

# Deep Learning with Python

## Chapter 4: Fundamentals of machine learning

**4.1.1** *Supervised learning*

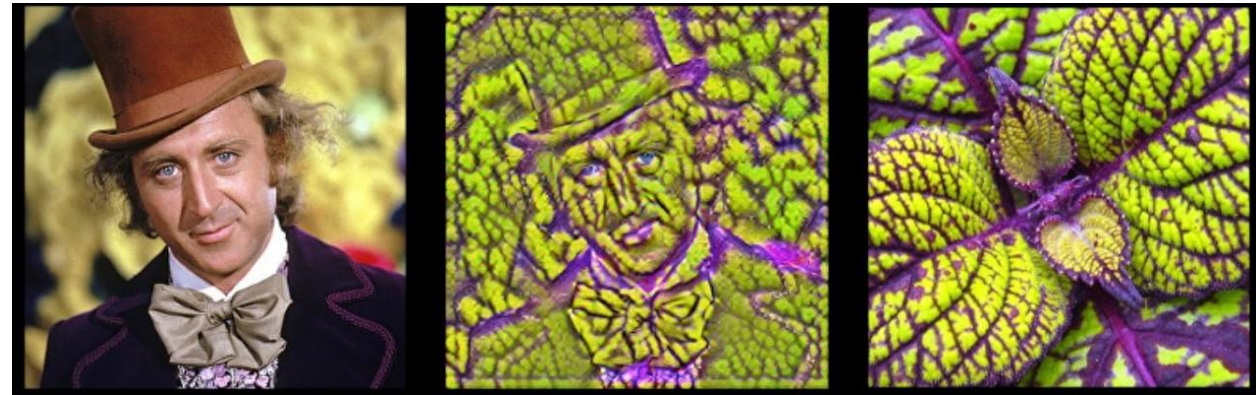
**4.1.2** *Unsupervised learning*

**4.1.3** *Self-supervised learning*

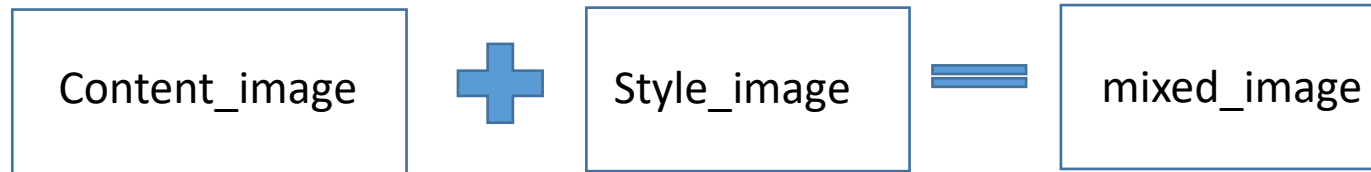
**4.1.4** *Reinforcement learning*



example\_1



example\_2



## **4.2    *Evaluating machine-learning models***

### **4.2.1    *Training, validation, and test sets***

## Listing 4.1 Hold-out validation

```
num_validation_samples = 10000
```

```
np.random.shuffle(data)
```

Shuffling the data is usually appropriate.

```
validation_data = data[:num_validation_samples]
```

Defines the validation set

```
data = data[num_validation_samples:]
```

Defines the training set

```
training_data = data[:]
```

```
model = get_model()
```

```
model.train(training_data)
```

```
validation_score = model.evaluate(validation_data)
```

Trains a model on the training data, and evaluates it on the validation data

```
# At this point you can tune your model,  
# retrain it, evaluate it, tune it again...
```

```
model = get_model()
```

```
model.train(np.concatenate([training_data,  
                             validation_data]))
```

```
test_score = model.evaluate(test_data)
```

Once you've tuned your hyperparameters, it's common to train your final model from scratch on all non-test data available.

# K-FOLD VALIDATION

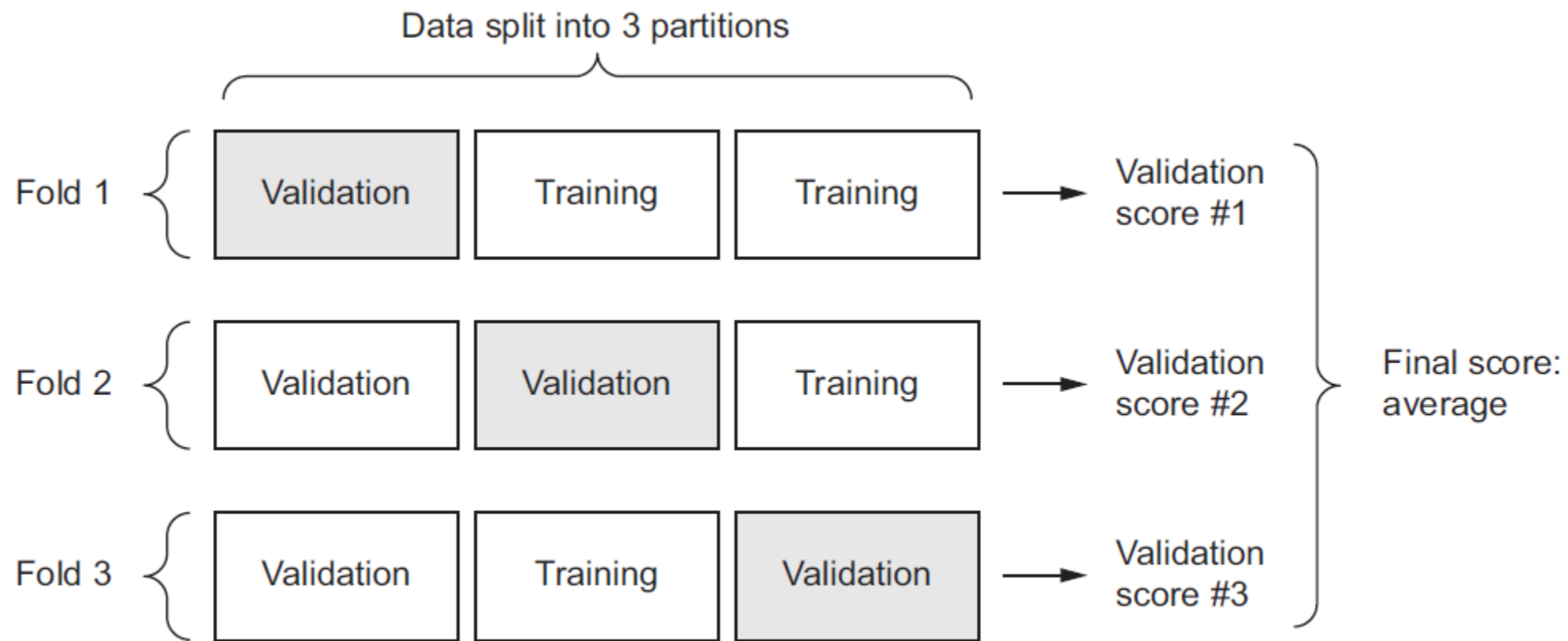


Figure 4.2 Three-fold validation

## Listing 4.2 K-fold cross-validation

```
k = 4
num_validation_samples = len(data) // k
```

```
np.random.shuffle(data)
```

```
validation_scores = []
```

```
for fold in range(k):
```

```
    validation_data = data[num_validation_samples * fold:
        num_validation_samples * (fold + 1)]
```

```
    training_data = data[:num_validation_samples * fold] +
        data[num_validation_samples * (fold + 1):]
```

```
    model = get_model()
```

```
    model.train(training_data)
```

```
    validation_score = model.evaluate(validation_data)
```

```
    validation_scores.append(validation_score)
```

```
validation_score = np.average(validation_scores)
```

```
model = get_model()
```

```
model.train(data)
```

```
test_score = model.evaluate(test_data)
```

**Selects the validation-  
data partition**

**Uses the remainder of the data  
as training data. Note that the  
+ operator is list concatenation,  
not summation.**

**Creates a brand-new instance  
of the model (untrained)**

**Trains the final  
model on all non-  
test data available**

**Validation score:  
average of the  
validation scores  
of the k folds**

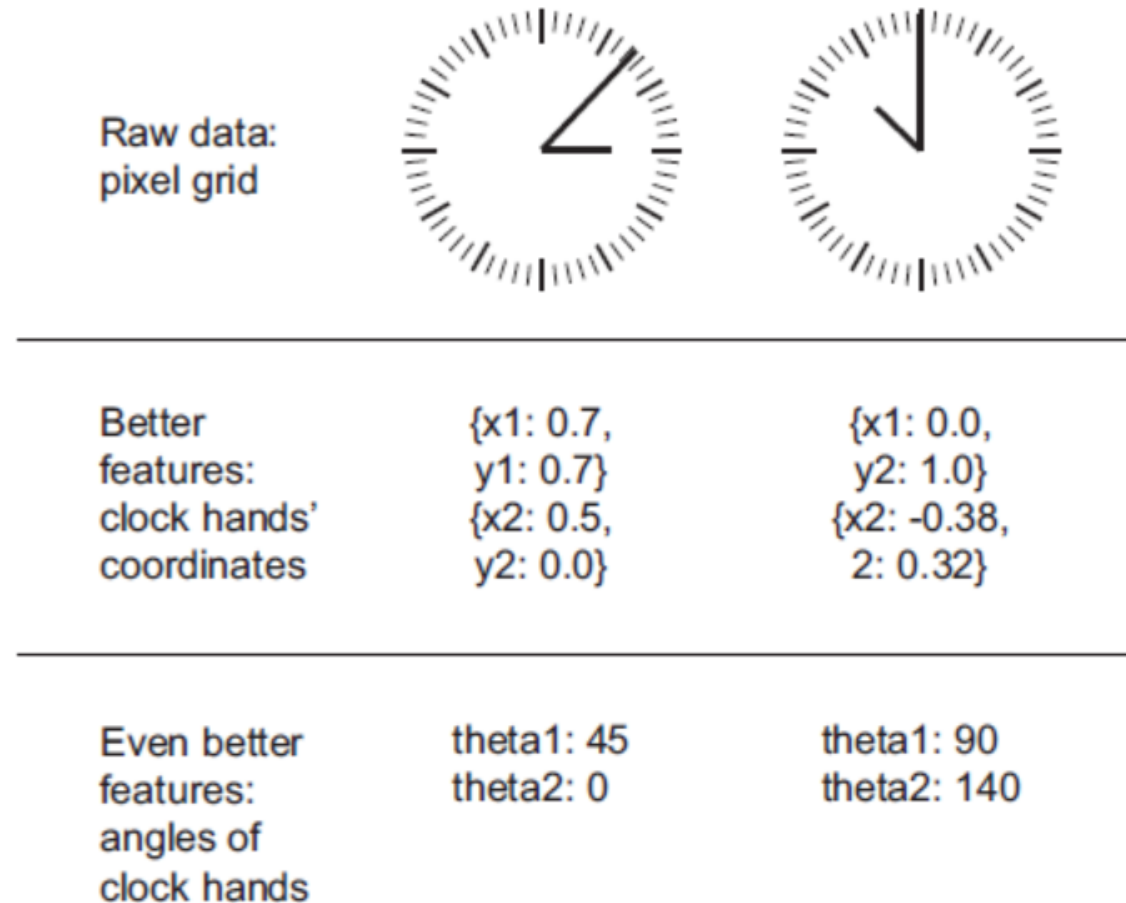
## HANDLING MISSING VALUES

You may sometimes have missing values in your data. For instance, in the house-price example, the first feature (the column of index 0 in the data) was the per capita crime rate. What if this feature wasn't available for all samples? You'd then have missing values in the training or test data.

In general, with neural networks, it's safe to input missing values as 0, with the condition that 0 isn't already a meaningful value. The network will learn from exposure to the data that the value 0 means *missing data* and will start ignoring the value.



## Feature engineering



**Figure 4.3** Feature engineering for reading the time on a clock

Fortunately, modern deep learning removes the need for most feature engineering, because neural networks are capable of automatically extracting useful features from raw data. Does this mean you don't have to worry about feature engineering as long as you're using deep neural networks? No, for two reasons:

- Good features still allow you to solve problems more elegantly while using fewer resources. For instance, it would be ridiculous to solve the problem of reading a clock face using a convolutional neural network.
- Good features let you solve a problem with far less data. The ability of deep-learning models to learn features on their own relies on having lots of training data available; if you have only a few samples, then the information value in their features becomes critical.

# *Overfitting and underfitting*

The fundamental issue in machine learning is the tension between optimization and generalization. *Optimization* refers to the process of adjusting a model to get the best performance possible on the training data (the *learning* in *machine learning*), whereas *generalization* refers to how well the trained model performs on data it has never seen before. The goal of the game is to get good generalization, of course, but you don't control generalization; you can only adjust the model based on its training data.

The processing of fighting overfitting this way is called *regularization*. Let's review some of the most common regularization techniques and apply them in practice to improve the movie-classification model from section 3.4.

## *Reducing the network's size*

Let's try this on the movie-review classification network. The original network is shown next.

#### Listing 4.3 Original model

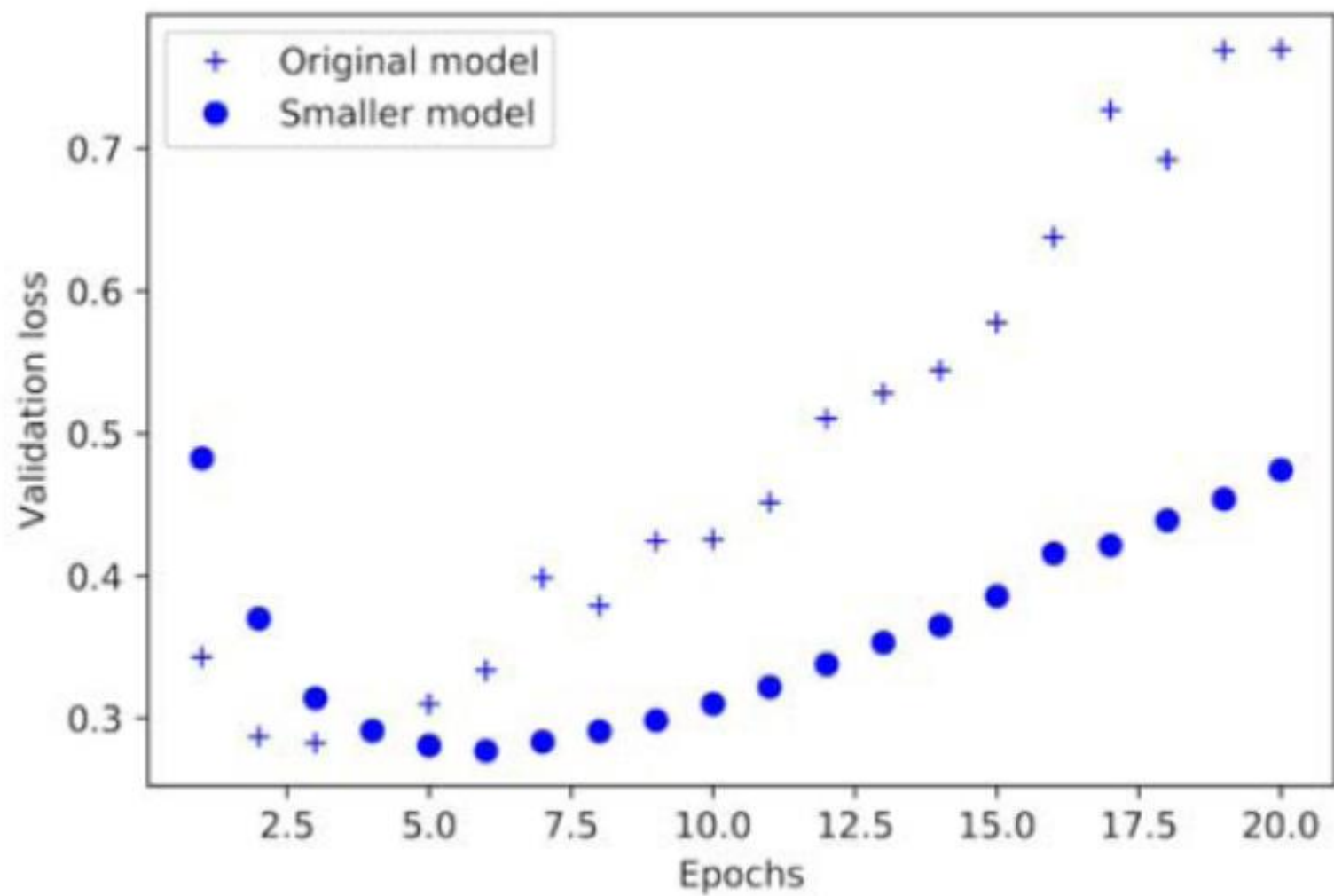
```
from keras import models
from keras import layers

model = models.Sequential()
model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(16, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))
```

Now let's try to replace it with this smaller network.

#### Listing 4.4 Version of the model with lower capacity

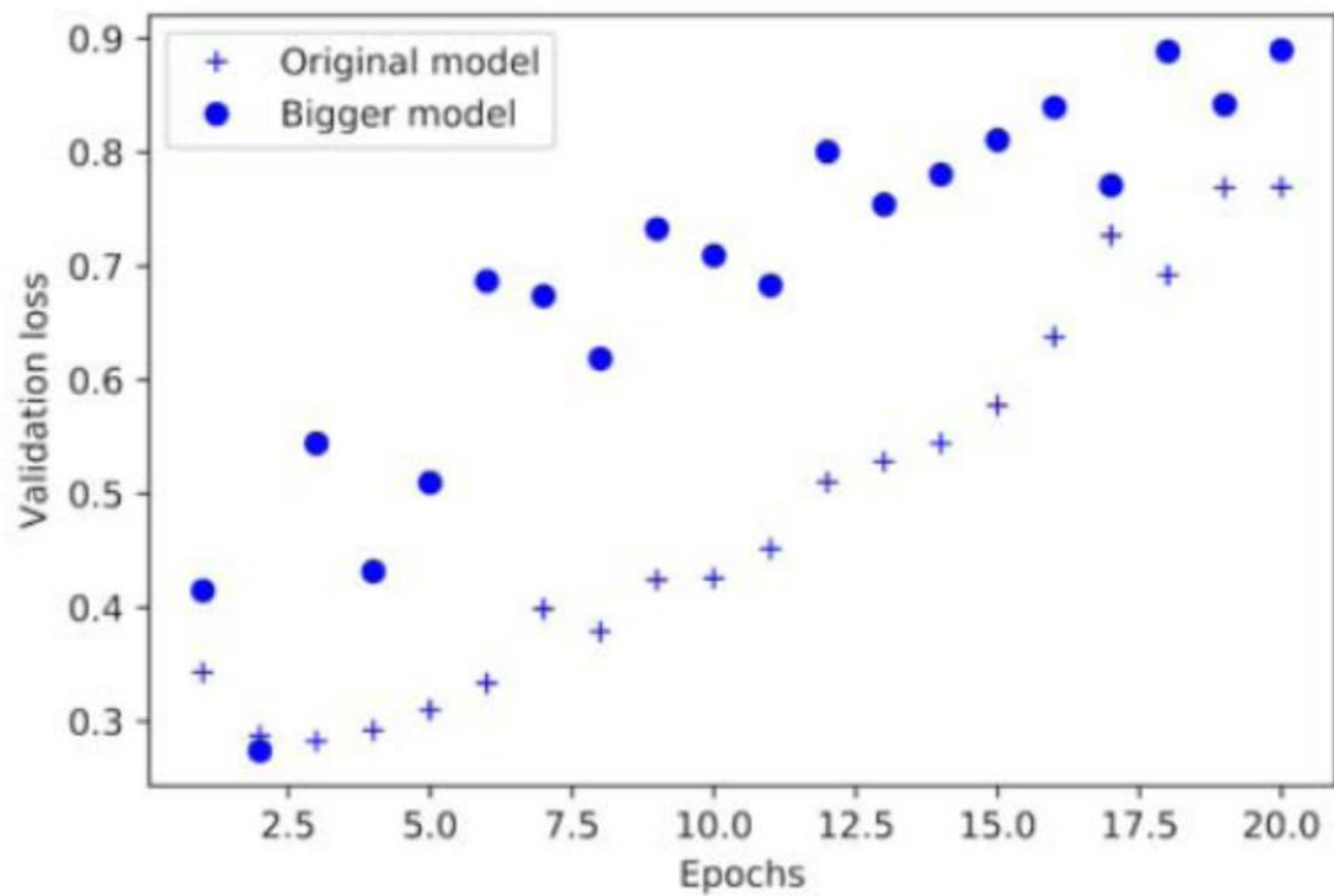
```
model = models.Sequential()
model.add(layers.Dense(4, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(4, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))
```



**Figure 4.4** Effect of model capacity on validation loss: trying a smaller model

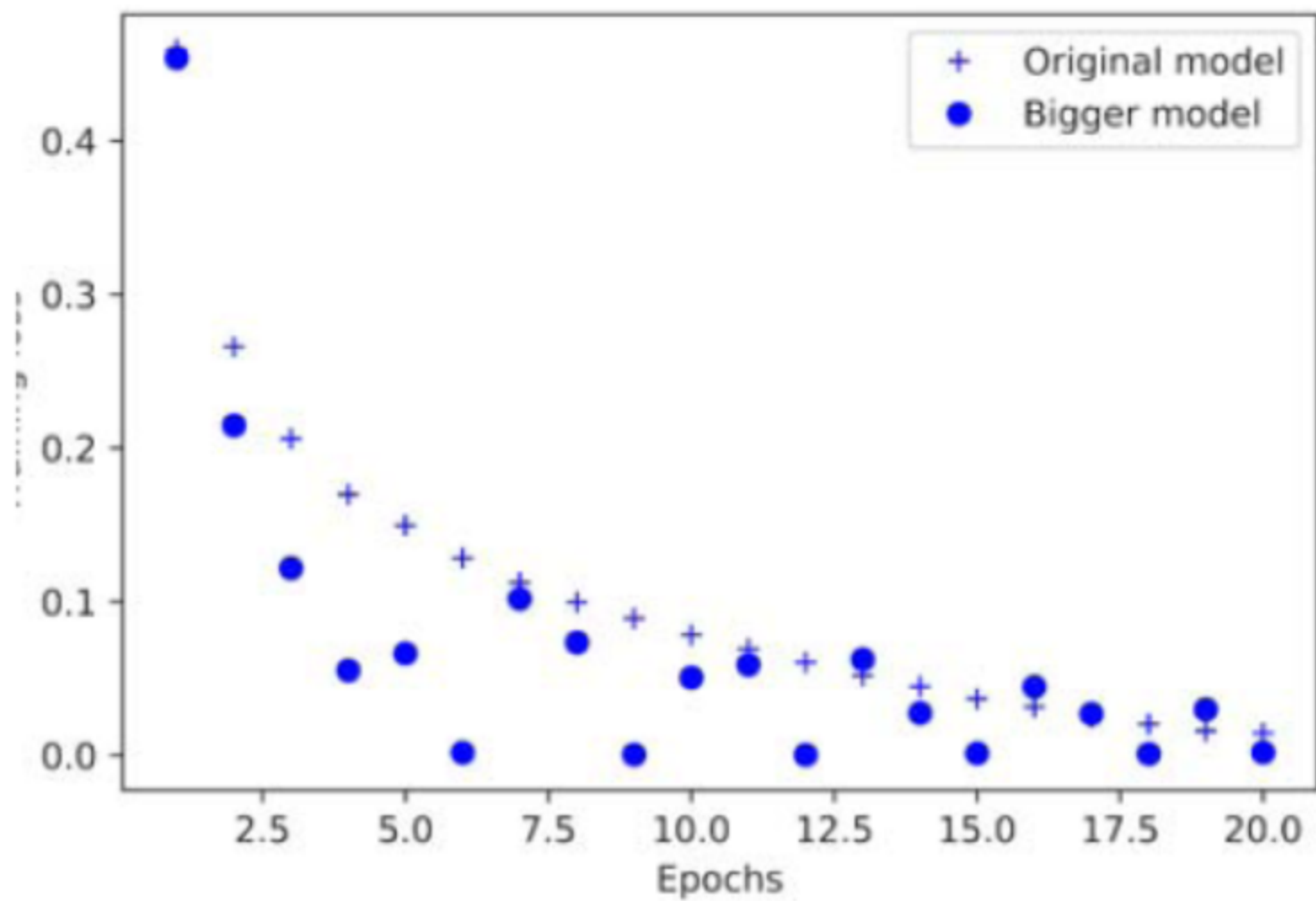
#### Listing 4.5 Version of the model with higher capacity

```
model = models.Sequential()  
model.add(layers.Dense(512, activation='relu', input_shape=(10000,)))  
model.add(layers.Dense(512, activation='relu'))  
model.add(layers.Dense(1, activation='sigmoid'))
```



**Figure 4.5** Effect of model capacity on validation loss: trying a bigger model





**Figure 4.6** Effect of model capacity on training loss: trying a bigger model

The processing of fighting overfitting this way is called *regularization*. Let's review some of the most common regularization techniques and apply them in practice to improve the movie-classification model from section 3.4.

*Reducing the network's size*

*Adding weight regularization*

A *simple model* in this context is a model where the distribution of parameter values has less entropy (or a model with fewer parameters, as you saw in the previous section). Thus a common way to mitigate overfitting is to put constraints on the complexity of a network by forcing its weights to take only small values, which makes the distribution of weight values more *regular*. This is called *weight regularization*, and it's done by adding to the loss function of the network a *cost* associated with having large weights. This cost comes in two flavors:

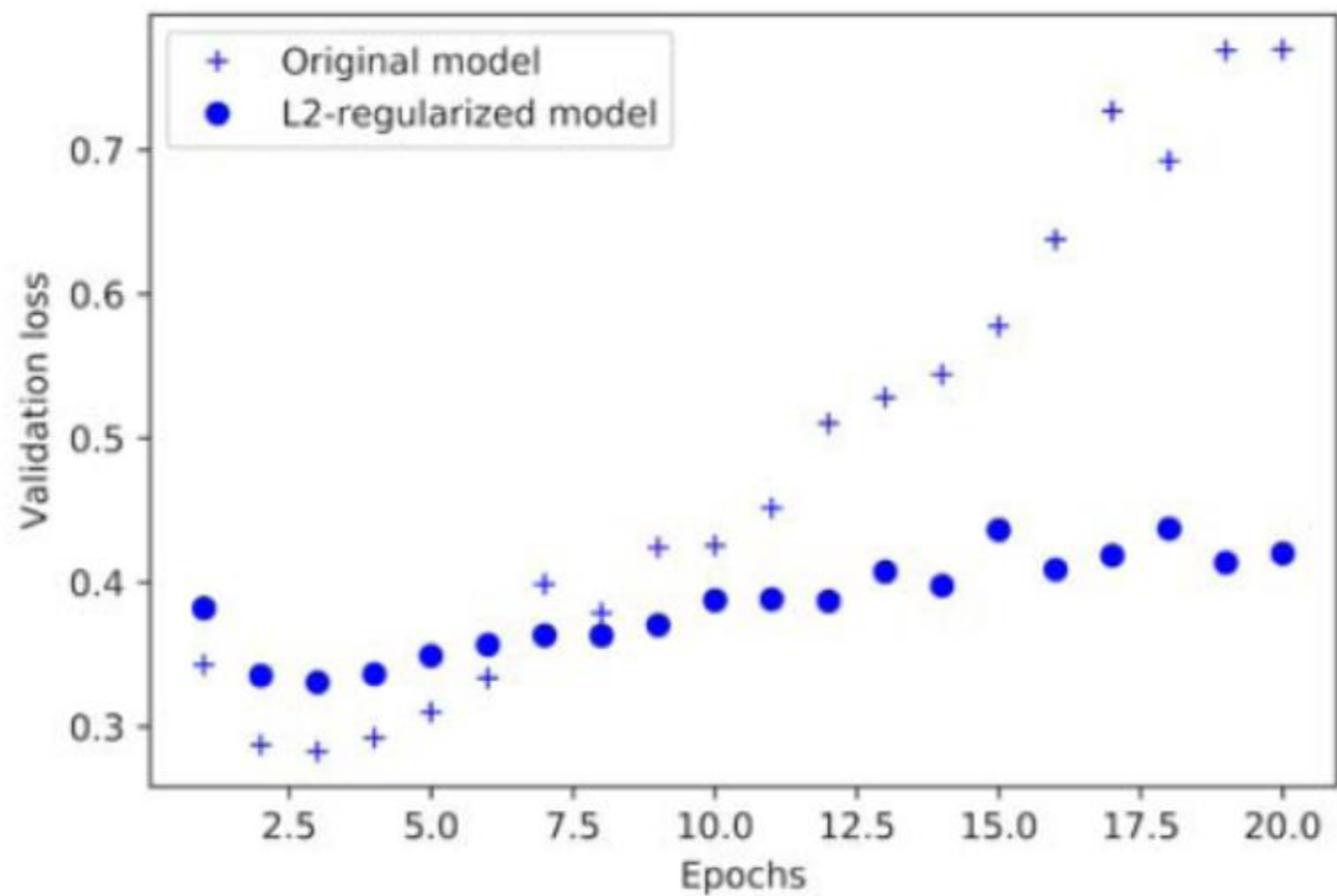
- *L1 regularization*—The cost added is proportional to the *absolute value of the weight coefficients* (the *L1 norm* of the weights).
- *L2 regularization*—The cost added is proportional to the *square of the value of the weight coefficients* (the *L2 norm* of the weights). L2 regularization is also called *weight decay* in the context of neural networks. Don't let the different name confuse you: weight decay is mathematically the same as L2 regularization.

## Listing 4.6 Adding L2 weight regularization to the model

```
from keras import regularizers

model = models.Sequential()
model.add(layers.Dense(16, kernel_regularizer=regularizers.l2(0.001),
                       activation='relu', input_shape=(10000,)))
model.add(layers.Dense(16, kernel_regularizer=regularizers.l2(0.001),
                       activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))
```

`l2(0.001)` means every coefficient in the weight matrix of the layer will add  $0.001 * \text{weight\_coefficient\_value}$  to the total loss of the network. Note that because this penalty is *only added at training time*, the loss for this network will be much higher at training than at test time.



**Figure 4.7** Effect of L2 weight regularization on validation loss

## Listing 4.7 Different weight regularizers available in Keras

```
from keras import regularizers
```

```
regularizers.l1(0.001) ← L1 regularization
```

```
regularizers.l1_l2(l1=0.001, l2=0.001)
```

**Simultaneous L1 and  
L2 regularization**

The processing of fighting overfitting this way is called *regularization*. Let's review some of the most common regularization techniques and apply them in practice to improve the movie-classification model from section 3.4.

*Reducing the network's size*

*Adding weight regularization*

*Adding dropout*

*Dropout* is one of the most effective and most commonly used regularization techniques for neural networks, developed by Geoff Hinton and his students at the University of Toronto. Dropout, applied to a layer, consists of randomly *dropping out* (setting to zero) a number of output features of the layer during training. Let's say a given layer would normally return a vector  $[0.2, 0.5, 1.3, 0.8, 1.1]$  for a given input sample during training. After applying dropout, this vector will have a few zero entries distributed at random: for example,  $[0, 0.5, 1.3, 0, 1.1]$ . The *dropout rate* is the fraction of the features that are zeroed out; it's usually set between 0.2 and 0.5. At test time, no units are dropped out; instead, the layer's output values are scaled down by a factor equal to the dropout rate, to balance for the fact that more units are active than at training time.



```

layer_output *= np.random.randint(0, high=2, size=layer_output.shape)
layer_output /= 0.5

```

Note that we're scaling up rather scaling down in this case.

At training time

0.3	0.2	1.5	0.0
0.6	0.1	0.0	0.3
0.2	1.9	0.3	1.2
0.7	0.5	1.0	0.0

50% dropout

0.0	0.2	1.5	0.0
0.6	0.1	0.0	0.3
0.0	1.9	0.3	0.0
0.7	0.0	0.0	0.0

\* 2

**Figure 4.8** Dropout applied to an activation matrix at training time, with rescaling happening during training. At test time, the activation matrix is unchanged.

## Listing 4.8 Adding dropout to the IMDB network

```
model = models.Sequential()  
model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))  
model.add(layers.Dropout(0.5))  
model.add(layers.Dense(16, activation='relu'))  
model.add(layers.Dropout(0.5))  
model.add(layers.Dense(1, activation='sigmoid'))
```

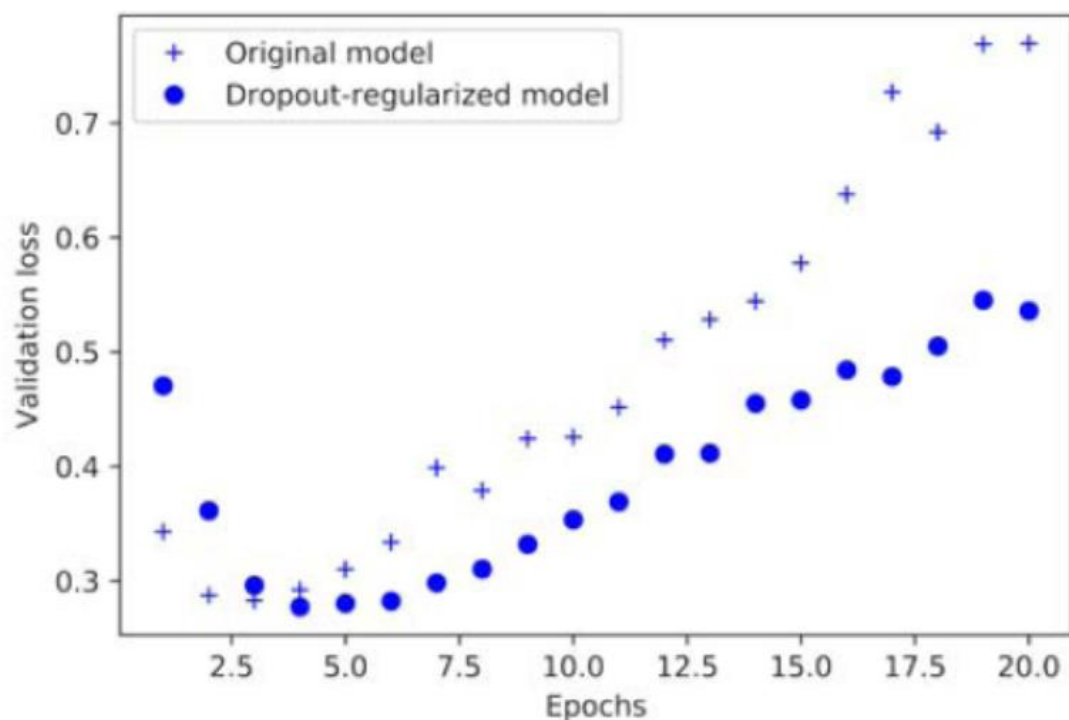


Figure 4.9 Effect of dropout on validation loss

**Table 4.1** Choosing the right last-layer activation and loss function for your model

Problem type	Last-layer activation	Loss function
Binary classification	<code>sigmoid</code>	<code>binary_crossentropy</code>
Multiclass, single-label classification	<code>softmax</code>	<code>categorical_crossentropy</code>
Multiclass, multilabel classification	<code>sigmoid</code>	<code>binary_crossentropy</code>
Regression to arbitrary values	None	<code>mse</code>
Regression to values between 0 and 1	<code>sigmoid</code>	<code>mse</code> or <code>binary_crossentropy</code>