# List of Figures

Chapter 1. A machine-learning odyssey

Chapter 4. A gentle introduction to classification

Chapter 5. Automatically clustering data

Figure 7.5. dot producttwo-input networkThe input dimension can be arbitrarily long. For example, each pixel in a grayscale image can have a corresponding input xi. This neural network uses all inputs to generate a single output number, which you might use for regression or classification. The notation wT means you're transposing w, which is an $n \times 1$ vector, into a $1 \times n$ vector. That way, you can properly multiply it with x (which has the dimensions $n \times 1$). Such a matrix multiplication is also called a dot product, and it yields a scalar (one-dimensional) value.

Figure 7.6. Nodes that don't interface to both the input and the output are called hidden neurons. A hidden layer is a collection of hidden units that aren't connected to each other.

Figure 7.7. If you want to create a network where the input equals the output, you can connect the corresponding nodes and set each parameter's weight to 1.

Figure 7.8. Here, you introduce a restriction to a network that tries to reconstruct its input. Data will pass through a narrow channel, as illustrated by the hidden layer. In this example, there's only one node in the hidden layer. This network is trying to encode (and decode) an n-dimensional input signal into just one dimension, which will likely be difficult in practice.

Figure 7.9. A colored image is composed of pixels, and each pixel contains values for red, green, and blue.

Figure 7.10. An image can be represented in row-major order. That way, you can represent a two-dimensional structure as a onedimensional structure.

## Chapter 8. Reinforcement learning

Figure 8.1. A person navigating to reach a destination in the midst of traffic and unexpected situations is a problem setup for reinforcement learning.

Figure 8.2. Actions are represented by arrows, and states are represented by circles. Performing an action on a state produces a reward. If you start at state s1, you can perform action a1 to obtain a reward r(s1, a1).

Figure 8.3. A policy suggests which action to take, given a state.

Figure 8.4. Given a state and the action taken, applying a utility function Q predicts the expected and the total rewards: the immediate reward (next state) plus rewards gained later by following an optimal policy.

Figure 8.5. Ideally, our algorithm should buy low and sell high. Doing so just once, as shown here, might yield a reward of around $160. But the real profit rolls in when you buy and sell more frequently. Ever heard the term high-frequency trading? It's about buying low and selling high as frequently as possible to maximize profits within a period of time.

Figure 8.6. This chart summarizes the opening stock prices of Microsoft (MSFT) from 7/22/1992 to 7/22/2016. Wouldn't it have been nice to buy around day 3000 and sell around day 5000? Let's see if our code can learn to buy, sell, and hold to make optimal gain.

Figure 8.7. Most reinforcement-learning algorithms boil down to just three main steps: infer, do, and learn. During the first step, the algorithm selects the best action (a), given a state (s), using the knowledge it has so far. Next, it does the action to find out the reward (r) as well as the next state (s'). Then it improves its understanding of the world by using the newly acquired knowledge (s, r, a, s').

Figure 8.8. A rolling window of a certain size iterates through the stock prices, as shown by the chart segmented to form states S1, S2, and S3. The policy suggests an action to take: you may either choose to exploit it or randomly explore another action. As you get rewards for performing an action, you can update the policy function over time.

Figure 8.9. The input is the state space vector, with three outputs: one for each output's Q-value.

Figure 8.10. The algorithm learns a good policy to trade Microsoft stocks.

## Chapter 9. Convolutional neural networks

Figure 9.1. In a fully connected network, each pixel of an image is treated as an input. For a grayscale image of size $256 \times 256$, that's $256 \times 256$ neurons! Connecting each neuron to 10 outputs yields $256 \times 256 \times 10 = 655,360$ weights.

Figure 9.2. Convolving a $5 \times 5$ patch over an image, as shown on the left, produces another image, as shown on the right. In this case, the produced image is the same size as the original. Converting an original image to a convolved image requires only $5 \times 5 = 25$ parameters!

Figure 9.3. Images from the CIFAR-10 dataset. Because they're only $32 \times 32$ in size, they're a bit difficult to see, but you can generally recognize some of the objects.

Figure 9.4. These are 32 randomly initialized matrices, each of size $5 \times 5$. They represent the filters you'll use to convolve an input image.

Figure 9.5. An example $24 \times 24$ image from the CIFAR-10 dataset

Figure 9.6. Resulting images from convolving the random filters on an image of a car

Figure 9.7. After you add a bias term and an activation function, the resulting convolutions can capture more-powerful patterns within images.

Figure 9.8. After running maxpool, the convolved outputs are halved in size, making the algorithm computationally faster without losing too much information.

Figure 9.9. An input image is convolved by multiple $5 \times 5$ filters. The convolution layer includes an added bias term with an activation function, resulting in $5 \times 5 + 5 = 30$ parameters. Next, a max-pooling layer reduces the dimensionality of the data (which requires no extra parameters).

## Chapter 10. Recurrent neural networks

Figure 10.1. A neural network with the input and output layers labeled as X(t) and Y(t), respectively

Chapter 11. Sequence-to-sequence models for chatbots

Figure A.5. A possible error message from running the TensorFlow container

Figure A.6. Listing and killing a Docker container to get rid of the error message in figure A.5