

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

A dense layer:

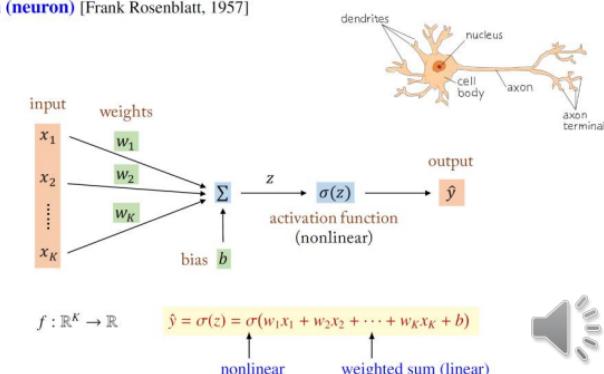
$$\text{output} = \text{relu}(\text{dot}(W, \text{input}) + b)$$

relu($w_i x + b_i$)

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

Gradual adjustment of them is called **training**

Perceptron (neuron) [Frank Rosenblatt, 1957]



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 .
- Update all weights of the network in a way that slightly reduces the loss on this batch.

What we will learn this chapter



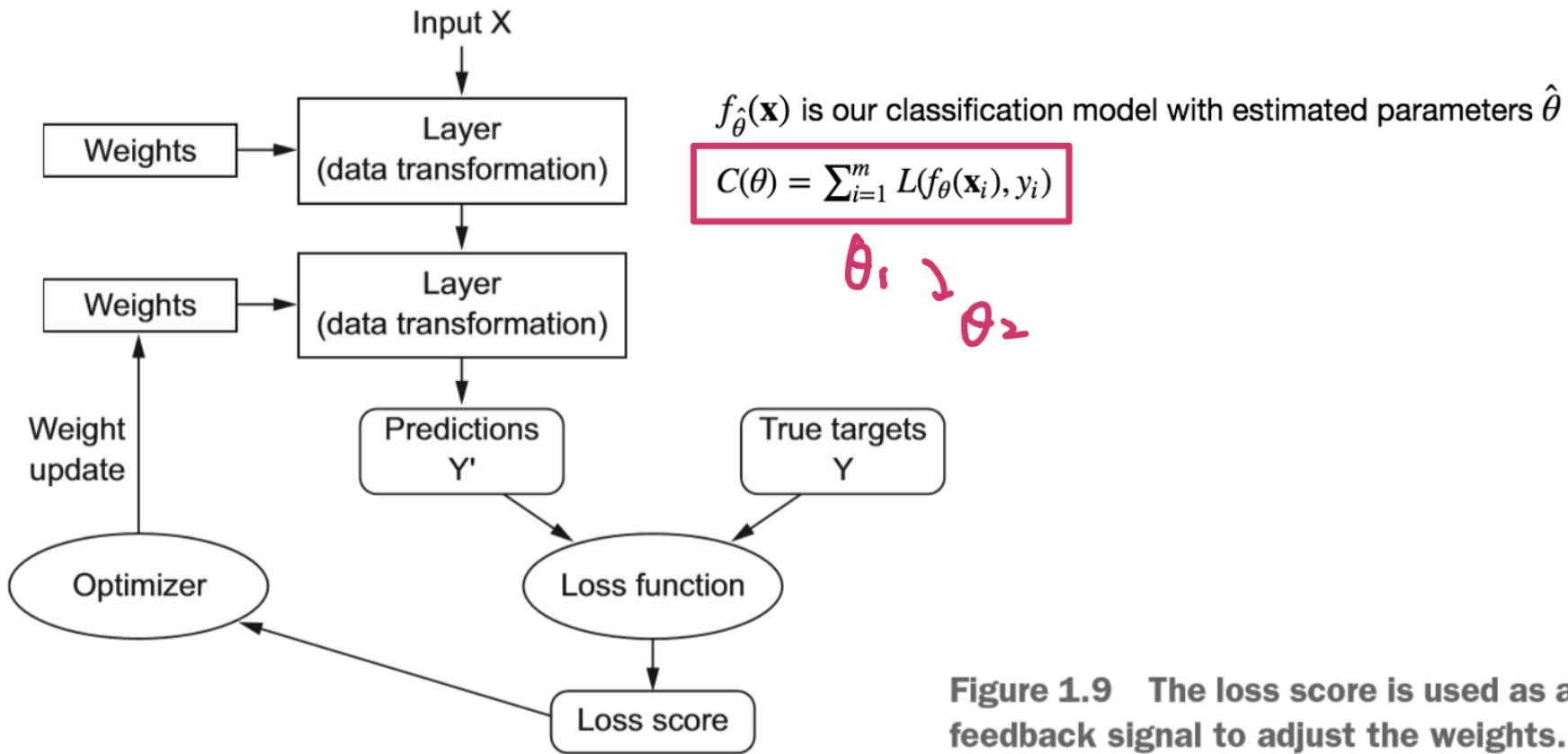
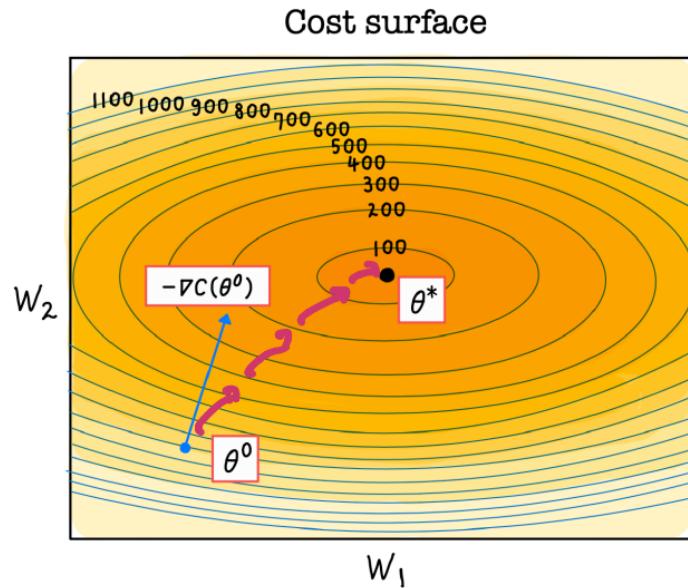


Figure 1.9 The loss score is used as a feedback signal to adjust the weights.



Gradient Descent 경사하강법

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



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

Randomly pick a starting point θ^0 .

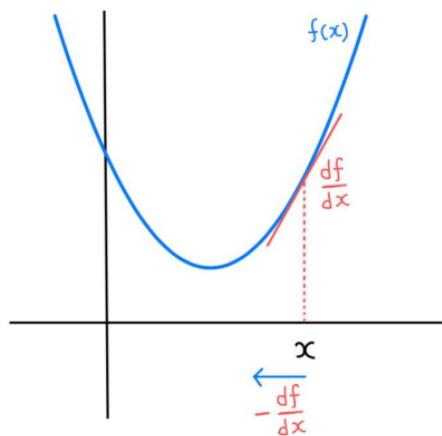
Compute the negative gradient at each θ^t .

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

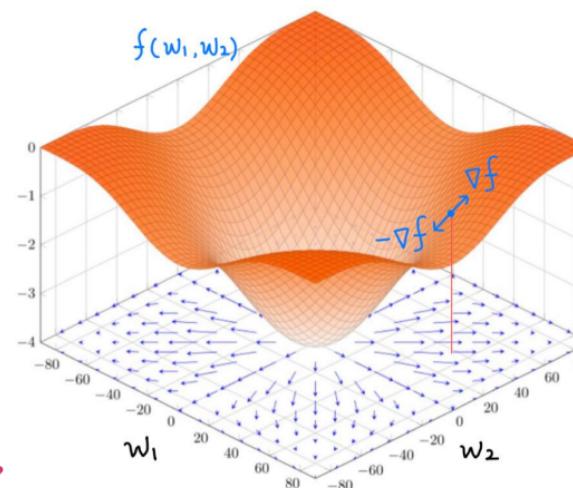
*

$$\nabla_{\theta} C(\theta^t) = \begin{bmatrix} \frac{\partial C(\theta^t)}{\partial w_1} \\ \frac{\partial C(\theta^t)}{\partial w_2} \end{bmatrix}$$

Derivative

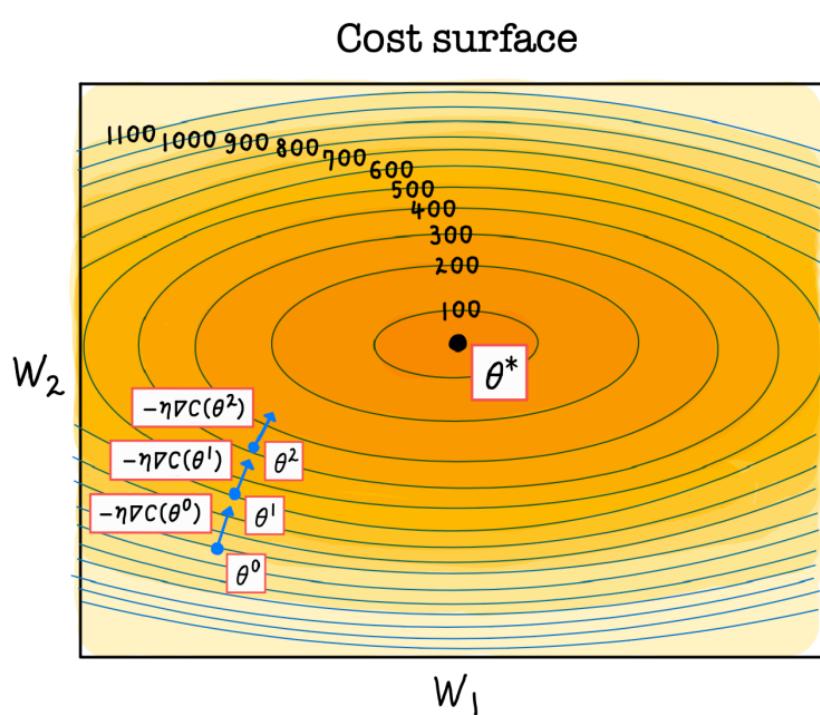


Gradient



$$* \text{ Gradient vector } \nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

indicates the direction and rate of fastest increase.



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

Randomly pick a starting point θ^0 .

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

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

learning rate

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 $C(\theta; X, y) = \sum_{i=1}^m L(\theta; x_i, y_i)$
- Minimize Cost w.r.t θ given data $\operatorname{argmin}_\theta C(\theta; X, y)$
- Iterate gradient updates $\hat{\theta}^* = \theta^t - \eta \nabla_\theta C(\theta^t)$
 $\theta^{t+1} = \theta^t - \eta \nabla_\theta C(\theta^t)$
- $\hat{f}_\theta(\mathbf{x})$ is your classification model



Ex) Binary Classification Model

- Minimize Loss w.r.t θ given data $\operatorname{argmin}_{\theta} C(\theta; \mathbf{X}, \mathbf{y})$
This part may suffer (**memory problem**)

- We can consider using **mini-batch Gradient Descent**, **Stochastic Gradient Descent**

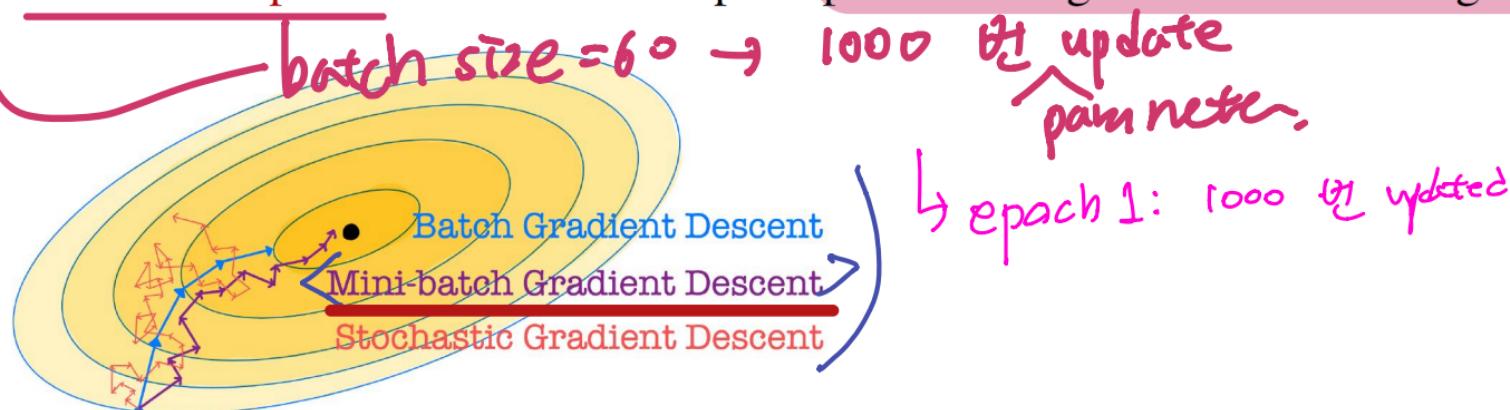




Gradient descent variants

- (1) Batch Gradient Descent: batch size = size of training dataset $\stackrel{60000}{\Rightarrow}$ iteration 1
- (2) Stochastic Gradient Descent (SGD): batch size = 1 $\stackrel{\text{전체 데이터}}{\Rightarrow}$
- * (3) Minibatch Gradient Descent $\stackrel{\text{Batch Gradient or Stochastic Gradient의 중간}}{\Rightarrow}$

- * 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*.

SGD ("Stochastic gradient descent")

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



Chaining derivatives: the Backpropagation algorithm

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

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

$$f_{\theta}(x) = f_3(f_2(f_1(x))),$$

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

$$L(\theta) = \sum_{i=1}^m L_i(\theta)$$

$$= \sum_{i=1}^m H(f_{\theta}(x_i), y_i)$$

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



Backpropagation algorithm

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

x_1

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

grad(loss_val, w)? grad(loss_val, b)?

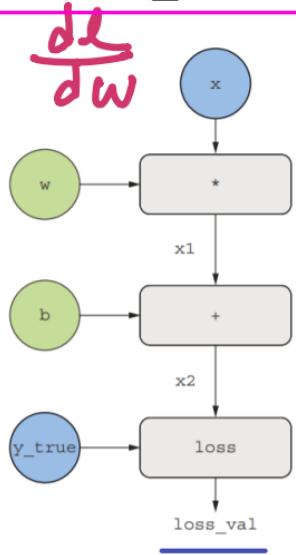


Figure 2.22 A basic example of a computation graph

$$\frac{dl}{db}$$

(이전 과정의)
backpropagation

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

$$x_2 = x_1 + b$$

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

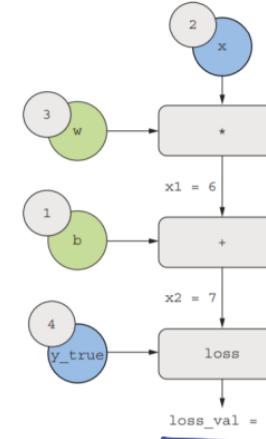


Figure 2.23 Running a forward pass

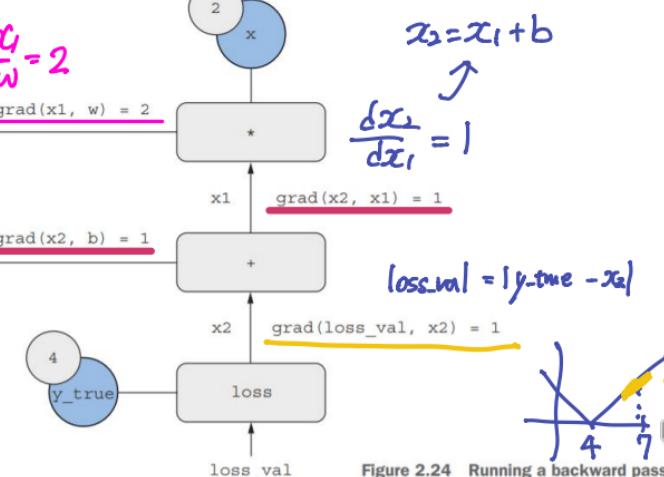
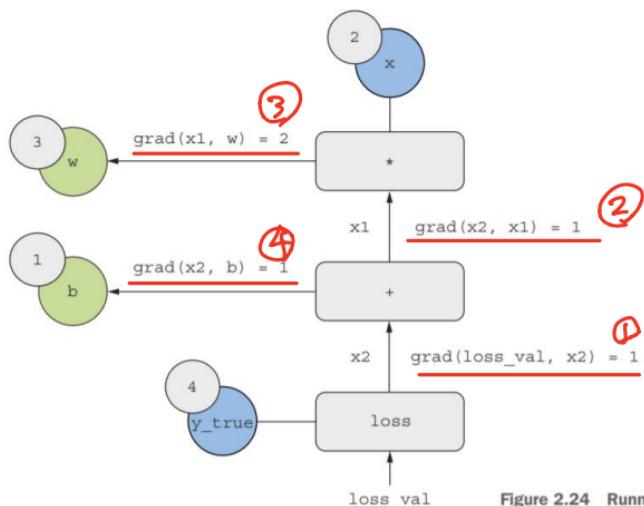


Figure 2.24 Running a backward pass

Backpropagation algorithm

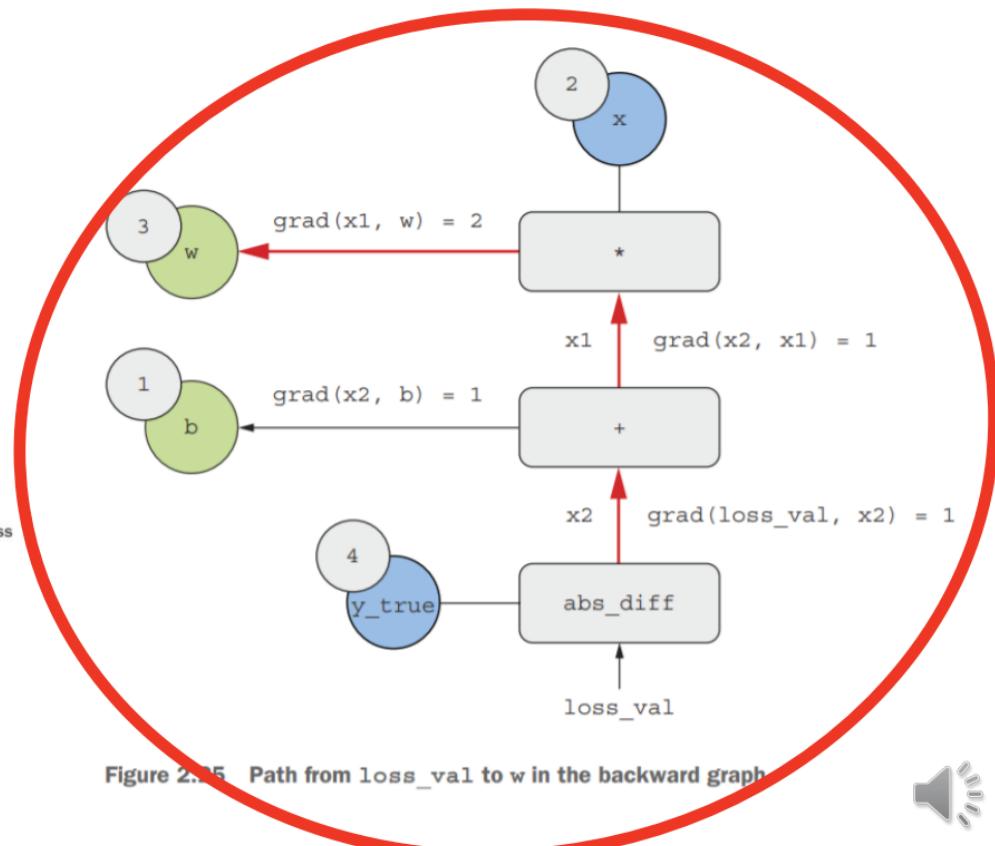


$$\text{grad}(\text{loss_val}, \text{w}) = 1 * 1 * 2 = 2$$

$$\frac{dh}{dw} = \frac{dl}{dx_2} \frac{dx_2}{dx_1} \frac{dx_1}{dw} = 1 \times 1 \times 2 = 2$$

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

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



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:

[\mathbb{C} , $\mathbb{J.C}$, \mathbb{J}]

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

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

```
grad_of_y_wrt_x = tape.gradient(y, x) ←
```

$\frac{dy}{dx}$

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.

$$y = Wx + b \quad \frac{dy}{dW} \quad \text{or} \quad \frac{dy}{db} \quad ??$$



chapter02_mathematical-building-blocks.i

File Edit View Insert Runtime Tools Help Cannot save changes

 Share

1



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](#)

✓ RAM Disk Editing ^

Automatic differentiation with computation graphs

+ Text

▼ The gradient tape in TensorFlow

```
✓ [16] import tensorflow as tf
      0s
      x = tf.Variable(0.)
      with tf.GradientTape() as tape:
          y = 2 * x + 3
      grad_of_y_wrt_x = tape.gradient(y, x)
```

```
✓ [18] 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)
```

```
✓ [20]: 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
grads_of_y_wrt_W_and_b = tape.gradient(y, [W, b])
```

- ▼ Looking back at our first example

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"),
    layers.Dense(10, activation="softmax")
])
```

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)
```



colab.research.google.com/github/fchollet/deep-learning-with-python-notebooks/blob/master/chapter02_mathematical-building-blocks.ipynb#scrollTo=-CTgC3HWb...

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

+ Code + Text 

```
<tf.Tensor: shape=(2,), dtype=float32, numpy=array([2., 2.], dtype=float32)>
```

RAM Disk 

Editing      

Looking back at our first example

```
[ ] (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
```

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

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

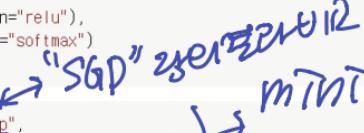
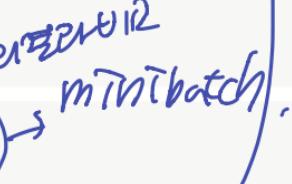
```
[ ] model.fit(train_images, train_labels, epochs=5, batch_size=128)
```

469 x 128 = 60000

469 times updated per epoch

Reimplementing our first example from scratch in TensorFlow

A simple Dense class

"SGD"  

검색하려면 여기에 입력하십시오.

0s completed at 5:54 PM

5°C     오전 5:59 2022-03-06

Reimplementing our first example from scratch in TensorFlow

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

Dense class $\text{output} = \text{activation}(\text{dot}(W, \text{input}) + b)$

```
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 $(\text{input_size}, \text{output_size})$, initialized with random values.

Create a vector, b , of shape $(\text{output_size},)$, initialized with zeros.

Apply the forward pass.

Convenience method for retrieving the layer's weights

$$\begin{matrix} Wx+b \\ \downarrow \quad \downarrow \quad \downarrow \\ q \times p \quad p \times 1 \quad q \times 1 \end{matrix}$$

$$\begin{matrix} XW \\ \downarrow \quad \downarrow \\ p \quad p \times q \end{matrix}$$



Sequential class

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

    def __call__(self, inputs):
        x = inputs
        for layer in self.layers:
            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")
])
```

$[l_1, l_2]$

$l(l_1(x))$

\curvearrowleft

Pizza toppings
28x28

```
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
    
```

) No shuffling in this batch generating

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).

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))
```

Handwritten notes:
- SGD
- Gradient Descent



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}")

```

loops

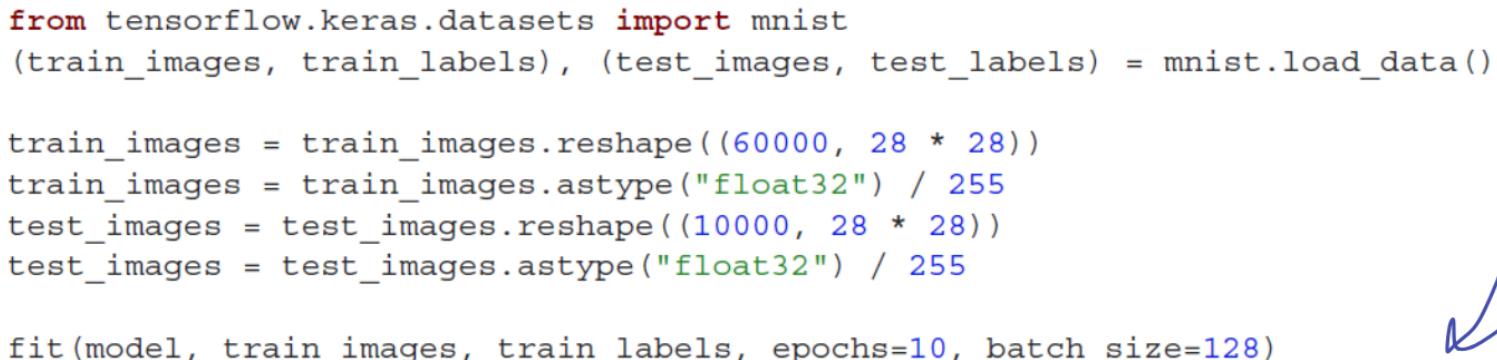
```

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)

```




← → C ⌂ colab.research.google.com/github/fchollet/deep-learning-with-pyth
Pause 00:00:00 Select Area Audio Record Pointer

news study data science 학회사이트 업무관련 기타 Bo

chapter02_mathematical-building-blocks.i

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

RAM Disk

60000/128 468.75

0s

▼ 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

4°C

검색하려면 여기에 입력하십시오.

뉴스 번역 SR Translate scm Data 기타 북마크 읽기 목록

오후 6:20 2022-03-06

Thank you!

