

CH2. Mathematical Building Blocks of Deep-learning (3)

Il-Youp Kwak
Chung-Ang University



Mathematical Building Blocs of Deep-learning

A first look on neural network

Data representations for neural network

The gear of neural networks: tensor operations

The engine of neural networks: gradient-based optimization

Looking back at our first example



The engine of neural networks: gradient-based optimization

(Perceptron model)

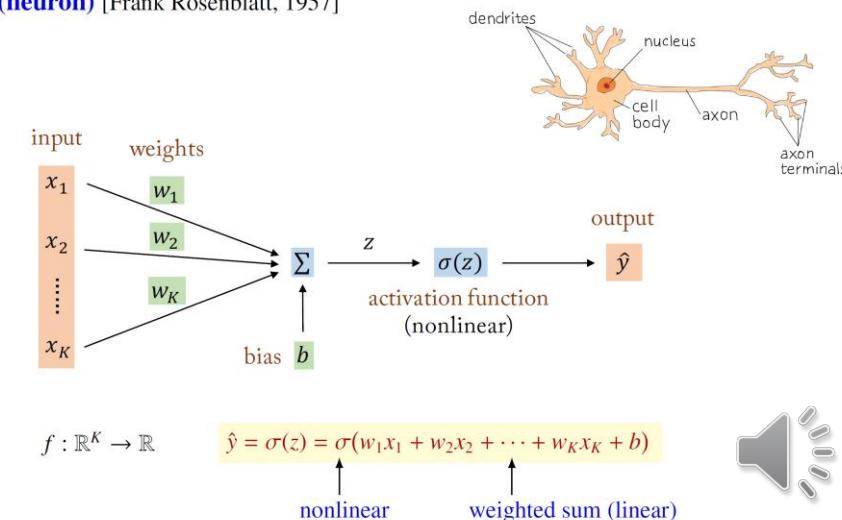
A dense layer:

$$\text{output} = \text{relu}(\text{dot}(W, \text{input}) + b)$$

W and b are **weights** or **trainable parameters**

Gradual adjustment of them is called **training** (\approx searching better and better W&b)

Perceptron (neuron) [Frank Rosenblatt, 1957]



whole data → 100 batches



→ train → loss / cost → optimize (Update W&b for this batch)

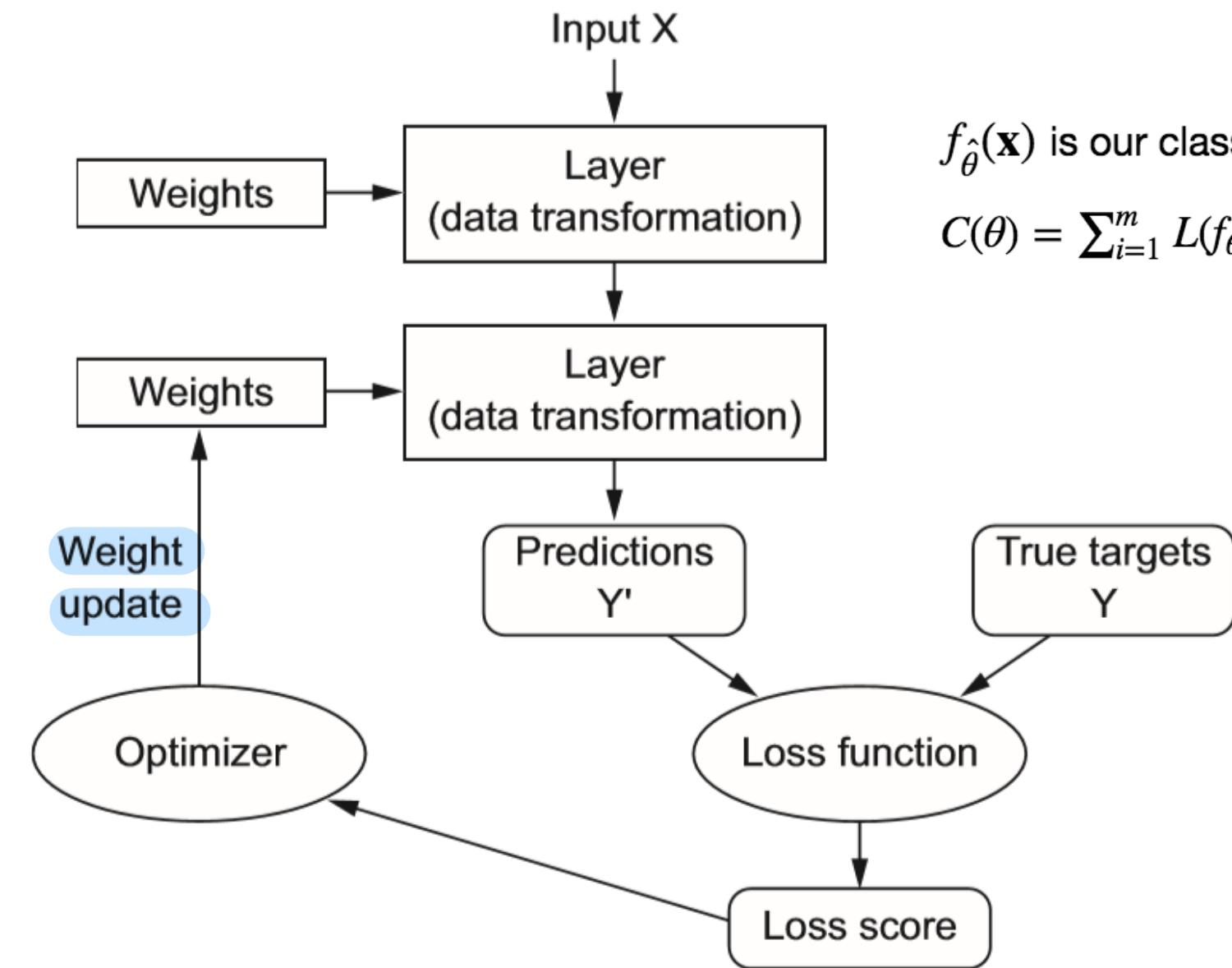
Training loop

What we learned

- Draw a batch of training samples x and corresponding targets y .
- Run the network on x (a step called the forward pass) to obtain predictions y_{pred} .
- Compute the loss of the network on the batch, a measure of the mismatch between y_{pred} and y . $cost$
- Update all weights of the network in a way that slightly reduces the loss on this batch.

What we will learn this chapter





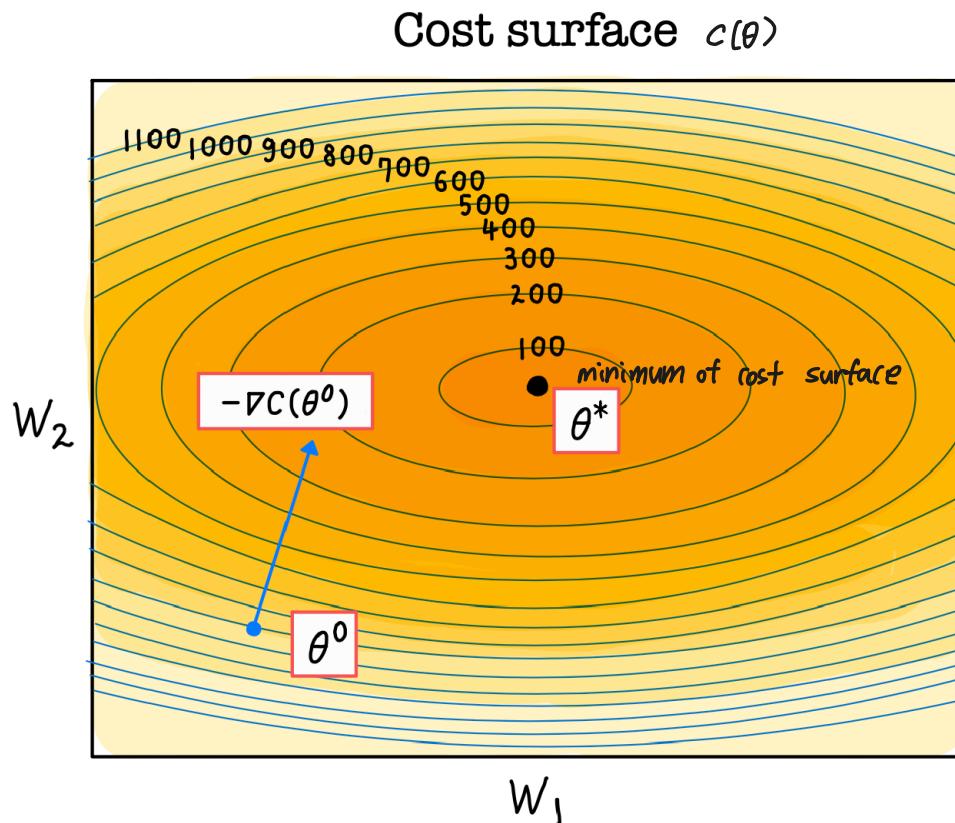
$\hat{f}_\theta(\mathbf{x})$ is our classification model with estimated parameters $\hat{\theta}$
 $C(\theta) = \sum_{i=1}^m L(f_\theta(\mathbf{x}_i), y_i)$

Figure 1.9 The loss score is used as a feedback signal to adjust the weights.



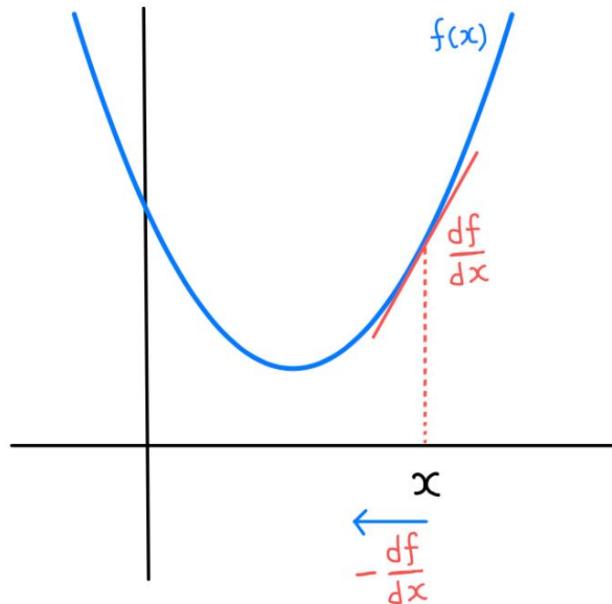
Gradient Descent

- Gradient descent is a first-order (requiring first-derivative/gradient) iterative optimization algorithm for finding the minimum of a function.**

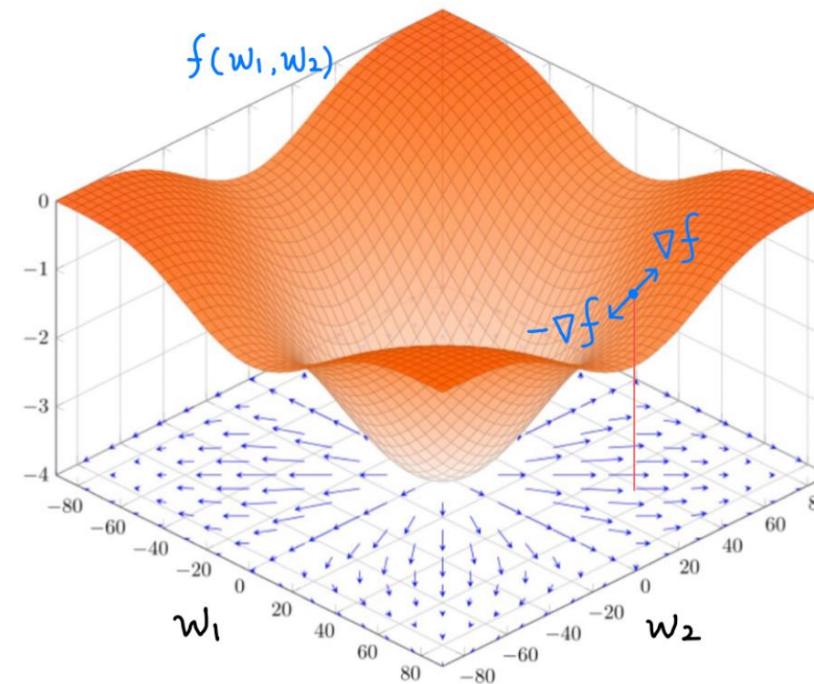


network with two parameters $\theta = \{w_1, w_2\}$

Derivative

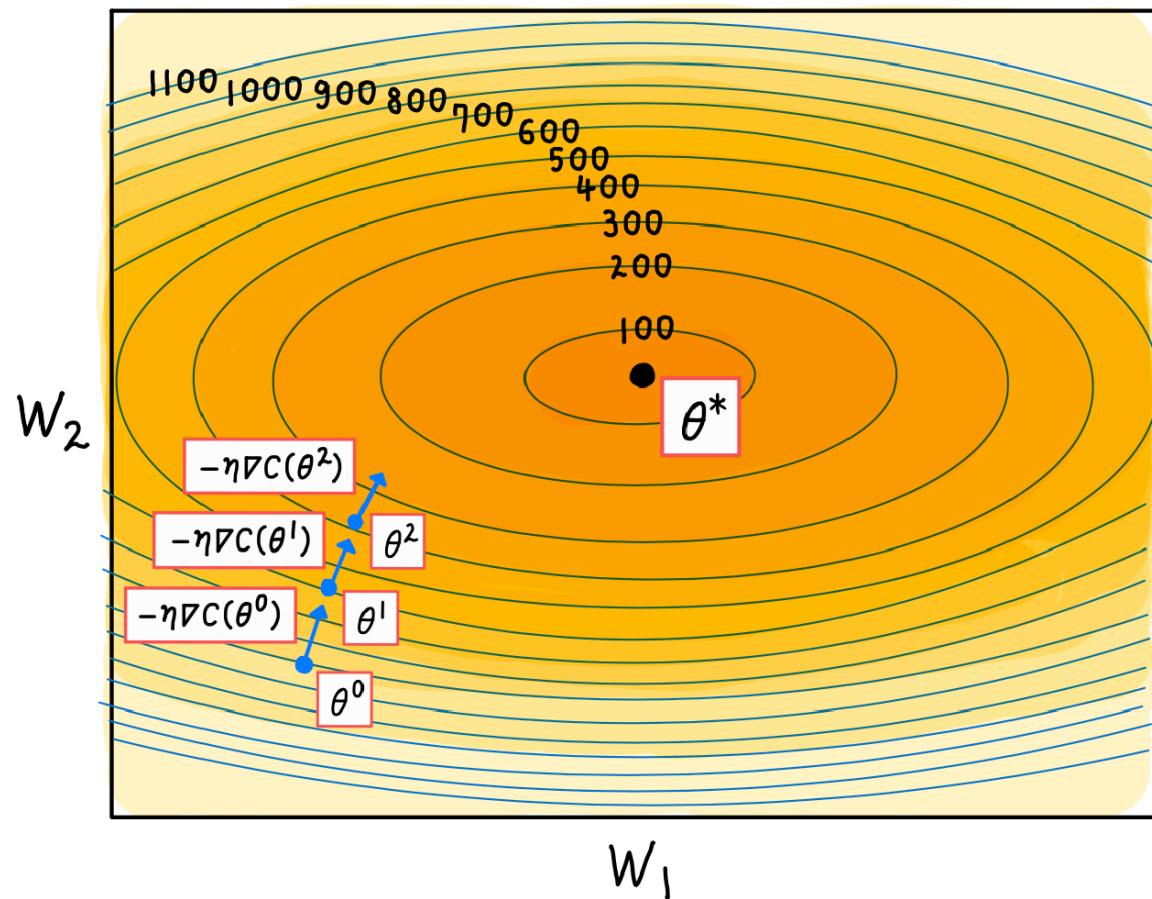


Gradient



* Gradient vector $\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$ indicates the direction and rate of fastest increase.

Cost surface



Randomly pick a starting point θ^0 .

Compute the negative gradient at each θ^t
followed by multiplying η .

$$\Rightarrow -\eta \nabla_{\theta} C(\theta^t)$$

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} C(\theta^t)$$

Ex) Binary Classification Model

- We are Classifying Cat or Dog $f : \mathbf{x} \xrightarrow{f_\theta} \mathbb{R}_{[0,1]}$ $\begin{cases} 0: \text{cat} \\ 1: \text{dog} \end{cases}$
- Define a Cost function $C(\theta; X, y) = \frac{m}{2} \sum_{i=1}^m L(\theta; x_i, y_i)$
- Minimize Cost w.r.t θ given data $\operatorname{argmin}_\theta C(\theta; X, y)$
- Iterate gradient updates $\hat{\theta}^{t+1} = \theta^t - \eta \nabla_\theta C(\theta^t)$ gradient descent algorithm
- $\hat{f}_\theta(\mathbf{x})$ is your classification model



Ex) Binary Classification Model

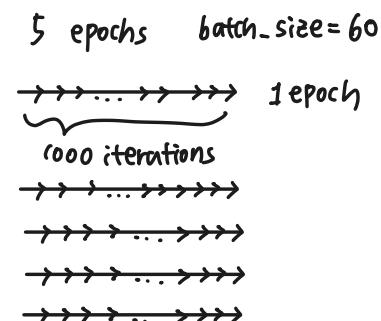
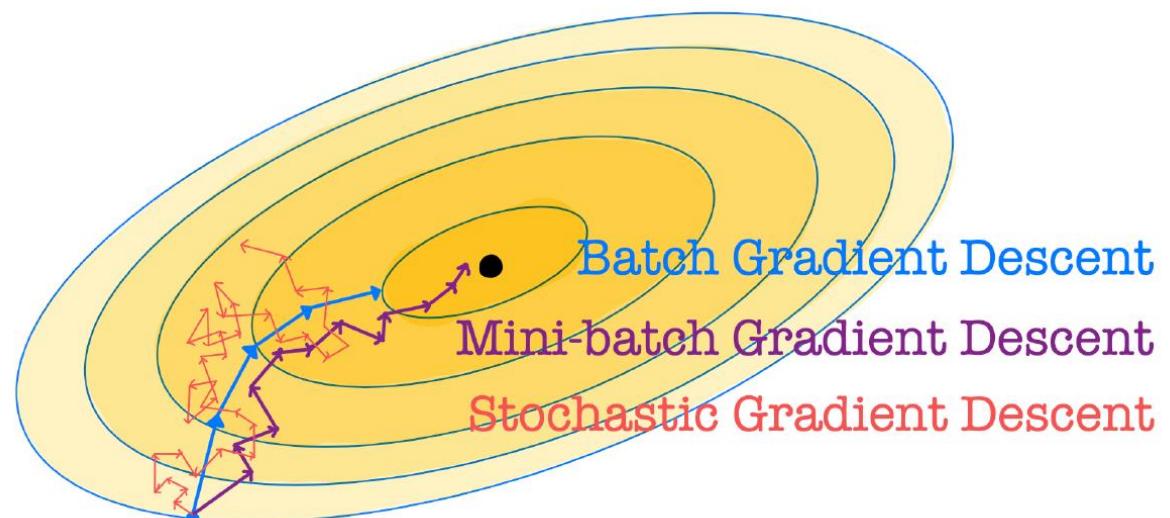
- Minimize Loss w.r.t θ given data $\text{argmin}_{\theta} C(\theta; X, y)$
This part may suffer (**memory problem**)
too many dataset
- We can consider using **mini-batch Gradient Descent**,
Stochastic Gradient Descent
(slower) *↓*
usually using it

Our programming is optimized with vectorized calculation.
→ mini-batch G-D is faster



Gradient descent variants

- (1) **Batch Gradient Descent**: batch size = size of training dataset (eg. 60,000) gradient update using all data iteration = 1
- (2) **Stochastic Gradient Descent (SGD)**: batch size = 1 θ^t is updated by 60,000 times.
- (3) **Minibatch Gradient Descent** eg) batch_size= 60 \rightarrow iteration= 1,000
 ↳ memory problem? we can control w/ batch size.
- * **Batch size**: number of samples processed before the model parameters are updated.
 - * Number of **epochs**: number of complete passes through the entire training dataset.



Variants of SGD

SGD with momentum as well as Adagrad, RMSProp, and several others.

Known as ***optimization methods*** or ***optimizers***.

```
network.compile(optimizer='rmsprop',  
                loss='categorical_crossentropy',  
                metrics=['accuracy'])
```



Chaining derivatives: the Backpropagation algorithm

chain rule: $f(g(x)) = f'(g(x)) * g'(x)$.

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} C(\theta^t)$$

How calculate?

$$C(\theta) = \sum_{i=1}^m L(\theta) = \sum_{i=1}^m L(f_\theta(x_i), y_i) \rightarrow \text{Need to derivative.}$$

D. Rumelhart, G. Hinton, and R. Williams (1986) proposed **backpropagation**, basically, it is gradient decent with chain rule

$$f_\theta(x) = f_3(f_2(f_1(x))), \quad \begin{array}{l} \text{memorize the results (values)} \\ \text{and using for calculating derivatives (Backpropagation algorithm)} \end{array}$$

$$f_\theta'(x) = f_3'(f_2(f_1(x)))(f_2'(f_1(x)))f_1'(x)$$



Backpropagation algorithm

$$y = w^*x + b$$

estimate parameter

$$\text{loss_val} = |y_{\text{true}} - y|$$

$$\text{grad}(\text{loss_val}, w) ? \text{grad}(\text{loss_val}, b) ?$$

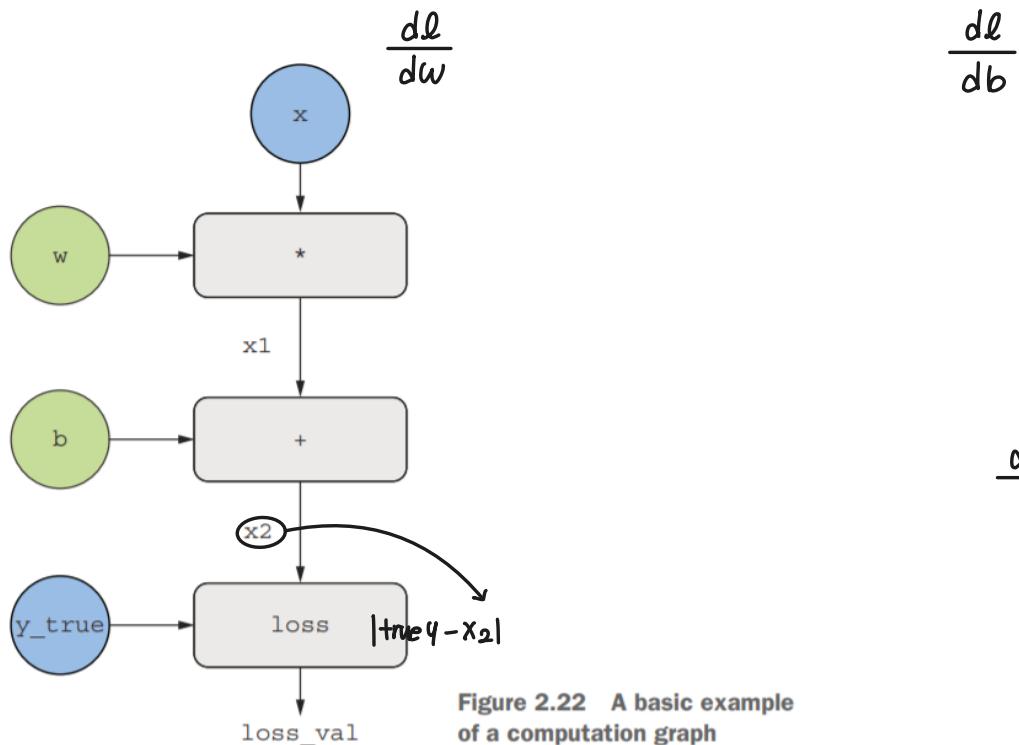


Figure 2.22 A basic example of a computation graph

$$\begin{aligned} \frac{dx_1}{dw} &= \frac{d w \cdot x}{dw} = x = 2 \\ \text{grad}(x_1, w) &= 2 \end{aligned}$$

$$\frac{d x_1 + b}{db} = 1$$

$$\text{grad}(x_2, b) = 1$$

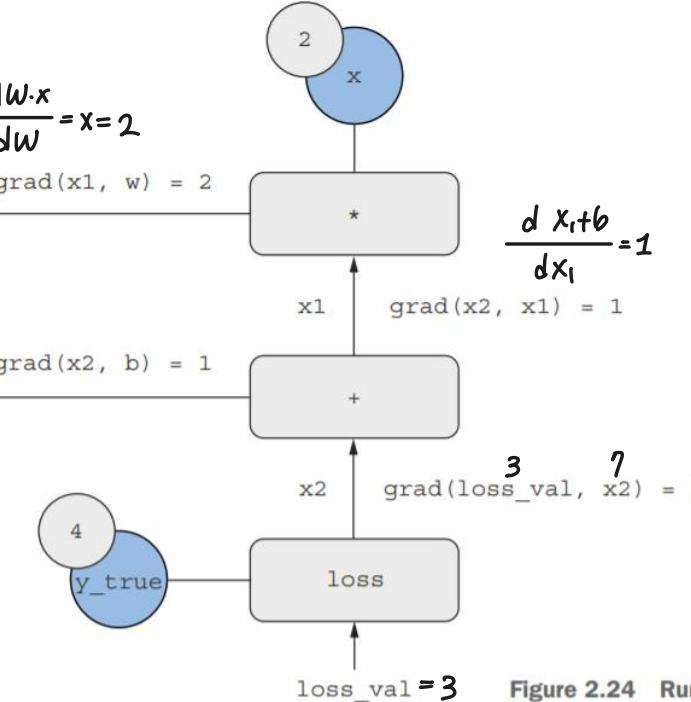


Figure 2.23 Running a forward pass

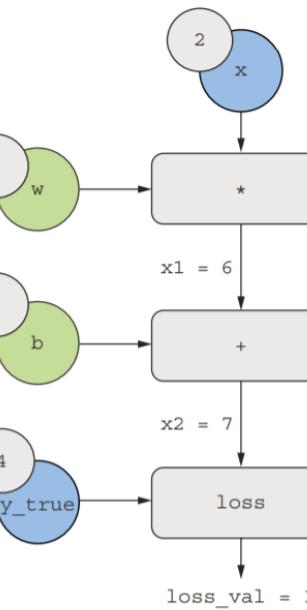
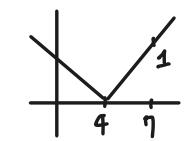


Figure 2.24 Running a backward pass



Backpropagation algorithm

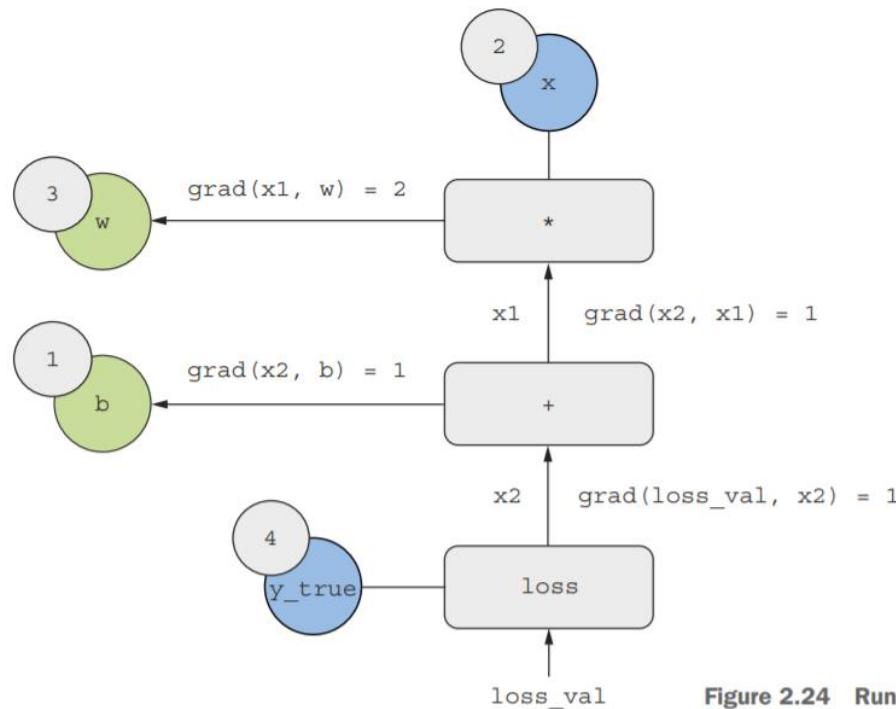


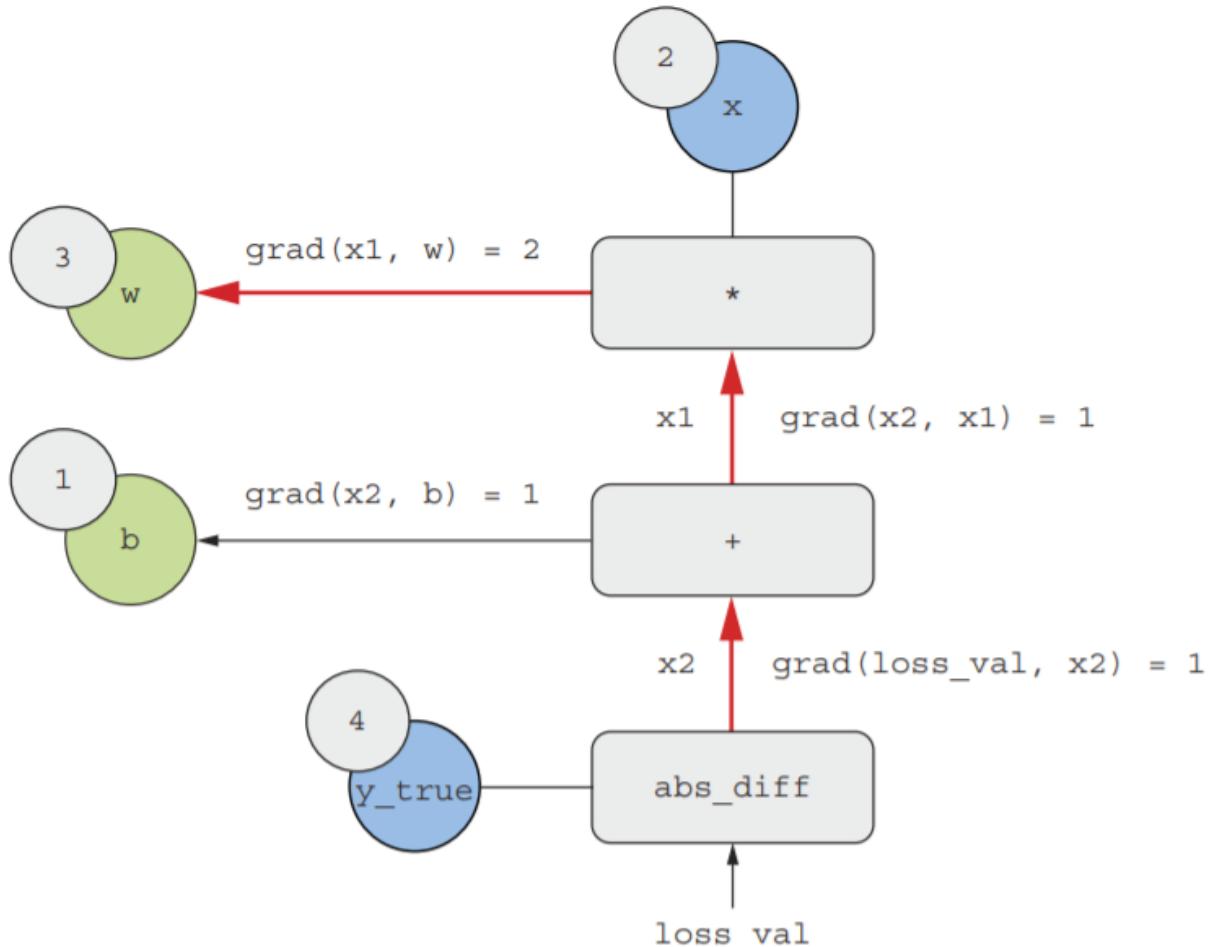
Figure 2.24 Running a backward pass

$$\text{grad}(\text{loss_val}, \text{w}) = 1 * 1 * 2 = 2$$

$$\frac{dl}{dw} = \frac{dl}{dx_2} \frac{dx_2}{dx_1} \frac{dx_1}{dw} = 2$$

$$\text{grad}(\text{loss_val}, \text{b}) = 1 * 1 = 1$$

$$\frac{dl}{db} = \frac{dl}{dx_2} \frac{dx_2}{db} = 1$$

Figure 2.25 Path from **loss_val** to **w** in the backward graph

Gradient tape in Tensorflow

Instantiate a scalar Variable with an initial value of 0.

```
import tensorflow as tf
x = tf.Variable(0.) ←
with tf.GradientTape() as tape: ←
    y = 2 * x + 3 ←
grad_of_y_wrt_x = tape.gradient(y, x) ←
     $\frac{dy}{dx}$ 
```

Open a GradientTape scope.

Inside the scope, apply some tensor operations to our variable.

Use the tape to retrieve the gradient of the output y with respect to our variable x.

The GradientTape works with tensor operations:

matrix 2D tensor

```
x = tf.Variable(tf.random.uniform((2, 2))) ←
with tf.GradientTape() as tape: ←
    y = 2 * x + 3 ←
grad_of_y_wrt_x = tape.gradient(y, x) ←
```

Instantiate a Variable with shape (2, 2) and an initial value of all zeros.

grad_of_y_wrt_x is a tensor of shape (2, 2) (like x) describing the curvature of $y = 2 * a + 3$ around $x = [[0, 0], [0, 0]]$.



Gradient tape in Tensorflow

It also works with lists of variables:

```
W = tf.Variable(tf.random.uniform((2, 2)))
b = tf.Variable(tf.zeros((2,)))
x = tf.random.uniform((2, 2))
with tf.GradientTape() as tape:
    y = tf.matmul(x, W) + b
grad_of_y_wrt_W_and_b = tape.gradient(y, [W, b])
```

matmul is how you say
“dot product” in TensorFlow.

grad_of_y_wrt_W_and_b is a
list of two tensors with the same
shapes as W and b, respectively.



Looking back at our first example

This was our input data:

```
(train_images, train_labels), (test_images, test_labels) = mnist.load_data()  
train_images = train_images.reshape((60000, 28 * 28))  
train_images = train_images.astype('float32') / 255  
test_images = test_images.reshape((10000, 28 * 28))  
test_images = test_images.astype('float32') / 255
```



This was our network:

```
model = keras.Sequential([
    layers.Dense(512, activation="relu"), → one hidden layer
    layers.Dense(10, activation="softmax") → one output layer
])
```

} both dense layer

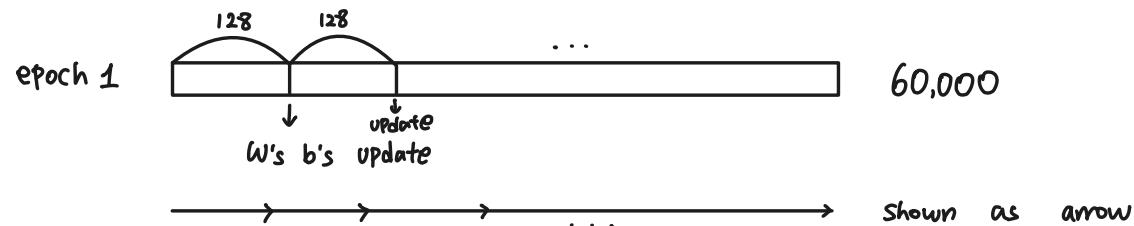
This was our network-compilation step:

```
model.compile(optimizer="rmsprop",
              loss="sparse_categorical_crossentropy",
              metrics=["accuracy"])
```



This was the training loop:

```
model.fit(train_images, train_labels, epochs=5, batch_size=128)
```



Reimplementing our first example from scratch in TensorFlow

```
model = keras.Sequential([
    layers.Dense(512, activation="relu"),
    layers.Dense(10, activation="softmax")
])
```

Dense class

$wx + b$
 $q \times p \times p \times q$

in programming xw
 $p \times p \times q$

```
import tensorflow as tf

class NaiveDense:
    def __init__(self, input_size, output_size, activation):
        self.activation = activation

        w_shape = (input_size, output_size)
        w_initial_value = tf.random.uniform(w_shape, minval=0, maxval=1e-1)
        self.W = tf.Variable(w_initial_value)

        b_shape = (output_size,)
        b_initial_value = tf.zeros(b_shape)
        self.b = tf.Variable(b_initial_value)

    def __call__(self, inputs):
        return self.activation(tf.matmul(inputs, self.W) + self.b)

    @property
    def weights(self):
        return [self.W, self.b]
```

Create a matrix, W , of shape $(input_size, output_size)$, initialized with random values.

Create a vector, b , of shape $(output_size)$, initialized with zeros.

Apply the forward pass.

Convenience method for retrieving the layer's weights



Sequential class

get the inputs as layers

```

class NaiveSequential:
    def __init__(self, layers):
        self.layers = layers

    def __call__(self, inputs):
        x = inputs
        for layer in self.layers: )  $\ell_2(\ell_1(x))$ 
            x = layer(x)
        return x

    @property
    def weights(self):
        weights = []
        for layer in self.layers:
            weights += layer.weights
        return weights

```

```

model = keras.Sequential([
    layers.Dense(512, activation="relu"),
    layers.Dense(10, activation="softmax")
])

```

```

model = NaiveSequential([
    NaiveDense(input_size=28 * 28, output_size=512, activation=tf.nn.relu),
    NaiveDense(input_size=512, output_size=10, activation=tf.nn.softmax)
])
assert len(model.weights) == 4

```



Batch Generator

```

import math

class BatchGenerator:
    def __init__(self, images, labels, batch_size=128):
        assert len(images) == len(labels)
        self.index = 0
        self.images = images
        self.labels = labels
        self.batch_size = batch_size
        self.num_batches = math.ceil(len(images) / batch_size)

    def next(self):
        images = self.images[self.index : self.index + self.batch_size]
        labels = self.labels[self.index : self.index + self.batch_size]
        self.index += self.batch_size
        return images, labels
    ) return batch samples
    ~~~~~ batch size number of objects

```

Run the “forward pass” (compute the model’s predictions under a GradientTape scope).

```

def one_training_step(model, images_batch, labels_batch):
    with tf.GradientTape() as tape:
        predictions = model(images_batch)
        per_sample_losses = tf.keras.losses.sparse_categorical_crossentropy(
            labels_batch, predictions)
        average_loss = tf.reduce_mean(per_sample_losses)
        gradients = tape.gradient(average_loss, model.weights)
        update_weights(gradients, model.weights)
    return average_loss

```

Update the weights using the gradients (we will define this function shortly).

batch size 10
↓
calculate 10 losses

Compute the gradient of the loss with regard to the weights. The output gradients is a list where each entry corresponds to a weight from the model.weights list.



Naive mini-batch gradient update

```
learning_rate = 1e-3

def update_weights(gradients, weights):
    for g, w in zip(gradients, weights):
        w.assign_sub(g * learning_rate)
```

assign_sub is the equivalent of -= for TensorFlow variables.

Using optimizer

```
from tensorflow.keras import optimizers

optimizer = optimizers.SGD(learning_rate=1e-3)

def update_weights(gradients, weights):
    optimizer.apply_gradients(zip(gradients, weights))
```



Full training loop

```
def fit(model, images, labels, epochs, batch_size=128):
    for epoch_counter in range(epochs):
        print(f"Epoch {epoch_counter}")
    batch_generator = BatchGenerator(images, labels)
    for batch_counter in range(batch_generator.num_batches):
        images_batch, labels_batch = batch_generator.next()
        loss = one_training_step(model, images_batch, labels_batch)
        if batch_counter % 100 == 0:
            print(f"loss at batch {batch_counter}: {loss:.2f}")
```

```
from tensorflow.keras.datasets import mnist
(train_images, train_labels), (test_images, test_labels) = mnist.load_data()

train_images = train_images.reshape((60000, 28 * 28))
train_images = train_images.astype("float32") / 255
test_images = test_images.reshape((10000, 28 * 28))
test_images = test_images.astype("float32") / 255

fit(model, train_images, train_labels, epochs=10, batch_size=128)
```



chapter02_mathematical-building-blocks.

File Edit View Insert Runtime Tools Help Cannot save changes

Share  

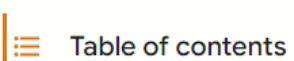


Table of contents

Derivative of a tensor operation: the gradient

Stochastic gradient descent

Chaining derivatives: The Backpropagation algorithm

The chain rule

Automatic differentiation with computation graphs

The gradient tape in TensorFlow

Looking back at our first example

Reimplementing our first example from scratch in TensorFlow

A simple Dense class

A simple Sequential class

A batch generator

Running one training step

The full training loop

Evaluating the model

Summary

+ Section

+ Code + Text | Copy to Drive

60000/12

468, 75

- ▼ Reimplementing our first example from scratch in TensorFlow

▼ A simple Dense class

```
[ ] import tensorflow as tf

class NaiveDense:
    def __init__(self, input_size, output_size, activation):
        self.activation = activation

        w_shape = (input_size, output_size)
        w_initial_value = tf.random.uniform(w_shape, minval=0, maxval=1e-1)
        self.W = tf.Variable(w_initial_value)

        b_shape = (output_size,)
        b_initial_value = tf.zeros(b_shape)
        self.b = tf.Variable(b_initial_value)

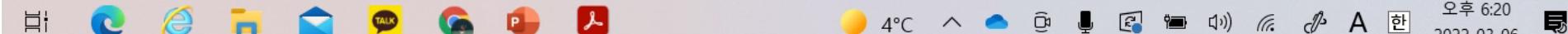
    def __call__(self, inputs):
        return self.activation(tf.matmul(inputs, self.W) + self.b)

    @property
    def weights(self):
        return [self.W, self.b]
```

✓ 9s completed at 6:01 PM



🔍 검색하려면 여기에 입력하십시오



Thank you!



CH3. Introduction to Keras and TensorFlow

Il-Youp Kwak
Chung-Ang University



Introduction to Keras and TensorFlow

What's Tensorflow

What's Keras

Keras and Tensorflow: A brief history

First steps with Tensorflow

Anatomy of a neural network: Understanding Keras APIs



What's TensorFlow?

TensorFlow is a Python-based, free, open source machine learning platform, developed by Google.

* Core Part

Compute the gradient of any differentiable expression

In ML,
one core part is estimating parameters
w/ optimization steps.

Run not only on **CPUs**, but also on **GPUs and TPUs**, parallel hardware accelerators.

Computation defined in TensorFlow can be easily distributed

TensorFlow programs can be exported to other runtimes (**C++**, **Java**)

TensorFlow is much more than a single library. It's really a platform, home to a vast ecosystem of components,



What's Keras?

Keras is a **deep learning API** for Python, built on top of **TensorFlow**

Provides a **convenient way to define** and train deep learning model.

Developed for research with the aim of **enabling fast experimentation**



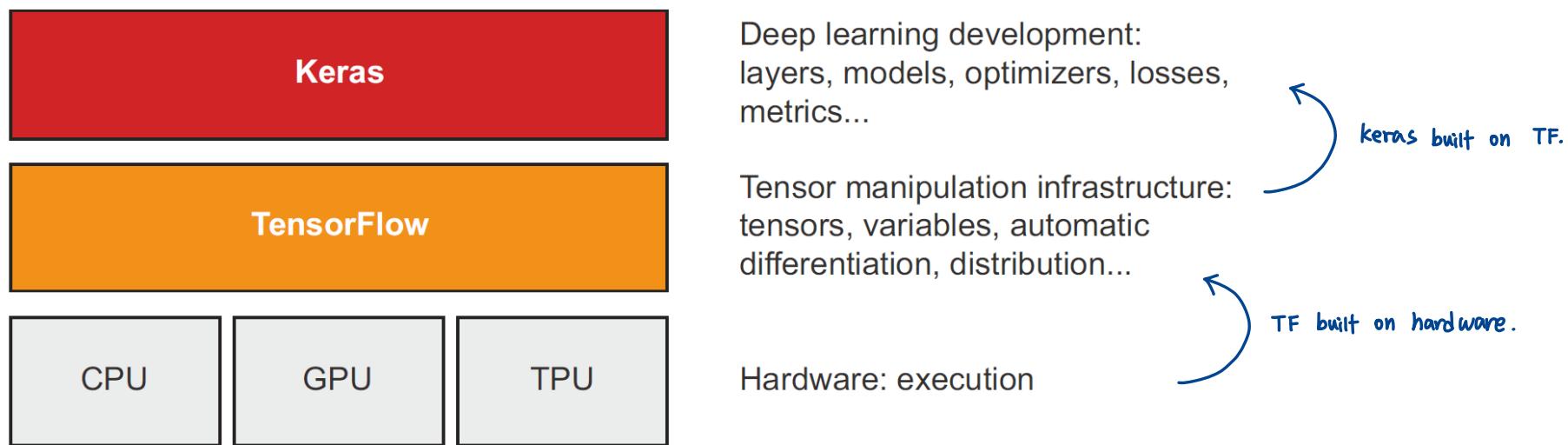


Figure 3.1 Keras and TensorFlow: TensorFlow is a low-level tensor computing platform, and Keras is a high-level deep learning API

- **Allows the same code to run seamlessly on CPU or GPU**
- **Has a user-friendly API that makes it easy to quickly prototype deep-learning models.**
- **Has built-in support for convolutional networks, recurrent networks, and any combination of both.**

MIT license, it can be freely used in commercial projects.



Keras and TensorFlow: A brief history

Keras predates TensorFlow by eight months. It was released in March 2015, and TensorFlow was released in November 2015.

Keras was originally built on top of Theano, another tensor-manipulation library that provided automatic differentiation and GPU support

In late 2015, after the release of TensorFlow, Keras was refactored to a multibackend architecture: Use Keras with either Theano or TensorFlow

* 딥러닝 프레임워크의 대표적 기능 : 자동미분, GPU의 활용

In 2017, two new additional backend options were added to Keras: CNTK and MXNet

마이크로소프트 딥러닝프레임워크 another software

In 2018, the TensorFlow leadership picked Keras as TensorFlow's official high-level API. As a result, the Keras API is front and center in TensorFlow 2.0, released in September 2019



First steps with TensorFlow

Training a neural network revolves around the following concepts:

- 1. low-level tensor manipulation: tensors, tensor operations, backpropagation**
- 2. high-level deep learning concepts: layers, loss function, metrics, training loop**



Constant tensors and variables

To do anything in TensorFlow, we're going to need some tensors

Defining constant tensors

Listing 3.1 All-ones or all-zeros tensors

```
>>> import tensorflow as tf
>>> x = tf.ones(shape=(2, 1))
>>> print(x)
tf.Tensor(
[[1.]
 [1.]], shape=(2, 1), dtype=float32)
>>> x = tf.zeros(shape=(2, 1))
>>> print(x)
tf.Tensor(
[[0.]
 [0.]], shape=(2, 1), dtype=float32)
```

Equivalent to `np.ones(shape=(2, 1))`

Equivalent to `np.zeros(shape=(2, 1))`

Defining random tensors

Listing 3.2 Random tensors

```
>>> x = tf.random.normal(shape=(3, 1), mean=0., stddev=1.)
>>> print(x)
tf.Tensor(
[[-0.14208166]
 [-0.95319825]
 [ 1.1096532 ]], shape=(3, 1), dtype=float32)
>>> x = tf.random.uniform(shape=(3, 1), minval=0., maxval=1.)
>>> print(x)
tf.Tensor(
```

Tensor of random values drawn from a normal distribution with mean 0 and standard deviation 1. Equivalent to `np.random.normal(size=(3, 1), loc=0., scale=1.)`.

Tensor of random values drawn from a uniform distribution between 0 and 1. Equivalent to `np.random.uniform(size=(3, 1), low=0., high=1.)`.



A significant difference between NumPy arrays and TensorFlow tensors is that TensorFlow tensors aren't assignable

Listing 3.3 NumPy arrays are assignable

```
import numpy as np  
x = np.ones(shape=(2, 2))  
x[0, 0] = 0.
```

Try to do the same thing in TensorFlow, and you will get an error: “EagerTensor object does not support item assignment.”

* Tensor 와 NP.array 은 차이점은 , Tensor는 값을 할당할 수 없다.

Listing 3.4 TensorFlow tensors are not assignable

```
x = tf.ones(shape=(2, 2))  
x[0, 0] = 0.
```

This will fail, as a tensor isn't assignable.



To train a model, we'll need to update its state. That's where variables come in. `tf.Variable` is the class meant to manage modifiable state in TensorFlow

Listing 3.5 Creating a TensorFlow variable

```
>>> v = tf.Variable(initial_value=tf.random.normal(shape=(3, 1)))
>>> print(v)
array([[-0.75133973],
       [-0.4872893 ],
       [ 1.6626885 ]], dtype=float32)>
```

State of a variable can be modified via its assign method

Listing 3.6 Assigning a value to a TensorFlow variable

```
>>> v.assign(tf.ones((3, 1)))
array([[1.],
       [1.],
       [1.]], dtype=float32)>
```

Listing 3.7 Assigning a value to a subset of a TensorFlow variable

```
>>> v[0, 0].assign(3.)
array([[3.],
       [1.],
       [1.]], dtype=float32)>
```



assign_add() and assign_sub() are efficient equivalents of += and -=**Listing 3.8 Using assign_add()**

```
v = []
>>> v.assign_add(tf.ones((3, 1)))
array([[2.],
       [2.],
       [2.]], dtype=float32) >
```

$a += 1$
 $a = a + 1$



Tensor operations: Doing math in TensorFlow

Just like NumPy, TensorFlow offers a large collection of tensor operations to express mathematical formulas.

Listing 3.9 A few basic math operations

```
a = tf.ones((2, 2))  
b = tf.square(a)           Take the square.  
c = tf.sqrt(a)            Take the square root.  
d = b + c                Add two tensors (element-wise).  
e = tf.matmul(a, b)       Take the product of two tensors  
e *= d                   (as discussed in chapter 2).  
  
Multiply two tensors      (element-wise).
```

Importantly, each of the preceding operations gets executed on the fly: at any point, you can print what the current result is, just like in NumPy. We call this **eager execution**. (in previous version of TF, it doesn't work)



A second look at the GradientTape API

TensorFlow seems to look a lot like NumPy. But here's something NumPy can't do:

Listing 3.10 Using the GradientTape

```
input_var = tf.Variable(initial_value=3.)  
with tf.GradientTape() as tape:  
    result = tf.square(input_var)  
gradient = tape.gradient(result, input_var)  $\frac{dy}{dx}$ 
```

Retrieve the gradients of the loss of a model with respect to its weights:

gradients = tape.gradient(loss, weights) $\frac{dl}{dw}$



Only trainable variables are tracked by default. With a constant tensor, you'd have to manually mark it as being tracked by calling `tape.watch()` on it.

Listing 3.11 Using GradientTape with constant tensor inputs

```
input_const = tf.constant(3.)
with tf.GradientTape() as tape:
    tape.watch(input_const)
    result = tf.square(input_const)
gradient = tape.gradient(result, input_const)
```



An end-to-end example: A linear classifier in pure TensorFlow

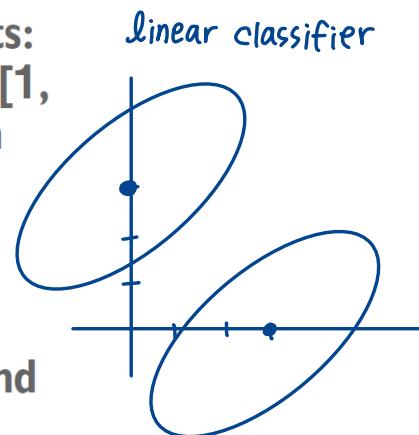
Synthetic data: two classes of points in a 2D plane.

Listing 3.13 Generating two classes of random points in a 2D plane

```
num_samples_per_class = 1000
negative_samples = np.random.multivariate_normal(
    mean=[0, 3],
    cov=[[1, 0.5], [0.5, 1]],
    size=num_samples_per_class)
positive_samples = np.random.multivariate_normal(
    mean=[3, 0],
    cov=[[1, 0.5], [0.5, 1]],
    size=num samples per class)
```

Generate the first class of points:
1000 random 2D points. $\text{cov}=[[1, 0.5], [0.5, 1]]$ corresponds to an oval-like point cloud oriented from bottom left to top right.

Generate the other class of points with a different mean and the same covariance matrix.



Listing 3.14 Stacking the two classes into an array with shape (2000, 2)

we make inputs data set

```
inputs = np.vstack((negative_samples, positive_samples)).astype(np.float32)
          vertically
          stacking
```



Listing 3.15 Generating the corresponding targets (0 and 1)

```
targets = np.vstack((np.zeros((num_samples_per_class, 1), dtype="float32"),
                     np.ones((num_samples_per_class, 1), dtype="float32")))
```

Listing 3.16 Plotting the two point classes (see figure 3.6)

```
import matplotlib.pyplot as plt
plt.scatter(inputs[:, 0], inputs[:, 1], c=targets[:, 0])
plt.show()
```

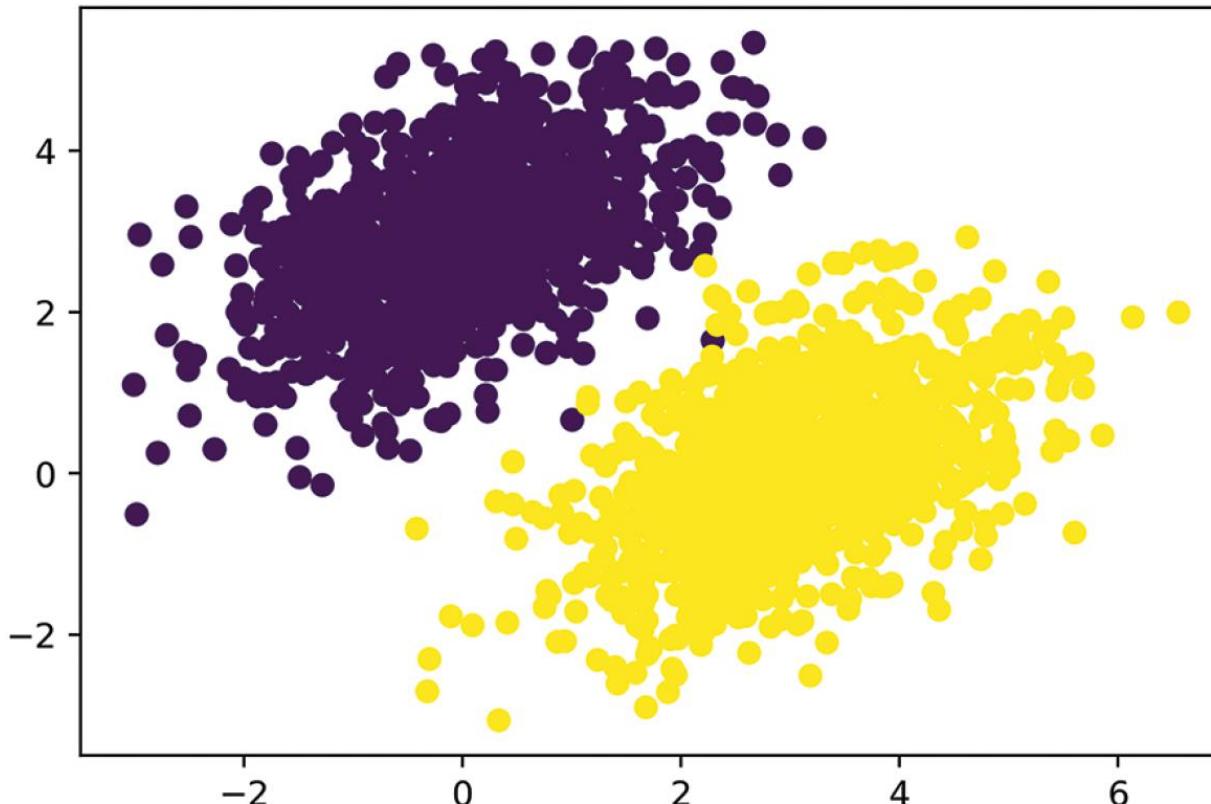


Figure 3.6 Our synthetic data: two classes of random points in the 2D plane



Create a linear classifier

Listing 3.17 Creating the linear classifier variables

The inputs will
be 2D points.

```
input_dim = 2
output_dim = 1
W = tf.Variable(initial_value=tf.random.uniform(shape=(input_dim, output_dim)))
b = tf.Variable(initial_value=tf.zeros(shape=(output_dim,)))
```

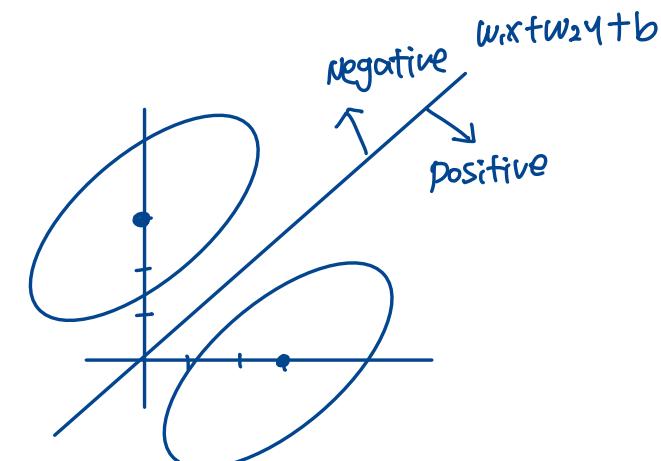
The output predictions will be a single score per sample (close to 0 if the sample is predicted to be in class 0, and close to 1 if the sample is predicted to be in class 1).

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} [x, y] + b$$

Listing 3.18 The forward pass function

```
def model(inputs):
    return tf.matmul(inputs, W) + b
```

prediction = $w_1 * x + w_2 * y + b$.



Define our loss function

Listing 3.19 The mean squared error loss function

per_sample_losses will be a tensor with the same shape as targets and predictions, containing per-sample loss scores.

```
def square_loss(targets, predictions):  
    per_sample_losses = tf.square(targets - predictions) ←  
    return tf.reduce_mean(per_sample_losses) ←
```

We need to average these per-sample loss scores into a single scalar loss value: this is what reduce_mean does.



Training step using gradient decent

Listing 3.20 The training step function

```
learning_rate = 0.1

def training_step(inputs, targets):
    with tf.GradientTape() as tape:
        predictions = model(inputs)model
        loss = square_loss(predictions, targets)
        grad_loss_wrt_W, grad_loss_wrt_b = tape.gradient(loss, [W, b])
        W.assign_sub(grad_loss_wrt_W * learning_rate)
        b.assign_sub(grad_loss_wrt_b * learning_rate)
    return loss
```

Retrieve the gradient of the loss with regard to weights.

Forward pass, inside a gradient tape scope

Update the weights.



Training step using gradient decent

Listing 3.20 The training step function

```
learning_rate = 0.1

def training_step(inputs, targets):
    with tf.GradientTape() as tape:
        predictions = model(inputs)
        loss = square_loss(predictions, targets)
    grad_loss_wrt_W, grad_loss_wrt_b = tape.gradient(loss, [W, b])
    W.assign_sub(grad_loss_wrt_W * learning_rate)
    b.assign_sub(grad_loss_wrt_b * learning_rate)
    return loss
```

Retrieve the gradient of the loss with regard to weights.

Forward pass, inside a gradient tape scope

Update the weights.

40 steps of training (40 epochs)

Listing 3.21 The batch training loop

```
for step in range(40):
    loss = training_step(inputs, targets)
    print(f"Loss at step {step}: {loss:.4f}")
```



$\text{prediction} = w_1 * x + w_2 * y + b$. Thus, class 0 is defined as $w_1 * x + w_2 * y + b < 0.5$, and class 1 is defined as $w_1 * x + w_2 * y + b > 0.5$

Rule $w_1 * x + w_2 * y + b = 0.5$ becomes $y = -w_1 / w_2 * x + (0.5 - b) / w_2$

Generate 100 regularly spaced numbers between -1 and 4, which we will use to plot our line.

```
x = np.linspace(-1, 4, 100)
y = -W[0] / W[1] * x + (0.5 - b) / W[1]
plt.plot(x, y, "-r")
plt.scatter(inputs[:, 0], inputs[:, 1], c=predictions[:, 0] > 0.5)
```

This is our line's equation.

Plot our line ("-r" means "plot it as a red line").

Plot our model's predictions on the same plot.

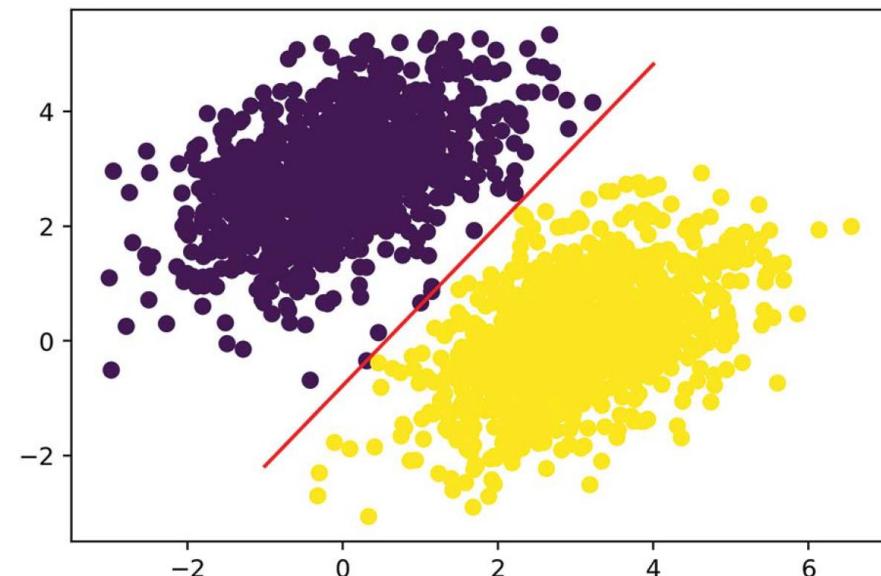


Figure 3.8 Our model, visualized as a line



Anatomy of a neural network: Understanding core Keras APIs

Layers: The building blocks of deep learning

A layer is a data processing module that takes as input tensors and that outputs tensors

Ex) densely connected layer, called **fully connected** or **dense** layer (Dense class in Keras)



Base Layer class in Keras

Layer is object that encapsulates some state (weights) and computation (forward pass)

weights are typically defined in a `build()` (although they could also be created in the constructor, `__init__()`), and the computation is defined in the `call()` method.

Listing 3.22 A Dense layer implemented as a Layer subclass

```
from tensorflow import keras

class SimpleDense(keras.layers.Layer):
    def __init__(self, units, activation=None):
        super().__init__()
        self.units = units # output dimension
        self.activation = activation

    def build(self, input_shape):
        input_dim = input_shape[-1]
        self.W = self.add_weight(shape=(input_dim, self.units),
                               initializer="random_normal")
        self.b = self.add_weight(shape=(self.units,), initializer="zeros")

    def call(self, inputs):
        y = tf.matmul(inputs, self.W) + self.b
        if self.activation is not None:
            y = self.activation(y)
        return y
```

We define the forward pass computation in the `call()` method.

All Keras layers inherit from the base Layer class.

Weight creation takes place in the `build()` method.

$W = [\dots]$
 $b = [x] \times b$ matrix
 $b = 1 \times b$ vector

`add_weight()` is a shortcut method for creating weights. It is also possible to create standalone variables and assign them as layer attributes, like `self.W = tf.Variable(tf.random.uniform(w_shape))`.

The `__call__()` method of the base layer looks like this:

```
def __call__(self, inputs):
    if not self.built:
        self.build(inputs.shape)
        self.built = True
    return self.call(inputs)
```



Base Layer class in Keras

Once instantiated, a layer like this can be used just like a function, taking as input a TensorFlow tensor:

```
>>> my_dense = SimpleDense(units=32, activation=tf.nn.relu)    ↪  
>>> input_tensor = tf.ones(shape=(2, 784))      ↪  
>>> output_tensor = my_dense(input_tensor)       ↪  
>>> print(output_tensor.shape)  
(2, 32)                                     ↪  
units=32
```

Instantiate our layer, defined previously.

Create some test inputs.

Call the layer on the inputs, just like a function.

```
from tensorflow.keras import layers  
layer = layers.Dense(32, activation="relu")
```

A dense layer with 32 output units



Automatic shape inference

layers didn't receive any information about the shape of their inputs—instead, they automatically inferred their input shape as being the shape of the first inputs they see.

```
from tensorflow.keras import models
from tensorflow.keras import layers
model = models.Sequential([
    layers.Dense(32, activation="relu"),
    layers.Dense(32)
])
```

Compare the Dense layer we implemented in chapter 2

```
model = NaiveSequential([
    NaiveDense(input_size=784, output_size=32, activation="relu"),
    NaiveDense(input_size=32, output_size=64, activation="relu"),
    NaiveDense(input_size=64, output_size=32, activation="relu"),
    NaiveDense(input_size=32, output_size=10, activation="softmax")
])
```

The shape is automatically matched.

just defining output shape

```
model = keras.Sequential([
    SimpleDense(32, activation="relu"),
    SimpleDense(64, activation="relu"),
    SimpleDense(32, activation="relu"),
    SimpleDense(10, activation="softmax")
])
```



From layers to models

eg) Transformer

A deep learning model is a graph of layers. So far we've learned **Sequential model**

There are much broader variety of network topologies: Two-branch networks, Multihead networks, Residual connections, etc.

There are generally two ways of building such models in Keras: you could directly subclass the **Model** class, or you could use the **Functional API**, which lets you do more with less code. (Will be covered from chapter 7)

Picking the right network architecture is more an art than a science, and although there are some best practices and principles you can rely on.

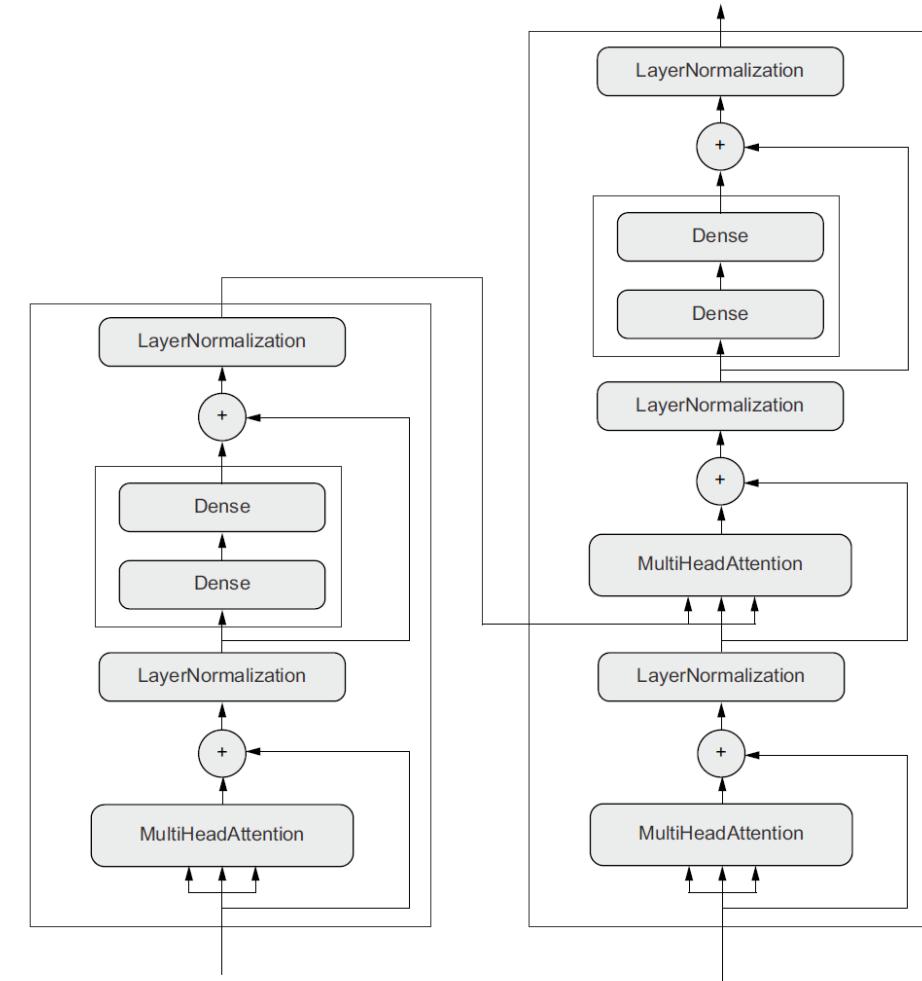


Figure 3.9 The Transformer architecture (covered in chapter 11). There's a lot going on here. Throughout the next few chapters, you'll climb your way up to understanding it.



Picking a loss function

Choosing the **right loss function for the right problem is extremely important**

There are simple guidelines you can follow to choose:

Binary crossentropy for a two-class classification

Categorical crossentropy for a many-class classification problem

Mean squared error for regression problem



Understanding the fit() method train the model

key arguments:

The data (inputs and targets) to train on: NumPy arrays or a TensorFlow Dataset object.

The number of epochs to train for: how many times the training loop iterate over the data passed.

The batch size to use within each epoch of mini-batch gradient descent: the number of training examples considered to compute the gradients for one weight update step.

Listing 3.23 Calling fit() with NumPy data

```
history = model.fit(  
    inputs,  
    targets,  
    epochs=5,  
    batch_size=128  
)  
  
The input examples,  
as a NumPy array  
The corresponding  
training targets, as  
a NumPy array  
The training loop will  
iterate over the data in  
batches of 128 examples.
```

```
>>> history.history  
{'binary_accuracy": [0.855, 0.9565, 0.9555, 0.95, 0.951],  
 "loss": [0.6573270302042366,  
         0.07434618508815766,  
         0.07687718723714351,  
         0.07412414988875389,  
         0.07617757616937161]}
```



Monitoring loss and metrics on validation data

To keep an eye on **how the model does on new data**, it's standard practice to reserve a subset of the training data as **validation data**: you won't be training the model on this data, but you will use it to compute a loss value and metrics value

Listing 3.24 Using the validation_data argument

```

model = keras.Sequential([keras.layers.Dense(1)])
model.compile(optimizer=keras.optimizers.RMSprop(learning_rate=0.1),
              loss=keras.losses.MeanSquaredError(),
              metrics=[keras.metrics.BinaryAccuracy()])
shuffle training data
indices_permutation = np.random.permutation(len(inputs))
shuffled_inputs = inputs[indices_permutation]
shuffled_targets = targets[indices_permutation]
train set ← 7:3 → valid set
num_validation_samples = int(0.3 * len(inputs))
val_inputs = shuffled_inputs[:num_validation_samples]
val_targets = shuffled_targets[:num_validation_samples]
training_inputs = shuffled_inputs[num_validation_samples:]
training_targets = shuffled_targets[num_validation_samples:]
model.fit(
    training_inputs,
    training_targets,
    epochs=5,
    batch_size=16,
    validation_data=(val_inputs, val_targets)
)
loss_and_metrics = model.evaluate(val_inputs, val_targets, batch_size=128)

```

To avoid having samples from only one class in the validation data, shuffle the inputs and targets using a random indices permutation.

Reserve 30% of the training inputs and targets for validation (we'll exclude these samples from training and reserve them to compute the validation loss and metrics).

Training data, used to update the weights of the model

Validation data, used only to monitor the validation loss and metrics accuracy check if overfitting

You can compute the validation loss and metrics after the training is complete



Inference: Using a model after training

Once you've trained your model, you're going to want to use it to make predictions on new data

```
predictions = model(new_inputs) ← if it is too large → batch_size  
Takes a NumPy array or  
TensorFlow tensor and returns  
a TensorFlow tensor
```

However, this will process all inputs in new_inputs at once, which may not be feasible if you're looking at a lot of data

```
predictions = model.predict(new_inputs, batch_size=128) ←  
Takes a NumPy array or  
a Dataset and returns  
a NumPy array
```



Thank you!

