

# **Advanced Data Analysis**

**DATA 71200**

Class 2

# Course Schedule

29-Jan      Introduction

**5-Feb**      **What is Machine Learning?**

12-Feb      No Class

19-Feb      Getting Started with Machine Learning

26-Feb      Inspecting Data

4-Mar      Representing Data

# Class Website

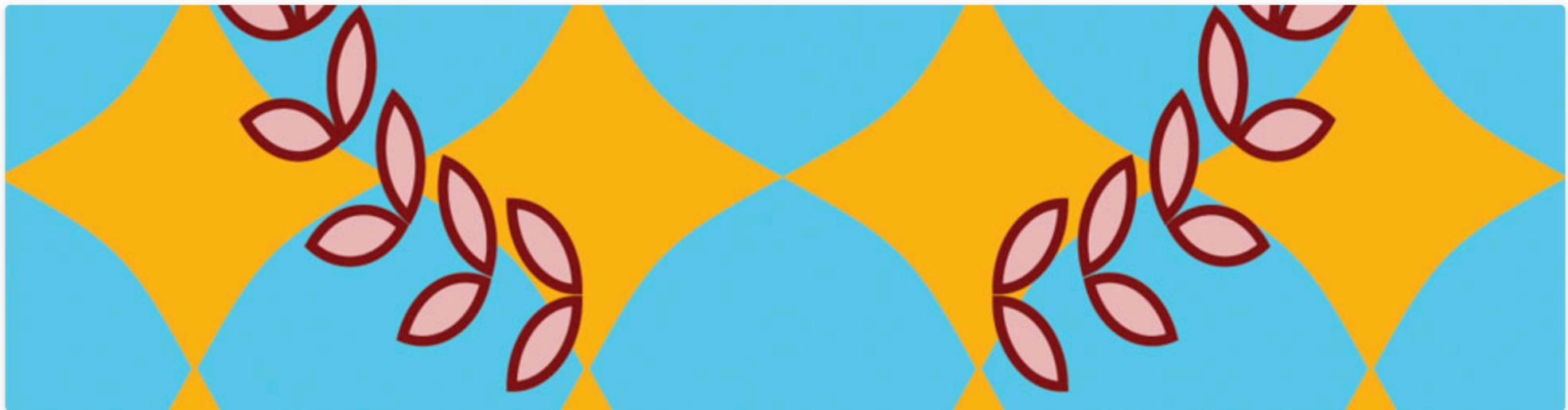
## DATA 71200 Advanced Data Analysis Methods

An introduction to supervised and unsupervised machine learning methods

---

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

---



**<https://data71200sp20.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.

# GitHub

 jcdevaney / **data71200sp20**

 Unwatch ▾ 1     Star 0     Fork 0

 Code     Issues 0     Pull requests 0     Actions     Projects 0     Wiki     Security     Insights     Settings

DATA 71200 Advanced Data Analysis Methods 

[Manage topics](#)

 1 commit     1 branch     0 packages     0 releases     1 contributor

Branch: master ▾    [New pull request](#)    [Create new file](#)    [Upload files](#)    [Find file](#)    [Clone or download](#) ▾

 jcdevaney Initial commit    Latest commit 02f220c 6 days ago

 README.md    Initial commit    6 days ago

 README.md 

## data71200sp20

DATA 71200 Advanced Data Analysis Methods

# GitHub - Forking versus Cloning

The screenshot shows a GitHub repository page for the project `onssen`. At the top right, the `Fork` button is highlighted with a black border. Below the header, there's a navigation bar with links for Code, Issues (1), Pull requests (0), Actions, Projects (0), Wiki, Security, and Insights. The main content area describes the project as "An open-source speech separation and enhancement library". It displays metrics: 28 commits, 2 branches, 0 packages, 0 releases, 1 contributor, and a license of GPL-3.0. A "Clone or download" button is also visible. The commit history shows two recent changes: one by `nateanl` creating a `LICENSE` file and another by `jcdevaney` adding `batch_norm` after `rnn`, refactoring training, and adding a `readme`. Both commits were made 3 months ago. A modal window titled "Fork onssen" is overlaid on the page, asking "Where should we fork onssen?" and showing a profile picture of `jcdevaney`.

Watch ▾ 11 Star 93 Fork 22

Code Issues 1 Pull requests 0 Actions Projects 0 Wiki Security Insights

An open-source speech separation and enhancement library

28 commits 2 branches 0 packages 0 releases 1 contributor GPL-3.0

Branch: master New pull request Create new file Upload files Find file Clone or download

nateanl Create LICENSE

jcdevaney Add batch\_norm after rnn, refactorize training, add readme

jcdevaney Add batch\_norm after rnn, refactorize training, add readme

Fork onssen

Where should we fork onssen?

jcdevaney

# GitHub - Forking versus Cloning

 jcdevaney / onssen

forked from speechLabBcCuny/onssen

 Watch ▼

0

 Star

0

 Fork

23

 Code

 Pull requests 0

 Actions

 Projects 0

 Wiki

 Security

 Insights

 Settings

An open-source speech separation and enhancement library

 Edit

[Manage topics](#)

 28 commits

 2 branches

 0 packages

 0 releases

 1 contributor

 GPL-3.0

Branch: master ▼

[New pull request](#)

[Create new file](#)

[Upload files](#)

[Find file](#)

[Clone or download ▼](#)

This branch is even with speechLabBcCuny:master.

 Pull request

 Compare

 nateanl Create LICENSE

Latest commit 0479d78 on Nov 29, 2019

 configs

Add batch\_norm after rnn, refactorize training, add readme

3 months ago

 data

Add batch\_norm after rnn, refactorize training, add readme

3 months ago

**Clone with HTTPS ?**

[Use SSH](#)

Use Git or checkout with SVN using the web URL.

<https://github.com/jcdevaney/onssen.git> 

[Open in Desktop](#)

[Download ZIP](#)

# GitHub Desktop

The screenshot shows the GitHub Desktop application interface. At the top, there's a dark header bar with three colored window control buttons (red, yellow, green) on the left. To the right of these are three dropdown menus: 'Current Repository' set to 'data71200sp20', 'Current Branch' set to 'master', and 'Fetch origin' with a note 'Last fetched 25 minutes ago'. Below the header is a light-colored main area. On the left side, there's a sidebar with tabs for 'Changes' (which is selected and highlighted in blue) and 'History'. Under the 'Changes' tab, it says '0 changed files'. In the center, the main content area has a large heading 'No local changes' and a subtext message: 'You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.' To the right of this text are two small, semi-transparent blue icons: one showing a map and another showing a person with a gear. Below the main message are three horizontal cards, each containing a suggestion and a 'View on GitHub' button. The first card says 'Open the repository in your external editor' (with 'preferences' link), 'Repository menu or ⌘ ⌘ A', and has a 'Open in Visual Studio Code' button. The second card says 'View the files in your repository in Finder' (with 'Repository menu or ⌘ ⌘ F' link) and has a 'Show in Finder' button. The third card says 'Open the repository page on GitHub in your browser' (with 'Repository menu or ⌘ ⌘ G' link) and has a 'View on GitHub' button. On the far left, there's a vertical sidebar with sections for 'Summary (required)', 'Description', and a 'Commit to master' button at the bottom.

An updated version of GitHub Desktop is available and will be installed at the next launch. See [what's new](#) or [restart GitHub Desktop](#).

Changes History

0 changed files

## No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

**Open the repository in your external editor**  
Configure which editor you wish to use in [preferences](#)

Repository menu or ⌘ ⌘ A

**Open in Visual Studio Code**

**View the files in your repository in Finder**  
Repository menu or ⌘ ⌘ F

**Show in Finder**

**Open the repository page on GitHub in your browser**  
Repository menu or ⌘ ⌘ G

**View on GitHub**

Summary (required)

Description

Commit to master

# 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**

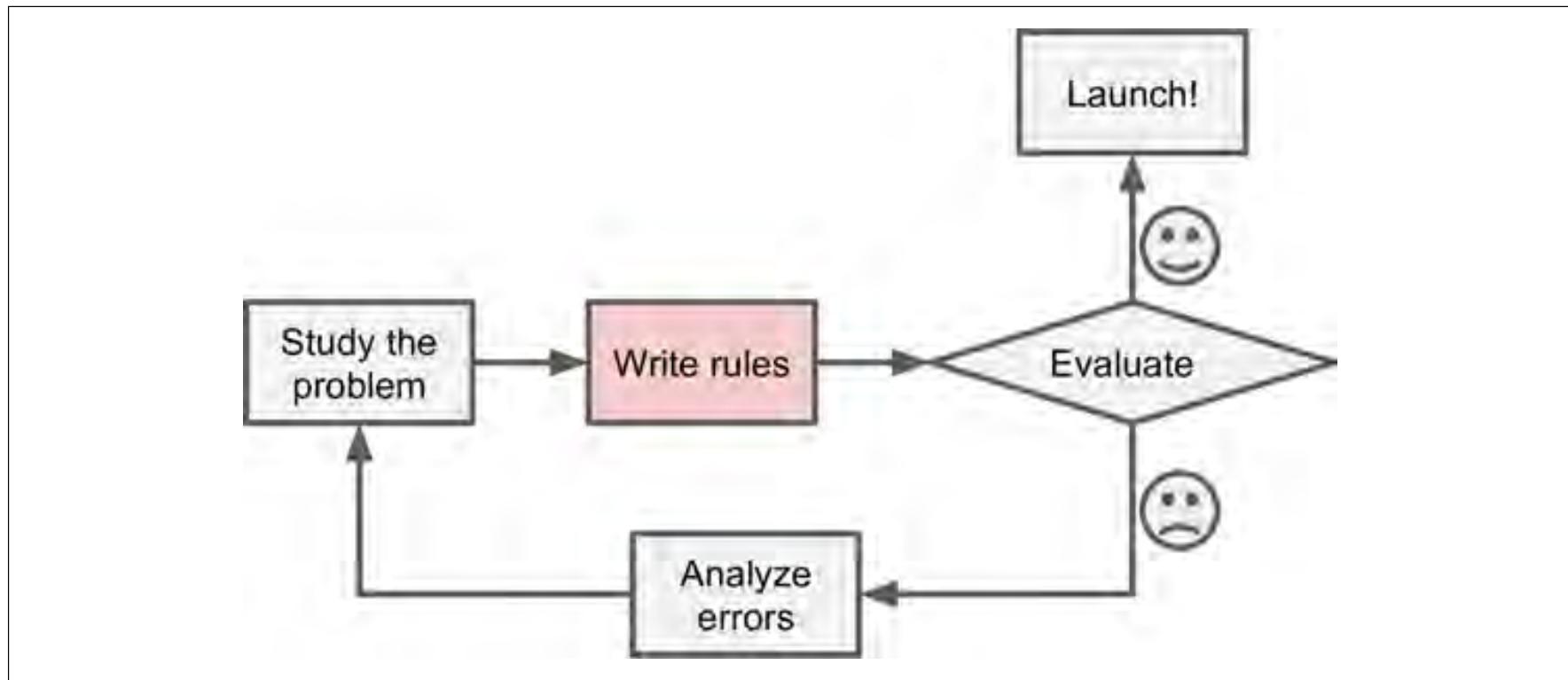


Tutorial: <https://machinelearningmastery.com/setup-python-environment-machine-learning-deep-learning-anaconda/>

# Key Questions

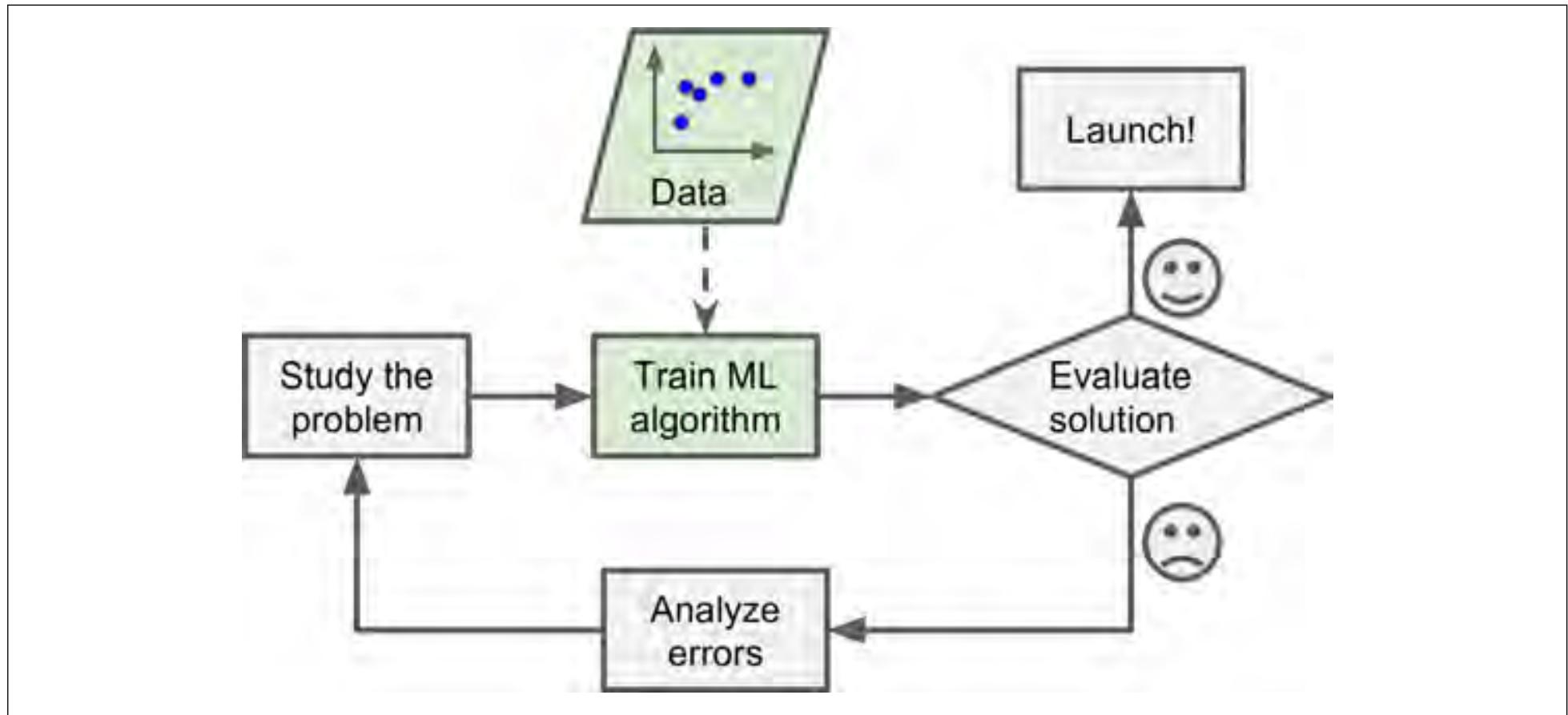
- ▶ “How can one construct computer systems that automatically improve through experience?”
- ▶ “What are the fundamental statistical-computational-information-theoretic laws that govern all learning systems, including computers, humans, and organizations?”
- ▶ “How accurately can the algorithm learn 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

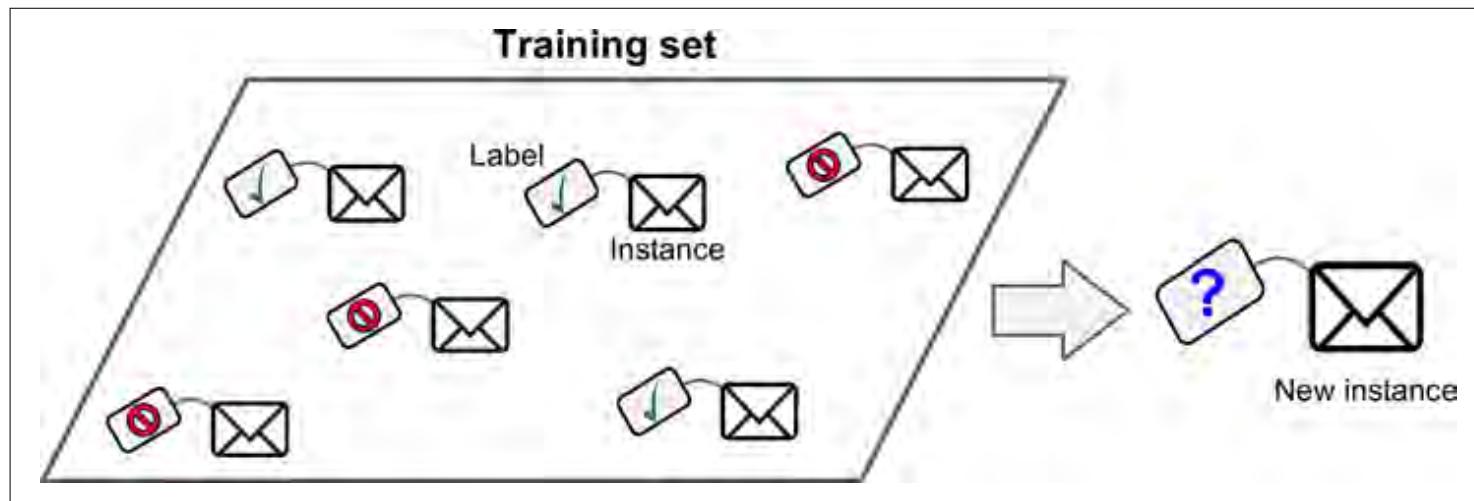


Figure 1-5. A labeled training set for supervised learning (e.g., spam classification)

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

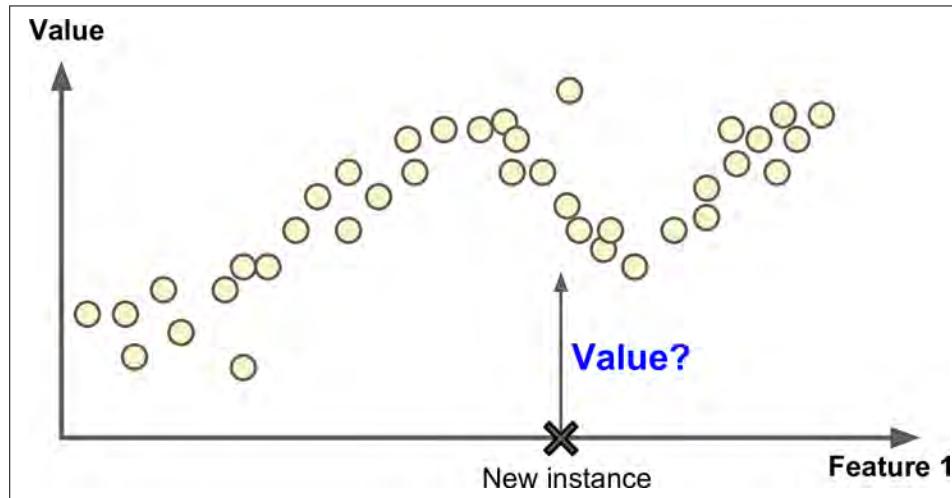
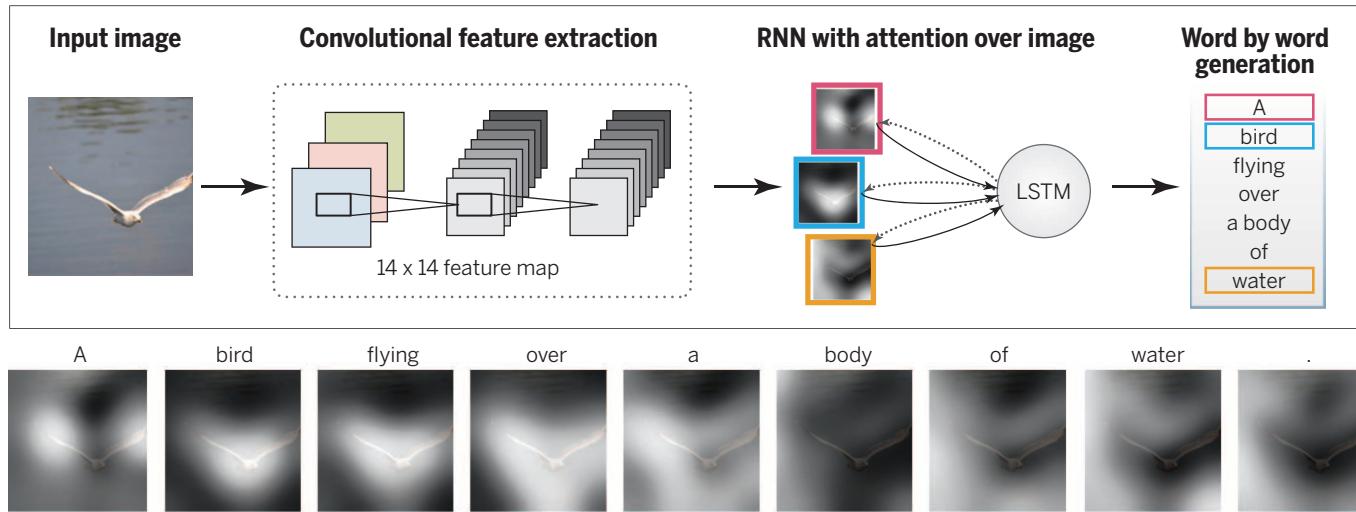


Figure 1-6. Regression



**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)”

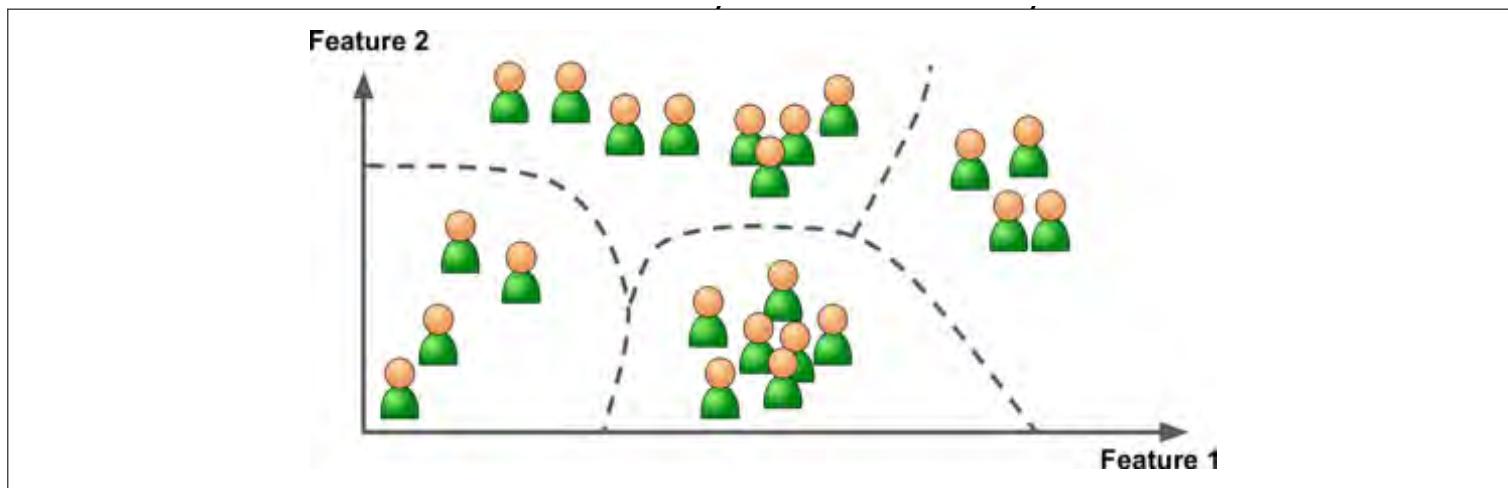
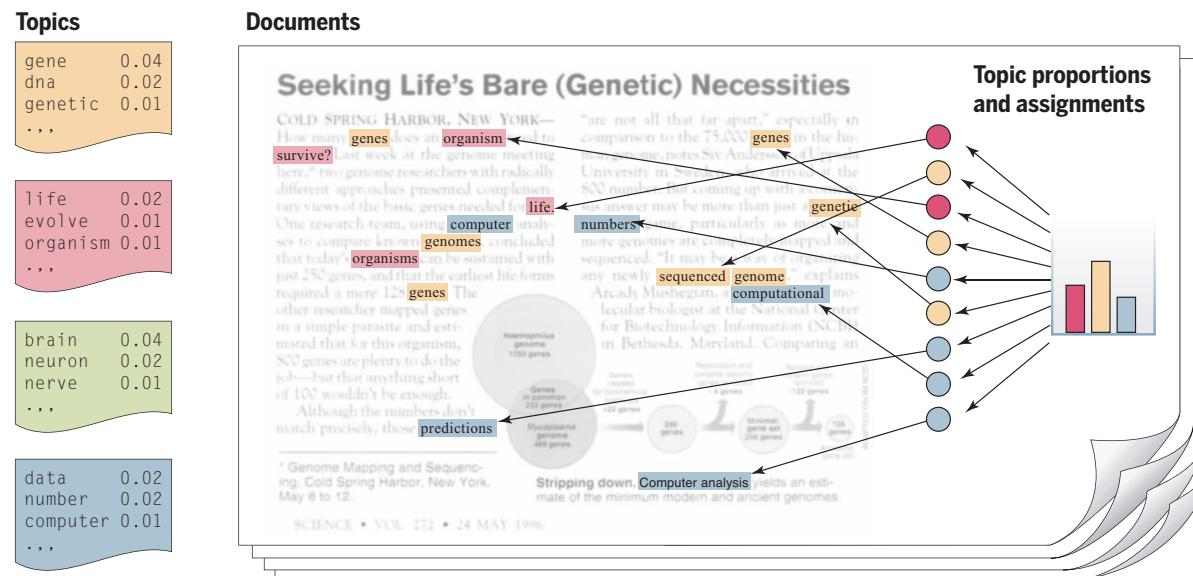


Figure 1-8. Clustering

# 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”

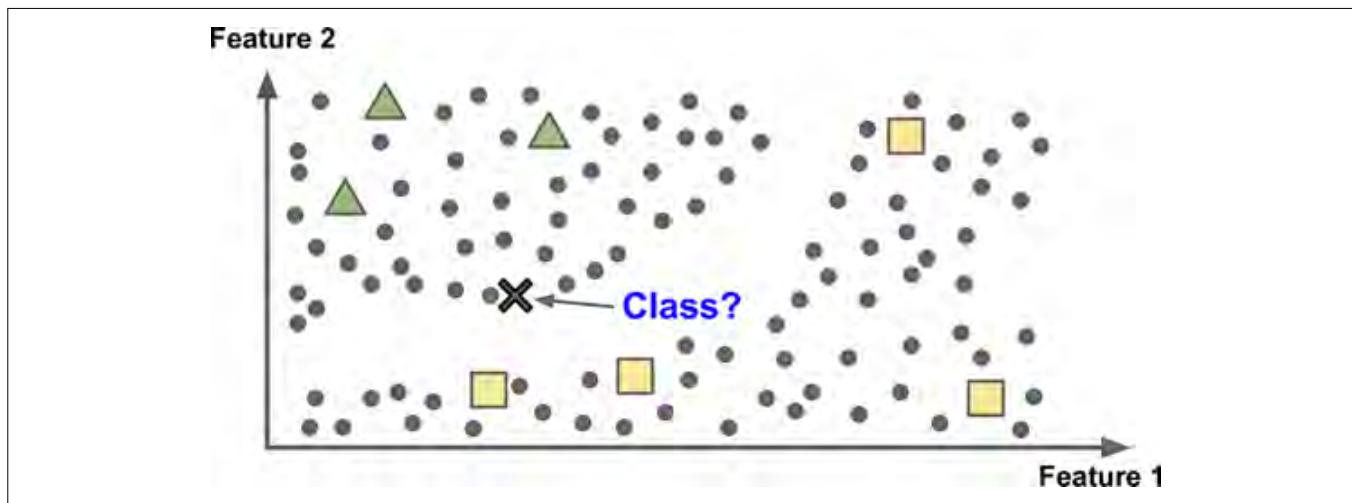


Figure 1-11. Semisupervised learning

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.

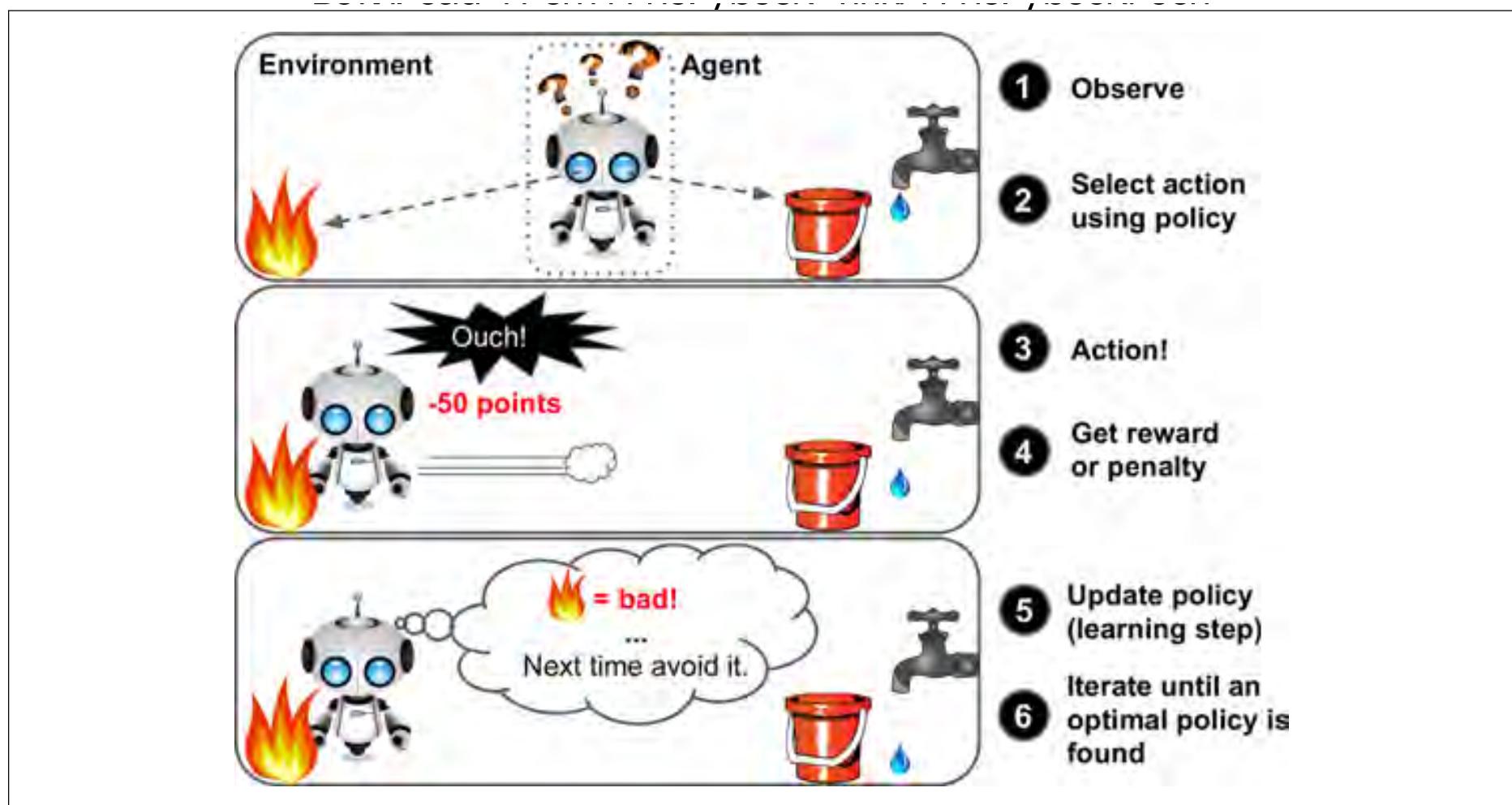
# Reinforcement Learning

- ▶ “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.”

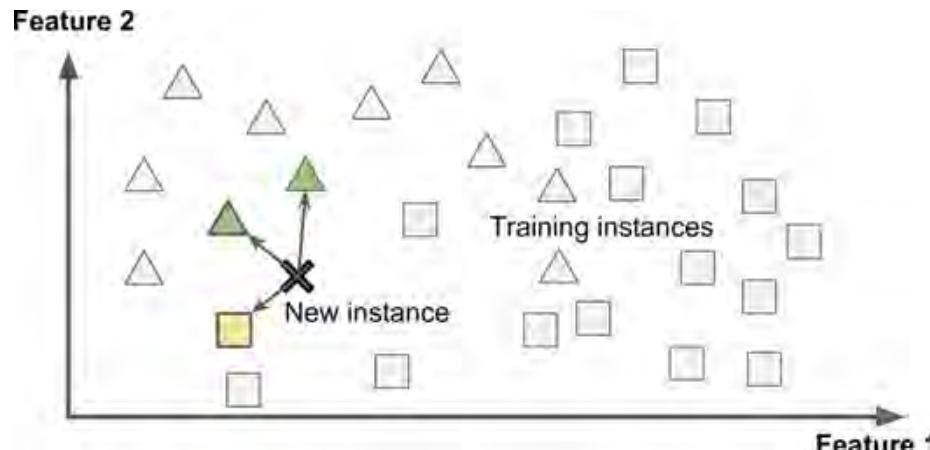
# Reinforcement Learning

- ▶ “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

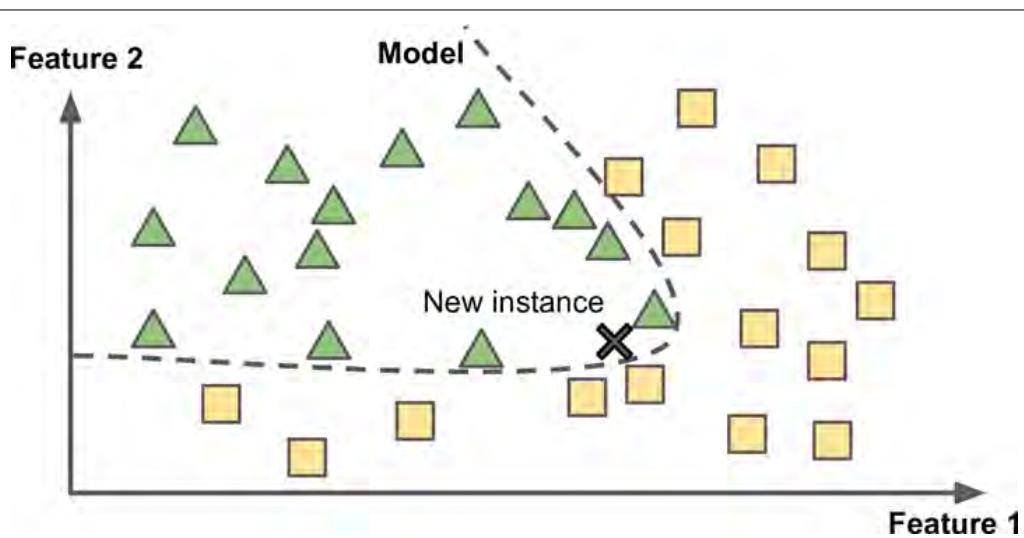


# Instance versus Model-Based Learning



**“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”**

# Review Questions 1

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

# Review Questions 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).”

# Review Questions 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.”

# Typical Machine Learning Project Steps

- ▶ “You studied the data.”
- ▶ “You selected a model.”
- ▶ Feature Engineering
- ▶ “You trained it on the training data (i.e., the learning algorithm searched for the model parameter values that minimize a cost function).”
- ▶ “Finally, you applied the model to make predictions on new cases (this is called inference), hoping that this model will generalize well.”

# Main Challenges

- ▶ **Insufficient training data**
  - Quantity and/or quality and/or non-representative
- ▶ **Irrelevant features**
- ▶ **Overfitting training data**
- ▶ **Under-fitting training data**

# Example: GDP and Happiness

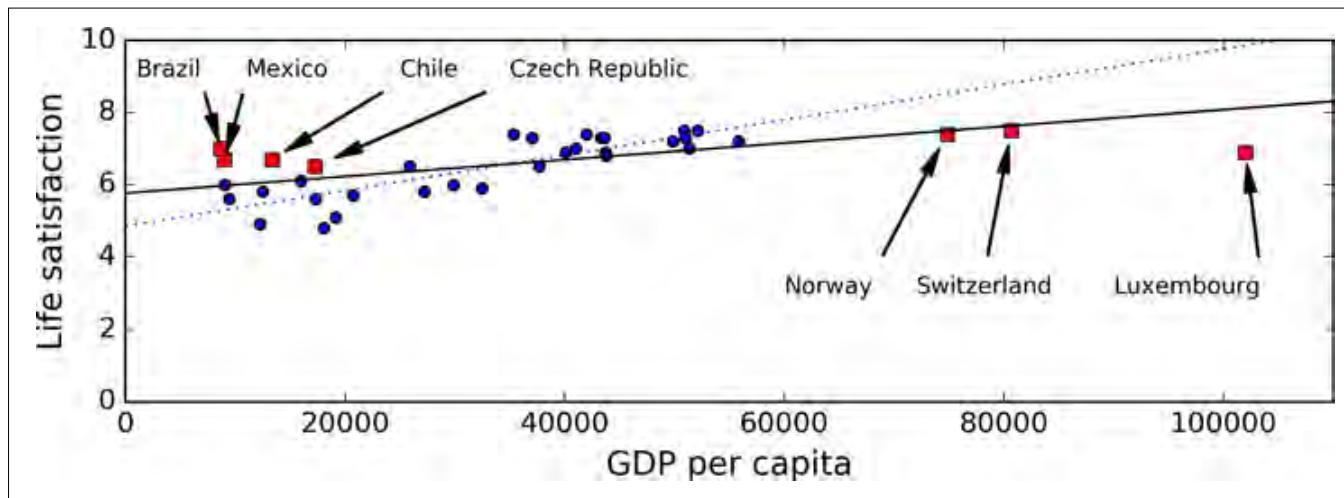


Figure 1-21. A more representative training sample

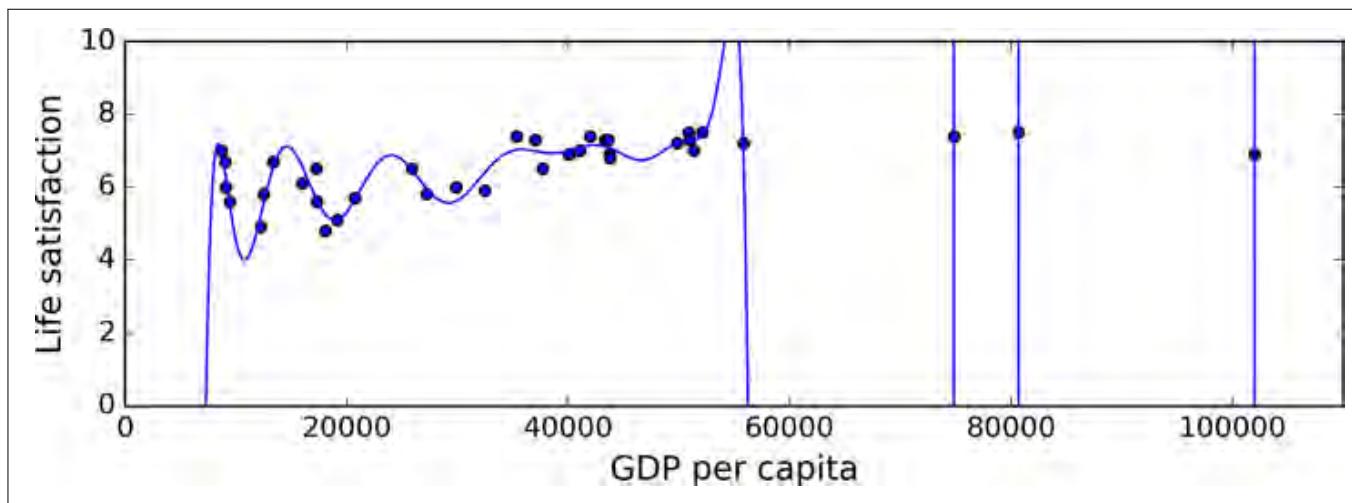


Figure 1-22. Overfitting the training data

# Example: GDP and Happiness

- ▶ **regularization**

- “constraining a model to make it simpler and reduce the risk of overfitting”

- ▶ **hyperparameter**

- “amount of regularization to apply during learning”
- “need to be set before training”

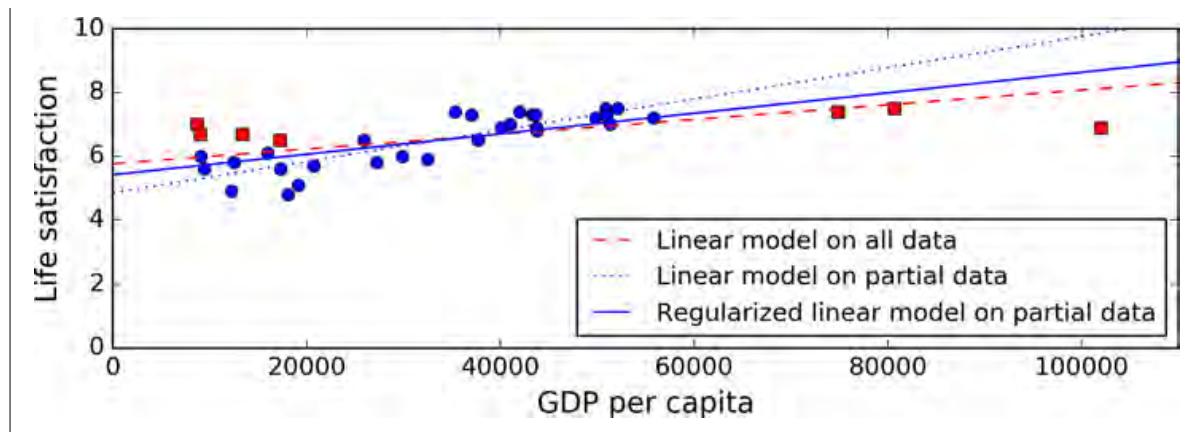
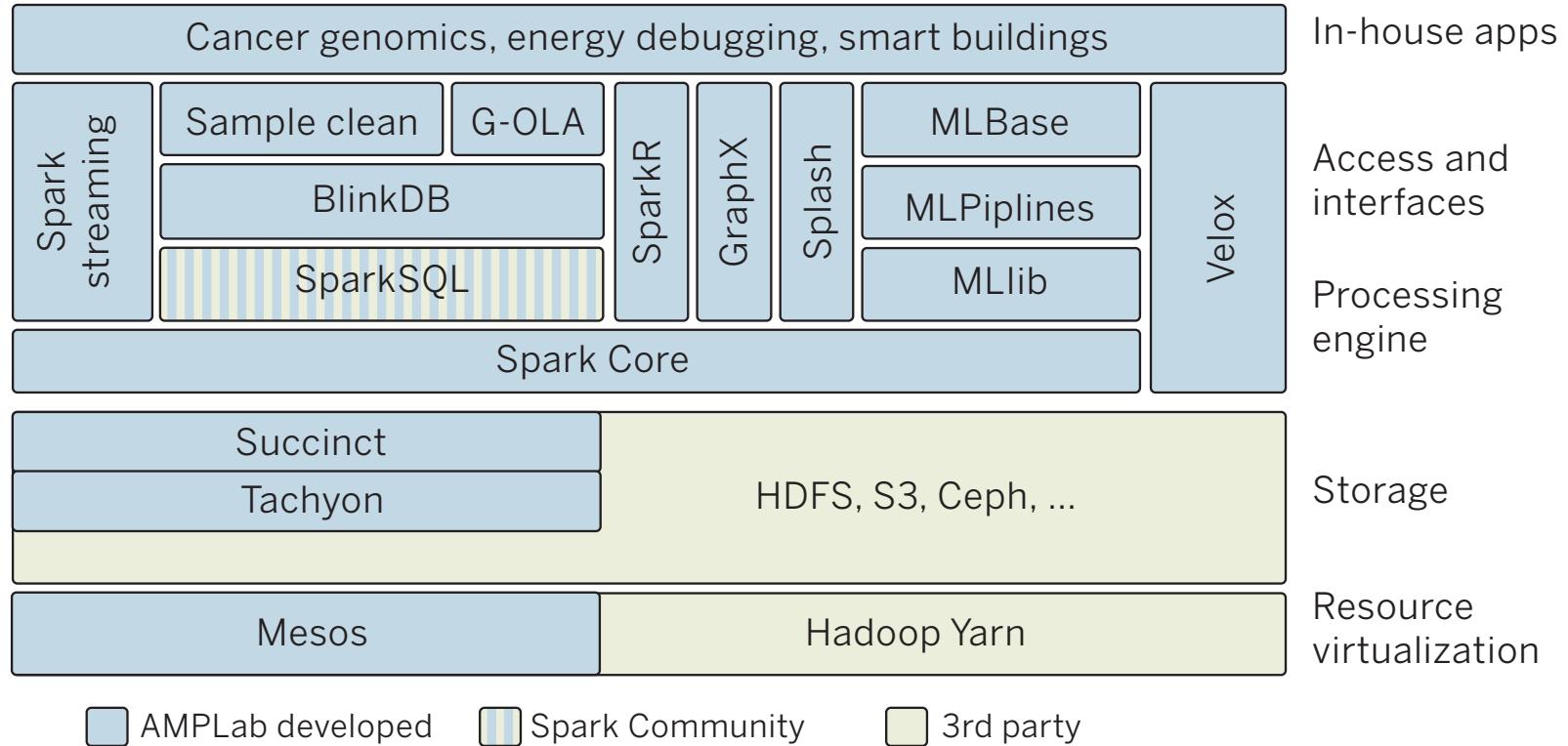


Figure 1-23. Regularization reduces the risk of overfitting

# Training/Testing/Validation

- ▶ **Training set - data used to train the model**
- ▶ **Testing set - hold out data used to estimate the generalization error on new data**
- ▶ **Validation set - used to compare models**
- ▶ **Cross-validation - iteratively holding out a subset of the data and testing on the rest (typically 80/20)**

# Data Analysis in Practice



**Fig. 5. Data analytics stack.** Scalable machine-learning systems are layered architectures that are built on parallel and distributed computing platforms. The architecture depicted here—an open-source data analysis stack developed in the Algorithms, Machines and People (AMP) Laboratory at the University of California, Berkeley—includes layers that interface to underlying operating systems; layers that provide distributed storage, data management, and processing; and layers that provide core machine-learning competencies such as streaming, subsampling, pipelines, graph processing, and model serving.

# Review Questions 2

- ▶ **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?**

# Review Questions 2

- ▶ **Can you name four of the main challenges in Machine Learning?**
  - “Some of the main challenges in Machine Learning are the lack of data, poor data quality, non-representative data, uninformative features, excessively simple models that underfit the training data, and excessively complex models that overfit the data.”

# Review Questions 2

- ▶ **If your model performs great on the training data but generalizes poorly to new instances, what is happening? Can you name three possible solutions?**
  - “If a model performs great on the training data but generalizes poorly to new instances, the model is likely overfitting the training data (or we got extremely lucky on the training data). Possible solutions to overfitting are getting more data, simplifying the model (selecting a simpler algorithm, reducing the number of parameters or features used, or regularizing the model), or reducing the noise in the training data.”

# Review Questions 2

- ▶ **What is a test set and why would you want to use it?**
  - “A test set is used to estimate the generalization error that a model will make on new instances, before the model is launched in production.”
- ▶ **What is the purpose of a validation set?**
  - “A validation set is used to compare models. It makes it possible to select the best model and tune the hyperparameters.”

# Review Questions 2

- ▶ **What can go wrong if you tune hyperparameters using the test set?**
  - “If you tune hyperparameters using the test set, you risk overfitting the test set, and the generalization error you measure will be optimistic (you may launch a model that performs worse than you expect).”
- ▶ **What is cross-validation and why would you prefer it to a validation set?**
  - “Cross-validation is a technique that makes it possible to compare models (for model selection and hyperparameter tuning) without the need for a separate validation set. This saves precious training data.”