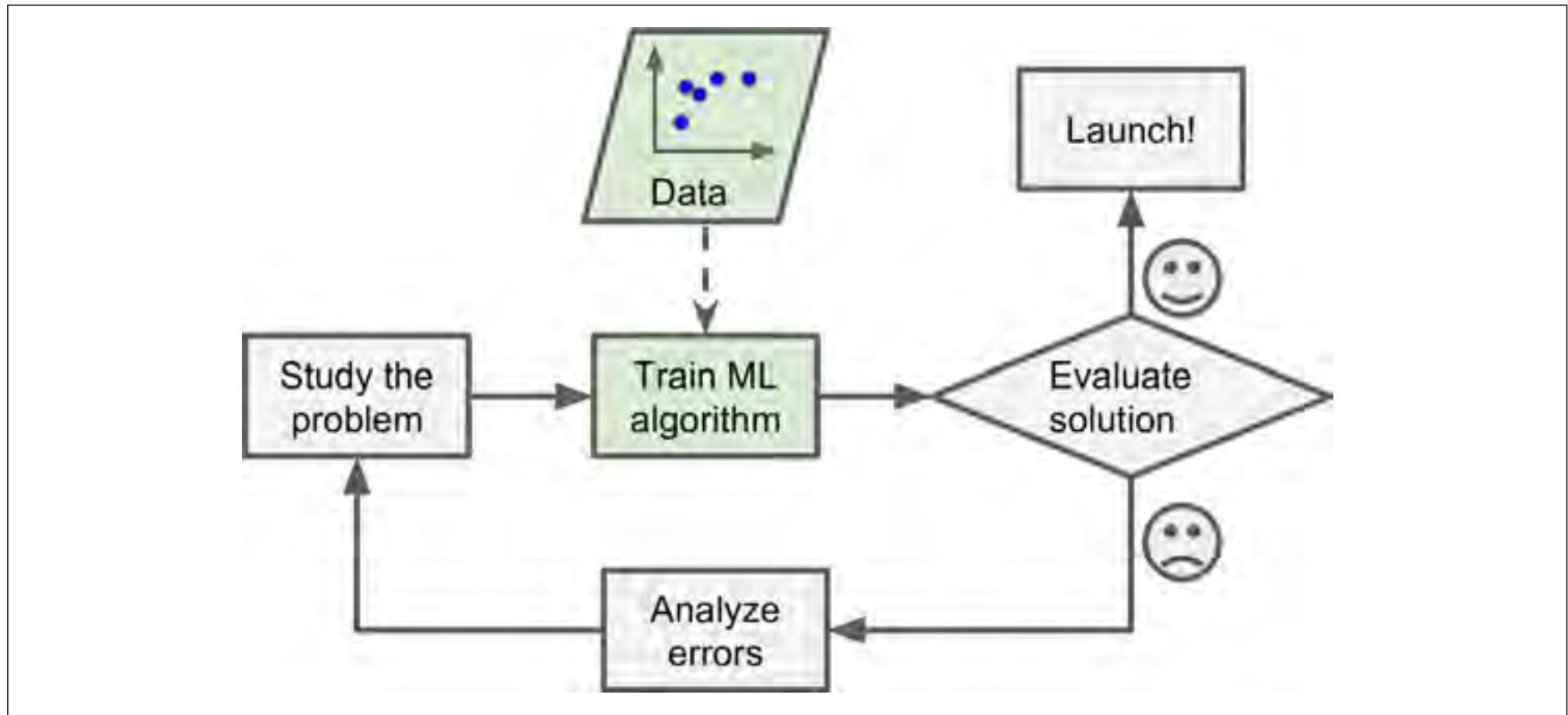


Figure 1-2. Machine Learning approach

Challenges

- ▶ “huge data sets require computationally tractable algorithms”
- ▶ “highly personal data raise the need for algorithms that minimize privacy effects”
- ▶ “the availability of huge quantities of unlabeled data raises the challenge of designing learning algorithms to take advantage of it”

Supervised Learning

► Function approximation problem

- “the training data take the form of a collection of (x, y) pairs and the goal is to produce a prediction y^* in response to a query x^* ”
- Task is to learn a mapping, $f(x)$, which outputs a y value for each inputted x value

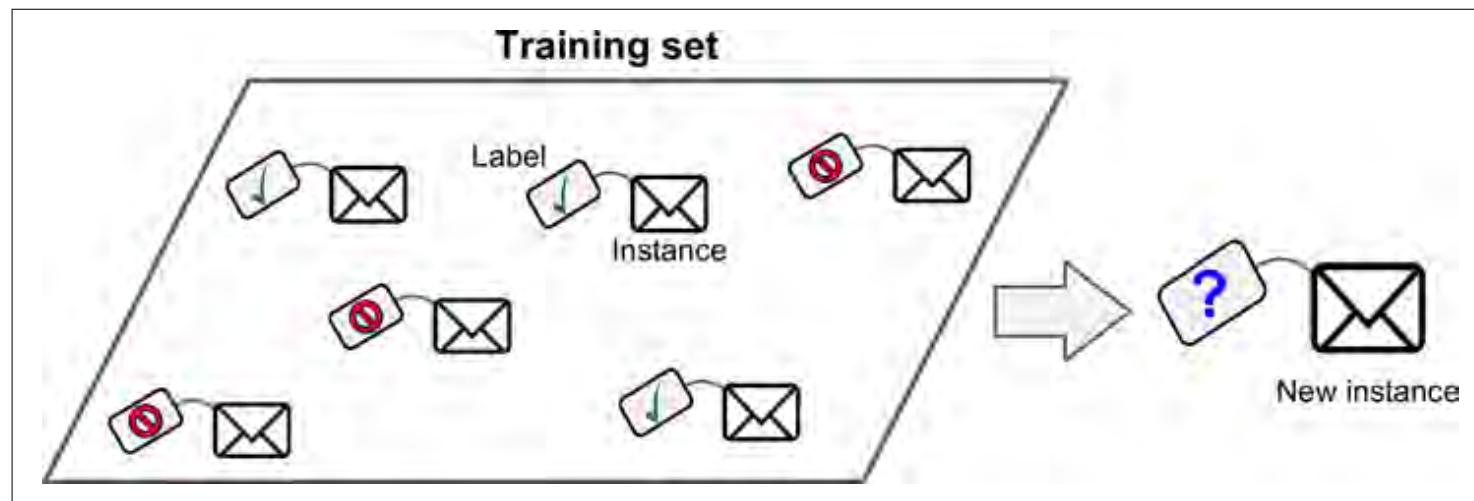


Figure 1-5. A labeled training set for supervised learning (e.g., spam classification)

Jordan, Michael I. and Tom M. Mitchell. (2015). “Machine Learning: Trends, perspectives, and prospects” *Science*.

Image from: Géron, Aurélien. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow* O'Reilly Media, Inc.

Supervised Learning

- ▶ **k -Nearest Neighbors**
- ▶ **Linear Regression**
- ▶ **Logistic Regression**
- ▶ **Support Vector Machines (SVMs)**
- ▶ **Decision Trees and Random Forests**
- ▶ **Naive Bayes Classifiers**
- ▶ **Neural networks**

Supervised Learning

- ▶ “**diversity of learning architectures and algorithms reflects the diverse needs of applications**”
 - “with different architectures capturing different kinds of mathematical structures, offering different levels of amenability to post-hoc visualization and explanation, and providing varying trade-offs between computational complexity, the amount of data, and performance.”

Supervised Learning



Figure 1-6. Regression

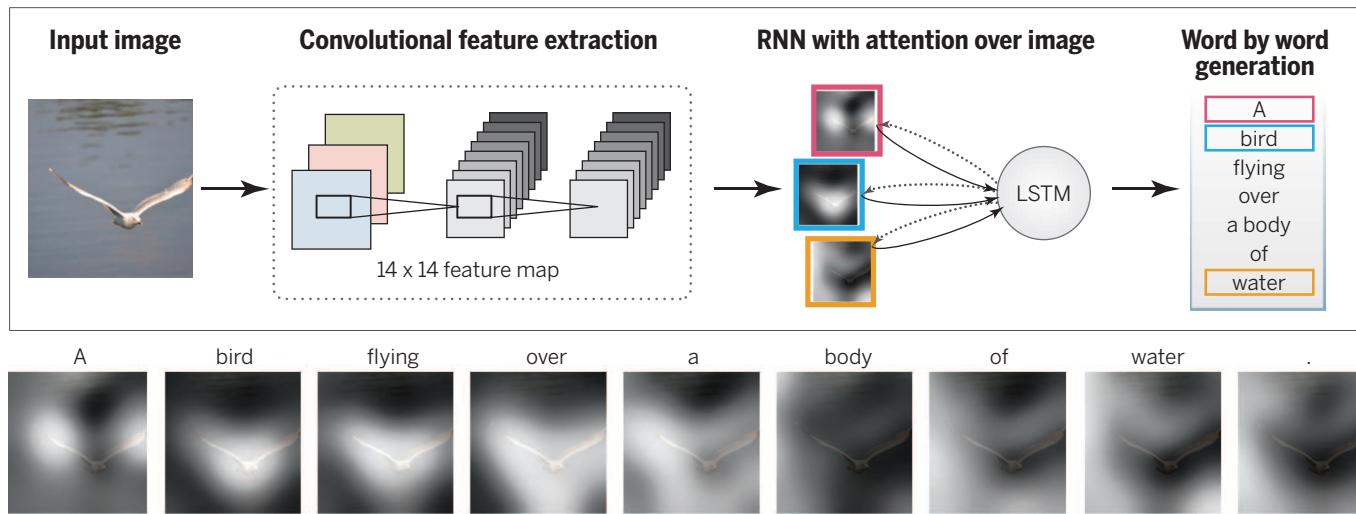


Fig. 2. Automatic generation of text captions for images with deep networks. A convolutional neural network is trained to interpret images, and its output is then used by a recurrent neural network trained to generate a text caption (top). The sequence at the bottom shows the word-by-word focus of the network on different parts of input image while it generates the caption word-by-word. [Adapted with permission from (30)]

Top image from: Géron, Aurélien. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow* O'Reilly Media, Inc.

Bottom image from: Jordan, Michael I. and Tom M. Mitchell. (2015). "Machine Learning: Trends, perspectives, and prospects" *Science*.

Unsupervised Learning

- ▶ “the analysis of unlabeled data under assumptions about structural properties of the data (e.g., algebraic, combinatorial, or probabilistic)”

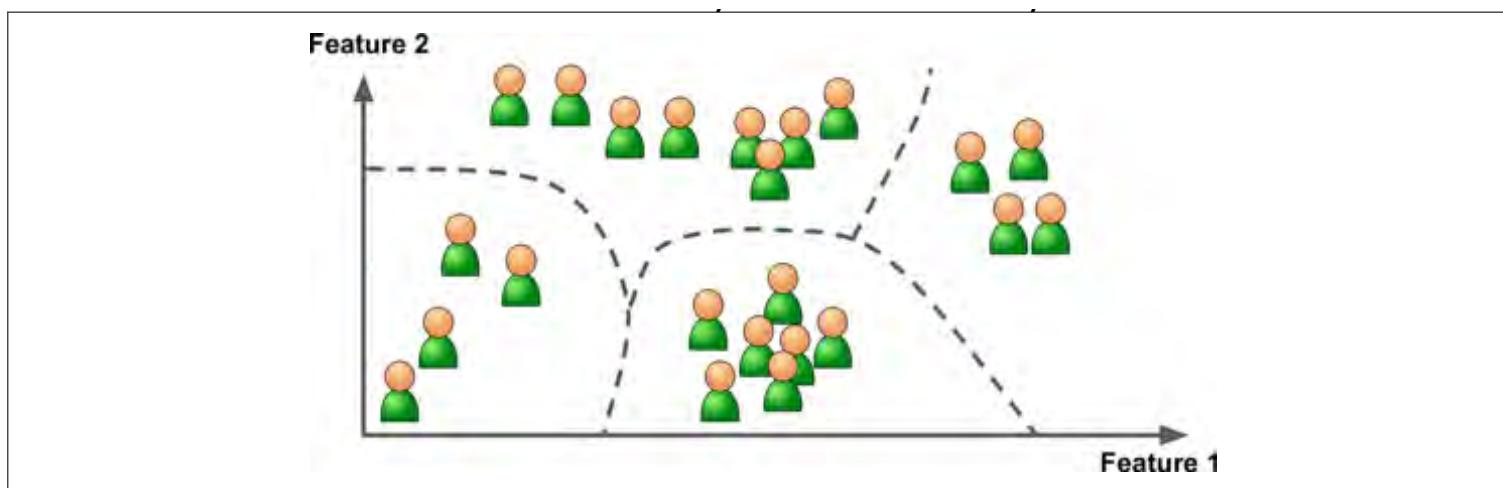


Figure 1-8. Clustering

Unsupervised Learning

- The models make the assumption “that data lie on a low-dimensional manifold and aim to identify that manifold explicitly from the data”
 - Dimensionality reduction (e.g., PCA)
 - Clustering (e.g., k -means)

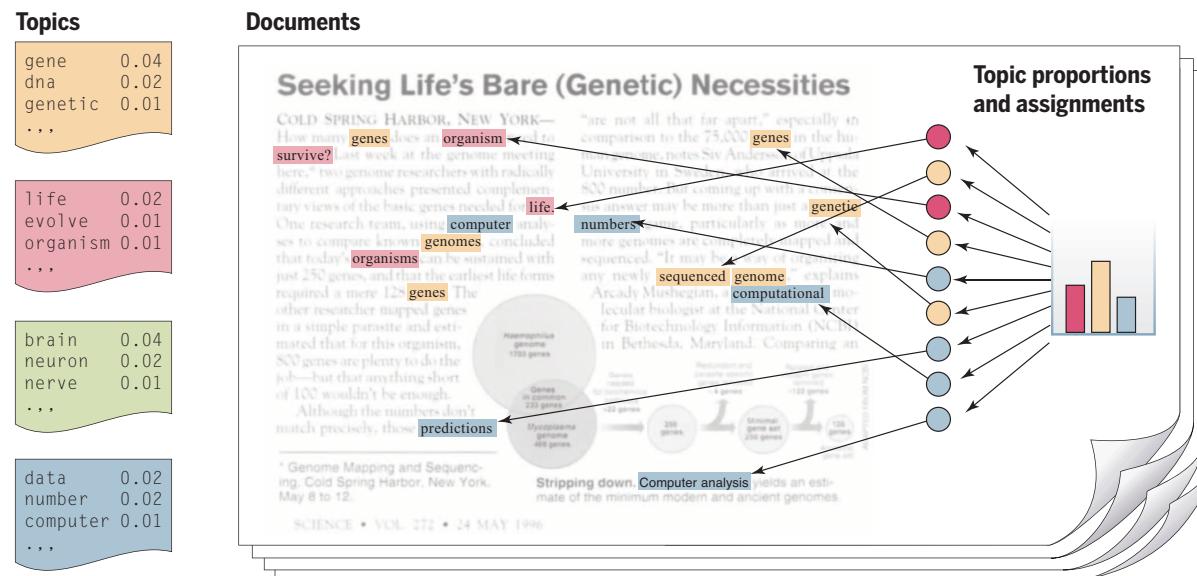


Fig. 3. Topic models. Topic modeling is a methodology for analyzing documents, where a document is viewed as a collection of words, and the words in the document are viewed as being generated by an underlying set of topics (denoted by the colors in the figure). Topics are probability distributions across words (leftmost column), and each document is characterized by a probability distribution across topics (histogram). These distributions are inferred based on the analysis of a collection of documents and can be viewed to classify, index, and summarize the content of documents. [From (31). Copyright 2012, Association for Computing Machinery, Inc. Reprinted with permission]

Semi-supervised Learning

- ▶ “makes use of unlabeled data to augment labeled data in a supervised learning context, and discriminative training blends architectures developed for unsupervised learning with optimization formulations that make use of labels”

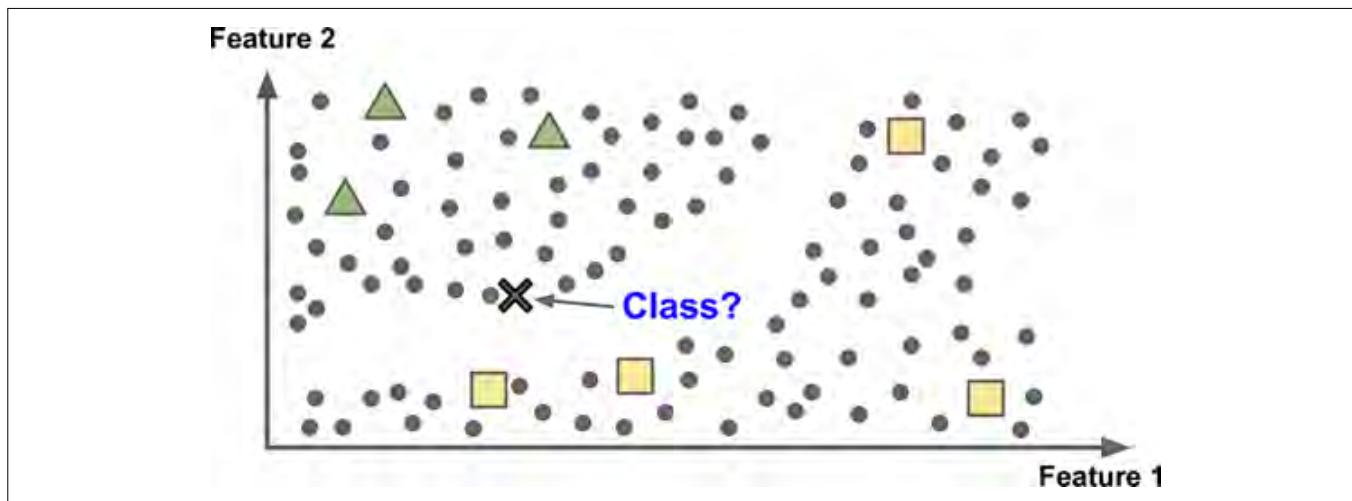


Figure 1-11. Semisupervised learning

Jordan, Michael I. and Tom M. Mitchell. (2015). “Machine Learning: Trends, perspectives, and prospects” *Science*.
Image from: Géron, Aurélien. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow* O'Reilly Media, Inc.

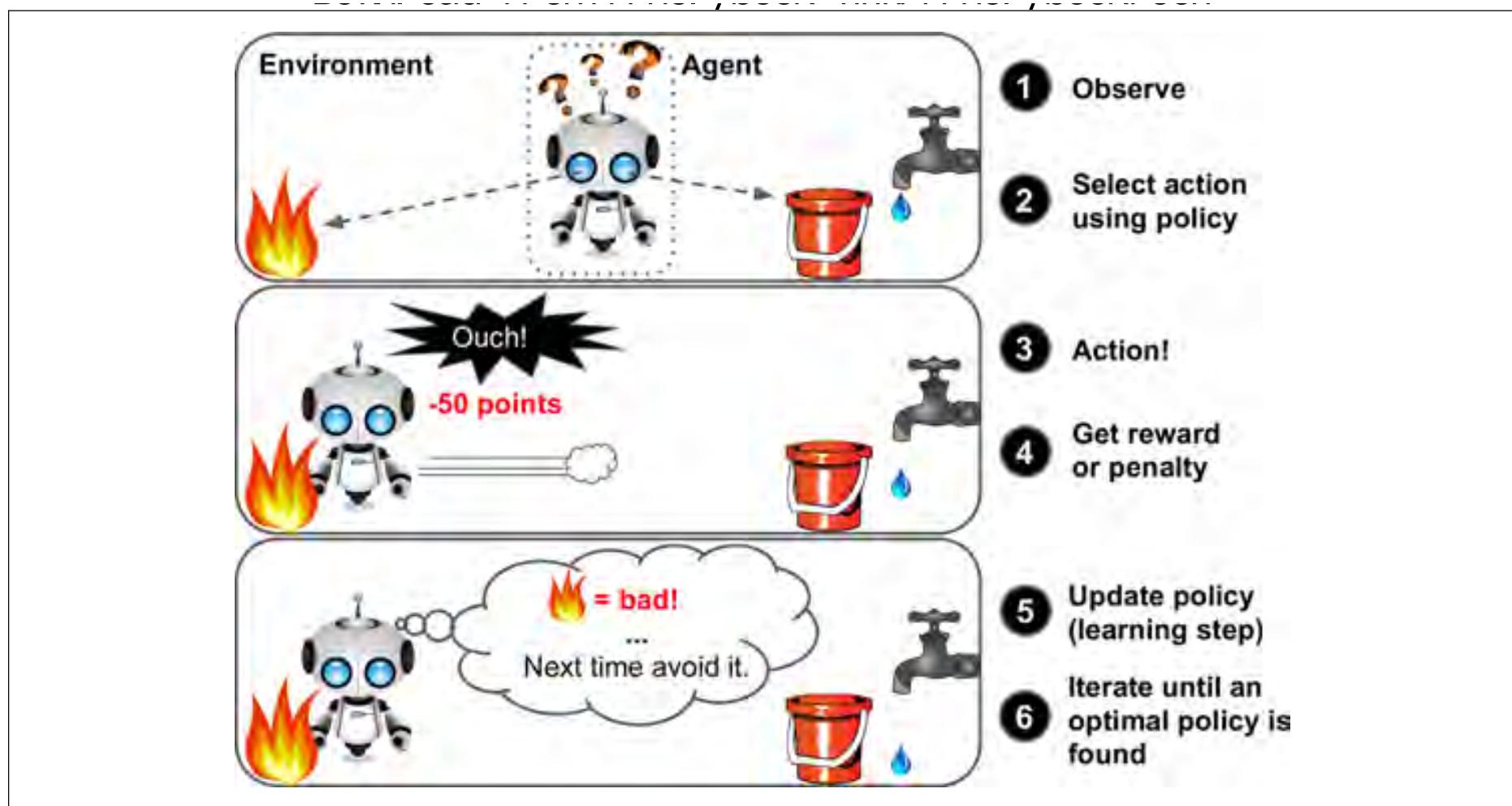
Reinforcement Learning

- ▶ “Instead of training examples that indicate the correct output for a given input, the training data in reinforcement learning are assumed to provide only an indication as to whether an action is correct or not; if an action is incorrect, there remains the problem of finding the correct action.”

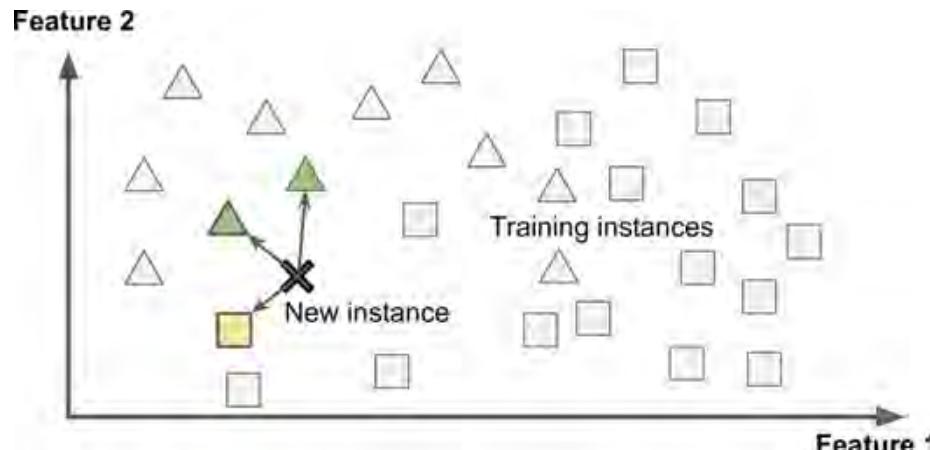
Reinforcement Learning

- ▶ “The learning system, called an agent in this context, can observe the environment, select and perform actions, and get rewards in return (or penalties in the form of negative rewards).”
- ▶ “It must then learn by itself what is the best strategy, called a policy, to get the most reward over time. A policy defines what action the agent should choose when it is in a given situation.”

Reinforcement Learning

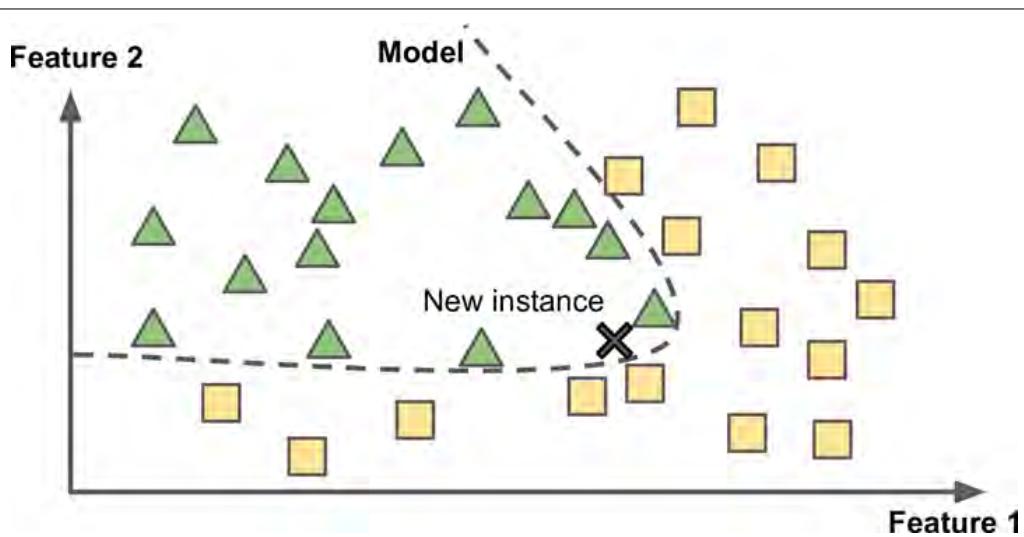


Instance versus Model-Based Learning



“the system learns the examples by heart, then generalizes to new cases using a similarity measure”

Figure 1-15. Instance-based learning



“another way to generalize from a set of examples is to build a model of these examples, then use that model to make predictions”

Figure 1-16. Model-based learning

Feature Engineering

- ▶ **Feature selection**
 - “selecting the most useful features to train on among existing features”
- ▶ **Feature extraction**
 - “combining existing features to produce a more useful one (as we saw earlier, dimensionality reduction algorithms can help)”
- ▶ **“Creating new features by gathering new data”**

Question Set 2

Géron (p. 22–30)

- ▶ **Can you name four of the main challenges in Machine Learning?**
- ▶ **If your model performs great on the training data but generalizes poorly to new instances, what is happening? Can you name three possible solutions?**
- ▶ **What is a test set and why would you want to use it?**
- ▶ **What is the purpose of a validation set?**
- ▶ **What can go wrong if you tune hyperparameters using the test set?**
- ▶ **What is cross-validation and why would you prefer it to a validation set?**