

Deep Learning Study

Week #4

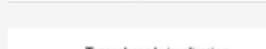
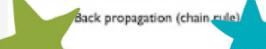
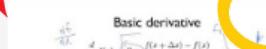
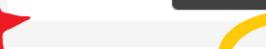
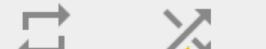
강준하

Fast Review

- Learning rate, Overfitting, Regularization
- Training set / Cross-validation set / Test set
- Perceptron 등장
- XOR Problem (1st winter of deep learning)
 - Backpropagation
- Gradient Vanishing Problem (2nd winter of deep learning)
 - ReLU (다음 주...)

모두를 위한 딥러닝 강좌 시즌 1

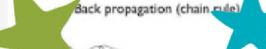
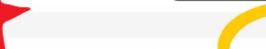
Sung Kim - 25 / 50



lec9-1: XOR 문제 딥러닝으로 풀기

Sung Kim

15:03

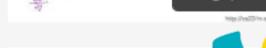


lec9-x: 특별편: 10분안에 미분 정리하기 (lec9-2)

이전에 보세요

Sung Kim

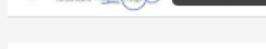
9:29



lec9-2: 딥넷트워크 학습 시키기
(backpropagation)

Sung Kim

18:28



ML lab 09-2: Tensorboard (Neural net for XOR)

Sung Kim

12:08



lec10-1: Sigmoid 보다 ReLU가 더 좋아

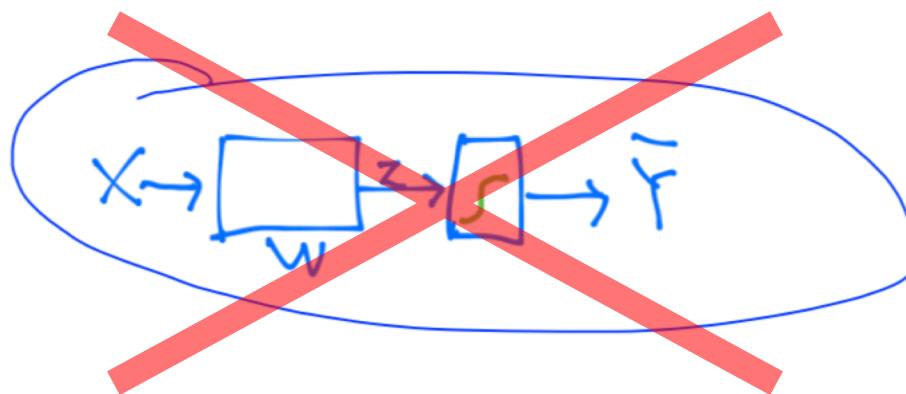
Sung Kim



CONGRATULATIONS

Lec 9-1

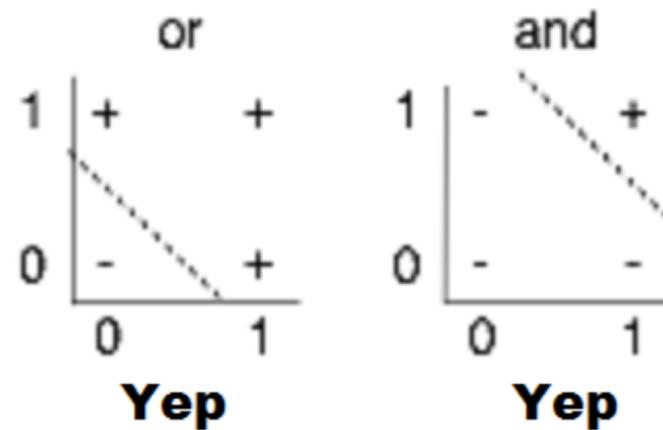
- 한 개의 Logistic Regression unit으로는 XOR 문제를 풀지 못했다...



Lec 9-1

Motivation

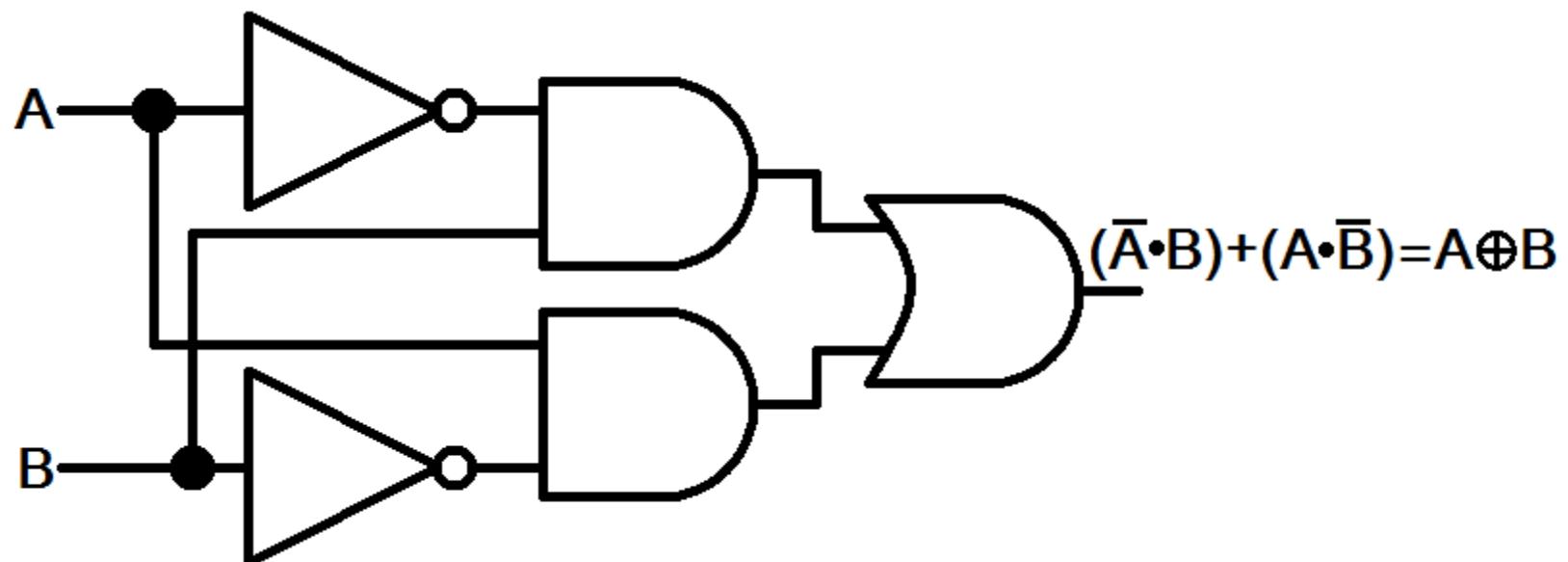
- AND gate와 OR gate는 하나의 regression unit으로도 만들 수 있다



- AND gate와 OR gate를 여러개 이용하면 XOR을 만들 수 있다...?

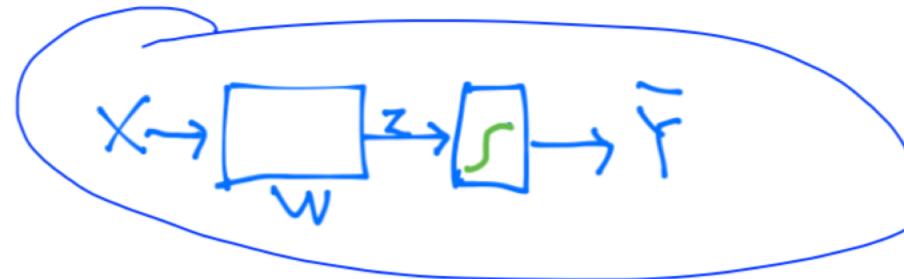
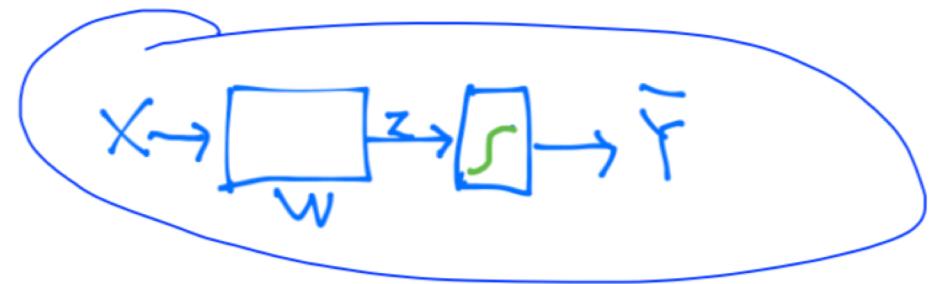
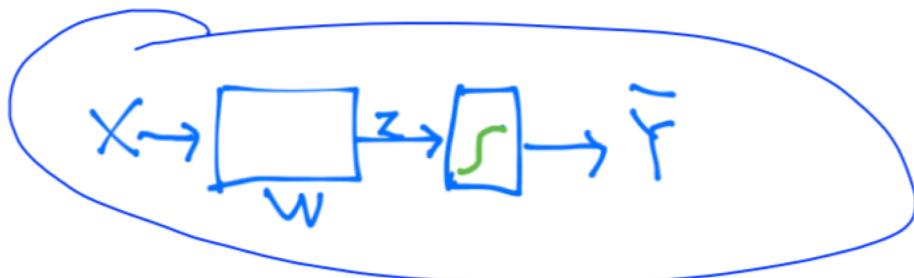
Lec 9-1

- XOR gate by AND gate and OR gate

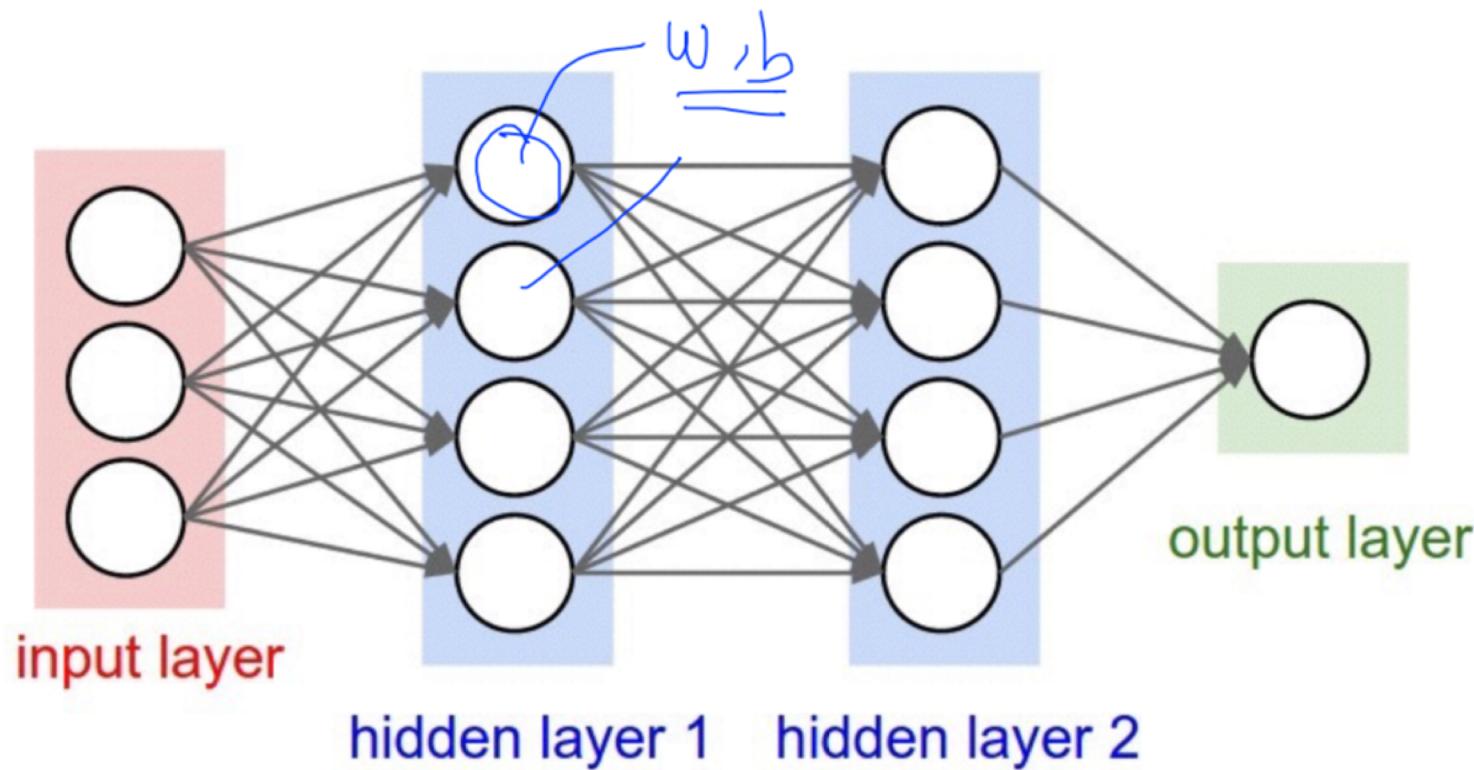


Lec 9-1

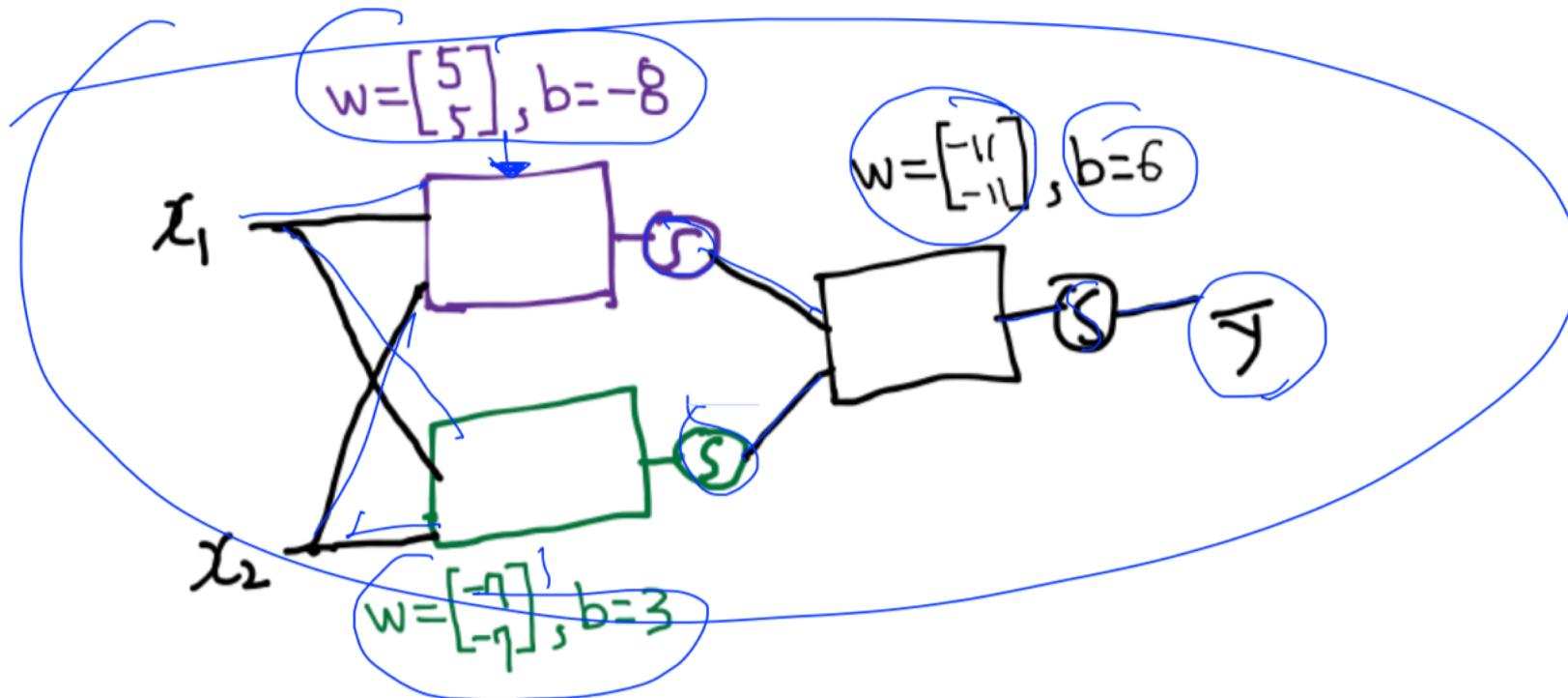
- Use multiple regression units!



Lec 9-1



Lec 9-1



Lec 9-1

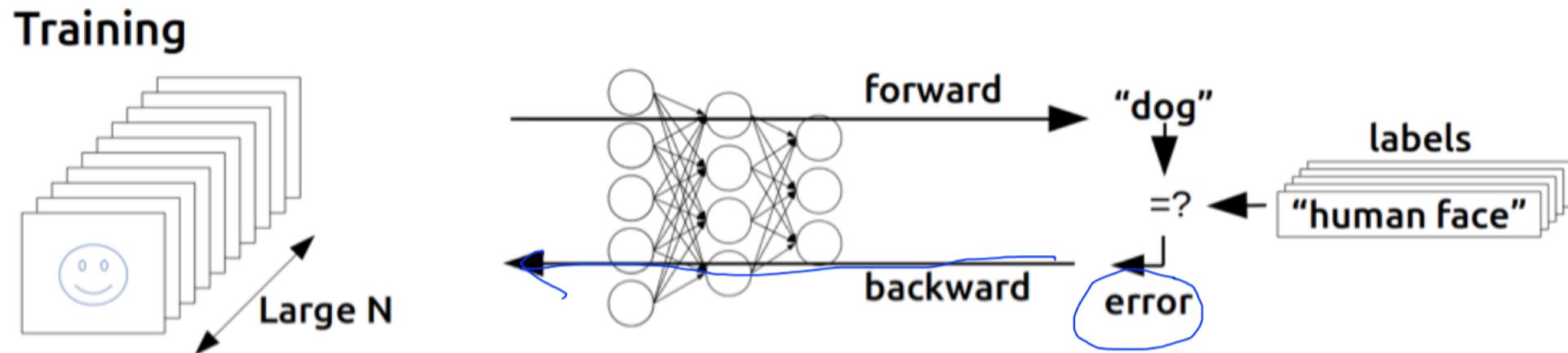
- How to learn multiple W and b?

Lec 9-X (미분 특별편)

- 스킵

Lec 9-2

- How to learn multiple W and b ?
- Backpropagation!



Lec 9-2

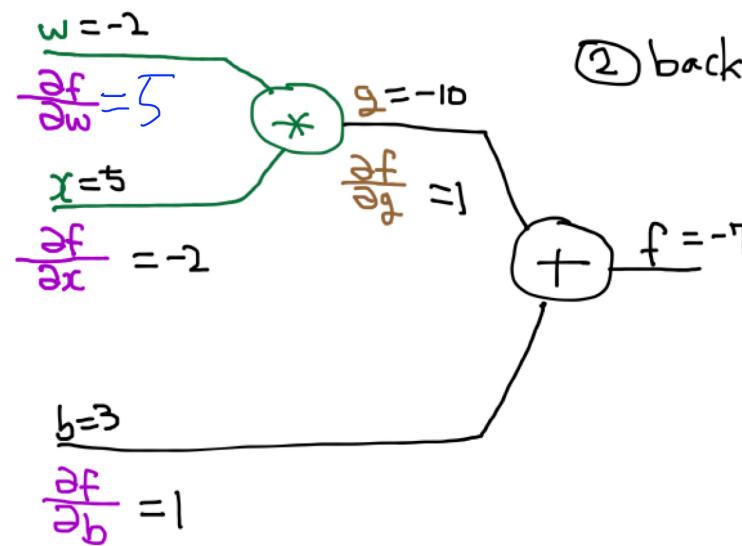
- 미분? 각각의 노드가 결과값에 얼마나/어떻게 영향을 미치는지를 알기 위한 수단!
- Chain Rule? 뒤에서부터 앞으로 거슬러가며 미분값을 계산하는 방법!

Back propagation (chain rule)

$$f = wx + b, g = wx, f = g + b$$

$\frac{\partial f}{\partial g} = 1, \frac{\partial f}{\partial b} = 1$
 $\frac{\partial g}{\partial w} = x, \frac{\partial g}{\partial x} = w$

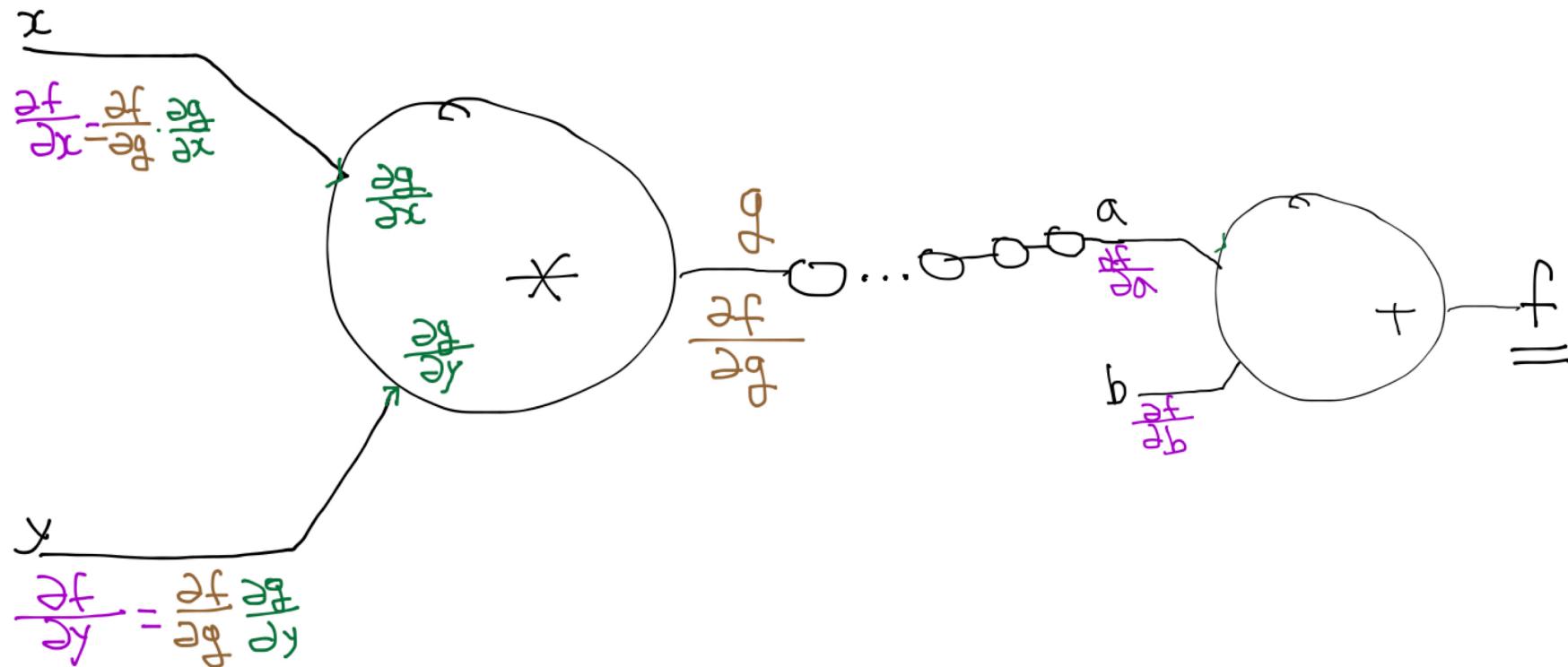
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} = 1 * w = -2$$
$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} = 1 * x = 5$$



① forward ($w = -2, x = 5, b = 3$)

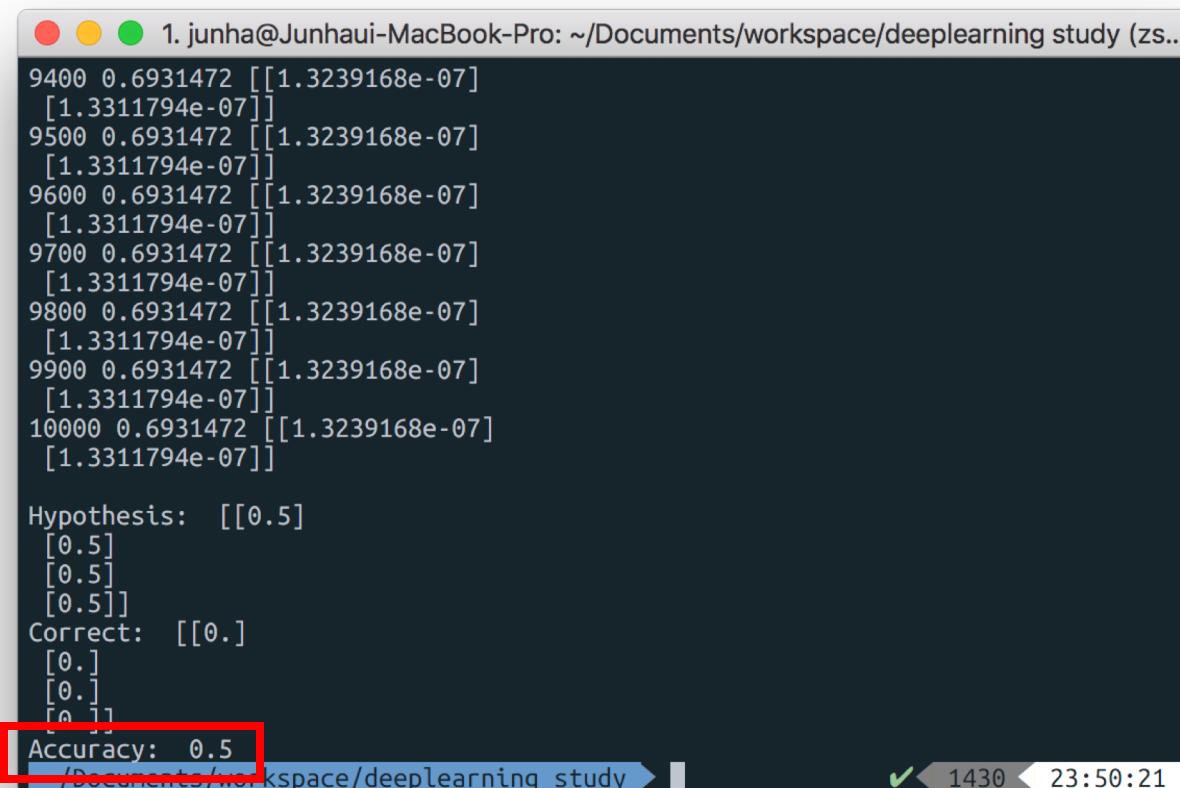
② backward

Back propagation (chain rule)



Lab 9-1

- 한 개의 Logistic Regression Unit을 사용했을 경우



```
1. junha@Junhai-MacBook-Pro: ~/Documents/workspace/deeplearning study (zs...
9400 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
9500 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
9600 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
9700 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
9800 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
9900 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]
10000 0.6931472 [[1.3239168e-07]
[1.3311794e-07]]

Hypothesis: [[0.5]
[0.5]
[0.5]
[0.5]]
Correct: [[0.]
[0.]
[0.]
[0.]]
Accuracy: 0.5
```

Lab 9-1

- Multiple Logistic Regression Unit을 사용했을 경우

```
 1. junha@Junhai-MacBook-Pro: ~/Documents/workspace/deeplearning study (zs...
9600 2.2202755e-06 [array([[-10.695417,  11.116092],
   [ 10.718395, -11.401869]], dtype=float32), array([[26.289412],
   [26.411997]], dtype=float32)]
9700 2.101066e-06 [array([[-10.710199,  11.130635],
   [ 10.733139, -11.416296]], dtype=float32), array([[26.38956 ],
   [26.512707]], dtype=float32)]
9800 2.011659e-06 [array([[-10.724949,  11.144869],
   [ 10.747823, -11.430691]], dtype=float32), array([[26.489662],
   [26.612396]], dtype=float32)]
9900 1.9073505e-06 [array([[-10.739635,  11.159071],
   [ 10.762475, -11.445048]], dtype=float32), array([[26.589813],
   [26.712193]], dtype=float32)]
10000 1.7881409e-06 [array([[-10.754227,  11.173415],
   [ 10.777167, -11.459353]], dtype=float32), array([[26.690191],
   [26.812716]], dtype=float32)]

Hypothesis: [[1.9362387e-06]
 [9.9999821e-01]
 [9.9999845e-01]
 [1.9069422e-06]]
Correct: [[0.]
 [1.]
 [1.]
 [0.]]
Accuracy: 1.0
```

Lab 9-1

- Multilayer structure

```
w1 = tf.Variable(tf.random_normal([2, 2]), name='weight1')
b1 = tf.Variable(tf.random_normal([2]), name='bias1')
layer1 = tf.sigmoid(tf.matmul(X, w1) + b1)

w2 = tf.Variable(tf.random_normal([2, 1]), name='weight2')
b2 = tf.Variable(tf.random_normal([1]), name='bias2')
hypothesis = tf.sigmoid(tf.matmul(layer1, w2) + b2)
```

Lab 9-1

- Multilayer structure
- More **wider**...

```
W1 = tf.Variable(tf.random_normal([2, 10]), name='weight1')
b1 = tf.Variable(tf.random_normal([10]), name='bias1')
layer1 = tf.sigmoid(tf.matmul(X, W1) + b1)

W2 = tf.Variable(tf.random_normal([10, 1]), name='weight2')
b2 = tf.Variable(tf.random_normal([1]), name='bias2')
hypothesis = tf.sigmoid(tf.matmul(layer1, W2) + b2)
```

Lab 9-1

- Multilayer structure
- More **deeper...**

```
w1 = tf.Variable(tf.random_normal([2, 10]), name='weight1')
b1 = tf.Variable(tf.random_normal([10]), name='bias1')
layer1 = tf.sigmoid(tf.matmul(X, w1) + b1)

w2 = tf.Variable(tf.random_normal([10, 10]), name='weight2')
b2 = tf.Variable(tf.random_normal([10]), name='bias2')
layer2 = tf.sigmoid(tf.matmul(layer1, w2) + b2)

w3 = tf.Variable(tf.random_normal([10, 10]), name='weight3')
b3 = tf.Variable(tf.random_normal([10]), name='bias3')
layer3 = tf.sigmoid(tf.matmul(layer2, w3) + b3)

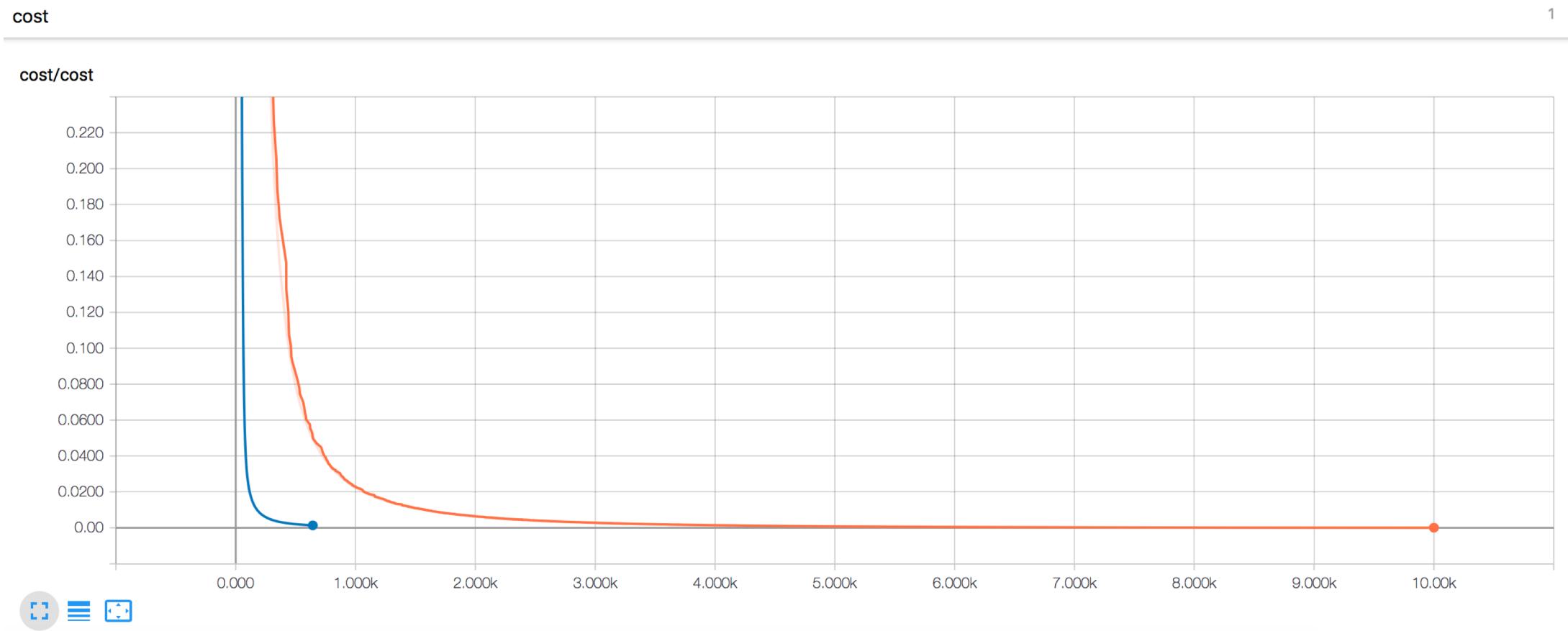
w4 = tf.Variable(tf.random_normal([10, 1]), name='weight4')
b4 = tf.Variable(tf.random_normal([1]), name='bias4')
hypothesis = tf.sigmoid(tf.matmul(layer3, w4) + b4)
```

Lab 9-2

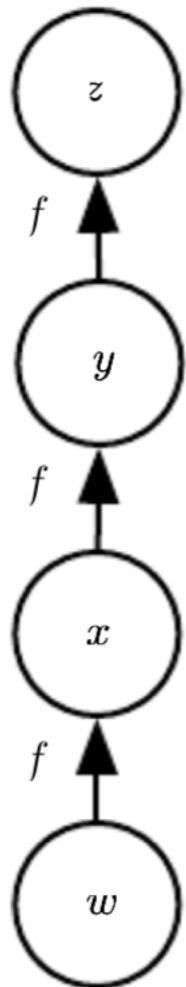
- About Tensorboard
- 자세한 문법은 생략...

Lab 9-2

- 파란색 – Learning Rate 0.1, 주황색 – Learning Rate 0.01



More about backpropagation



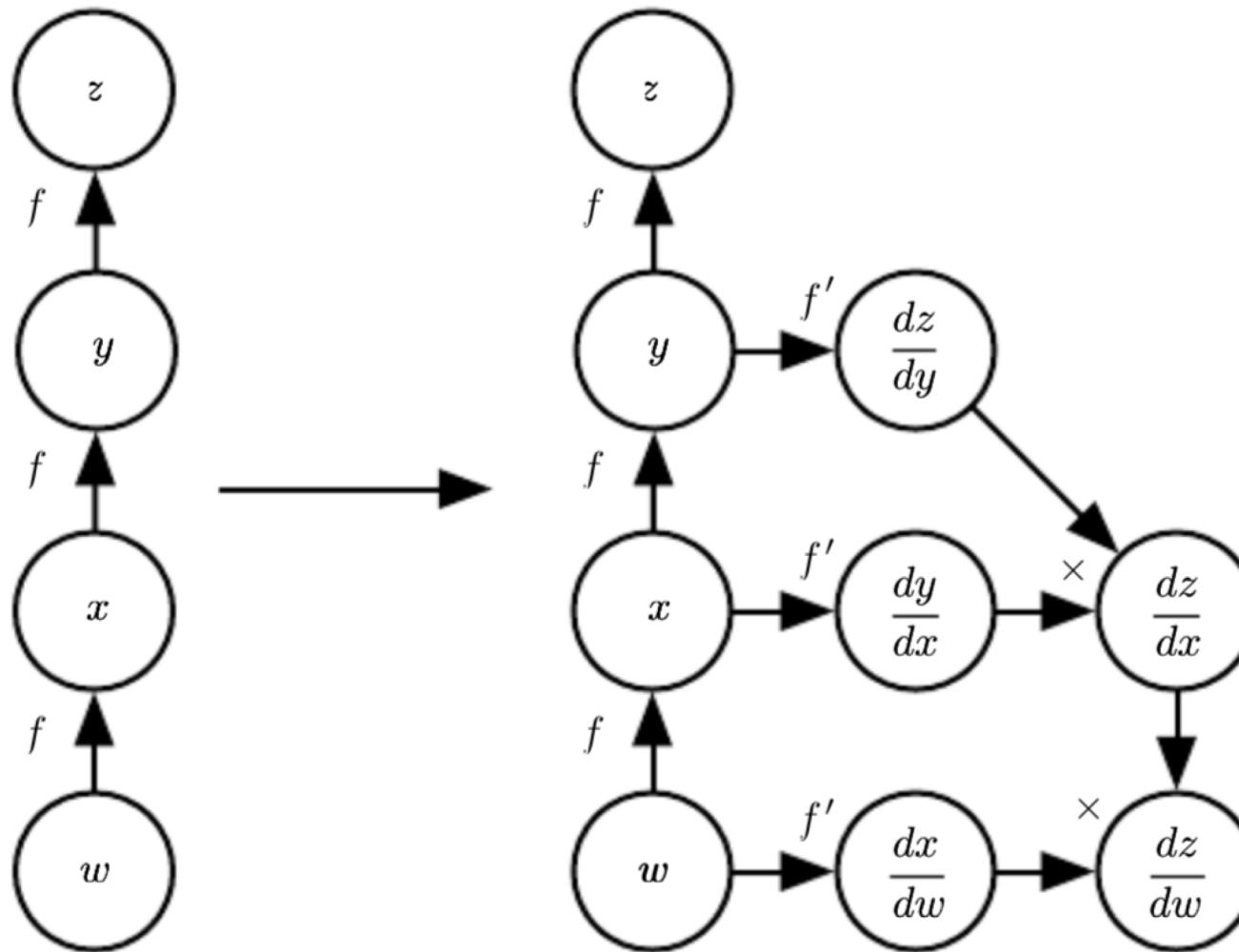
$$\frac{\partial z}{\partial w} \quad (6.50)$$

$$= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \quad (6.51)$$

$$= f'(y) f'(x) f'(w) \quad (6.52)$$

$$= f'(f(f(w))) f'(f(w)) f'(w) \quad (6.53)$$

More about backpropagation



More about backpropagation

Algorithm 6.3 Forward propagation through a typical deep neural network and the computation of the cost function. The loss $L(\hat{\mathbf{y}}, \mathbf{y})$ depends on the output $\hat{\mathbf{y}}$ and on the target \mathbf{y} (see section 6.2.1.1 for examples of loss functions). To obtain the total cost J , the loss may be added to a regularizer $\Omega(\theta)$, where θ contains all the parameters (weights and biases). Algorithm 6.4 shows how to compute gradients of J with respect to parameters \mathbf{W} and \mathbf{b} . For simplicity, this demonstration uses only a single input example \mathbf{x} . Practical applications should use a minibatch. See section 6.5.7 for a more realistic demonstration.

Require: Network depth, l
Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model
Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model
Require: \mathbf{x} , the input to process
Require: \mathbf{y} , the target output

```
 $\mathbf{h}^{(0)} = \mathbf{x}$ 
for  $k = 1, \dots, l$  do
     $\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$ 
     $\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$ 
end for
 $\hat{\mathbf{y}} = \mathbf{h}^{(l)}$ 
 $J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$ 
```

More about backpropagation

Algorithm 6.4 Backward computation for the deep neural network of algorithm 6.3, which uses in addition to the input \mathbf{x} a target \mathbf{y} . This computation yields the gradients on the activations $\mathbf{a}^{(k)}$ for each layer k , starting from the output layer and going backwards to the first hidden layer. From these gradients, which can be interpreted as an indication of how each layer's output should change to reduce error, one can obtain the gradient on the parameters of each layer. The gradients on weights and biases can be immediately used as part of a stochastic gradient update (performing the update right after the gradients have been computed) or used with other gradient-based optimization methods.

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

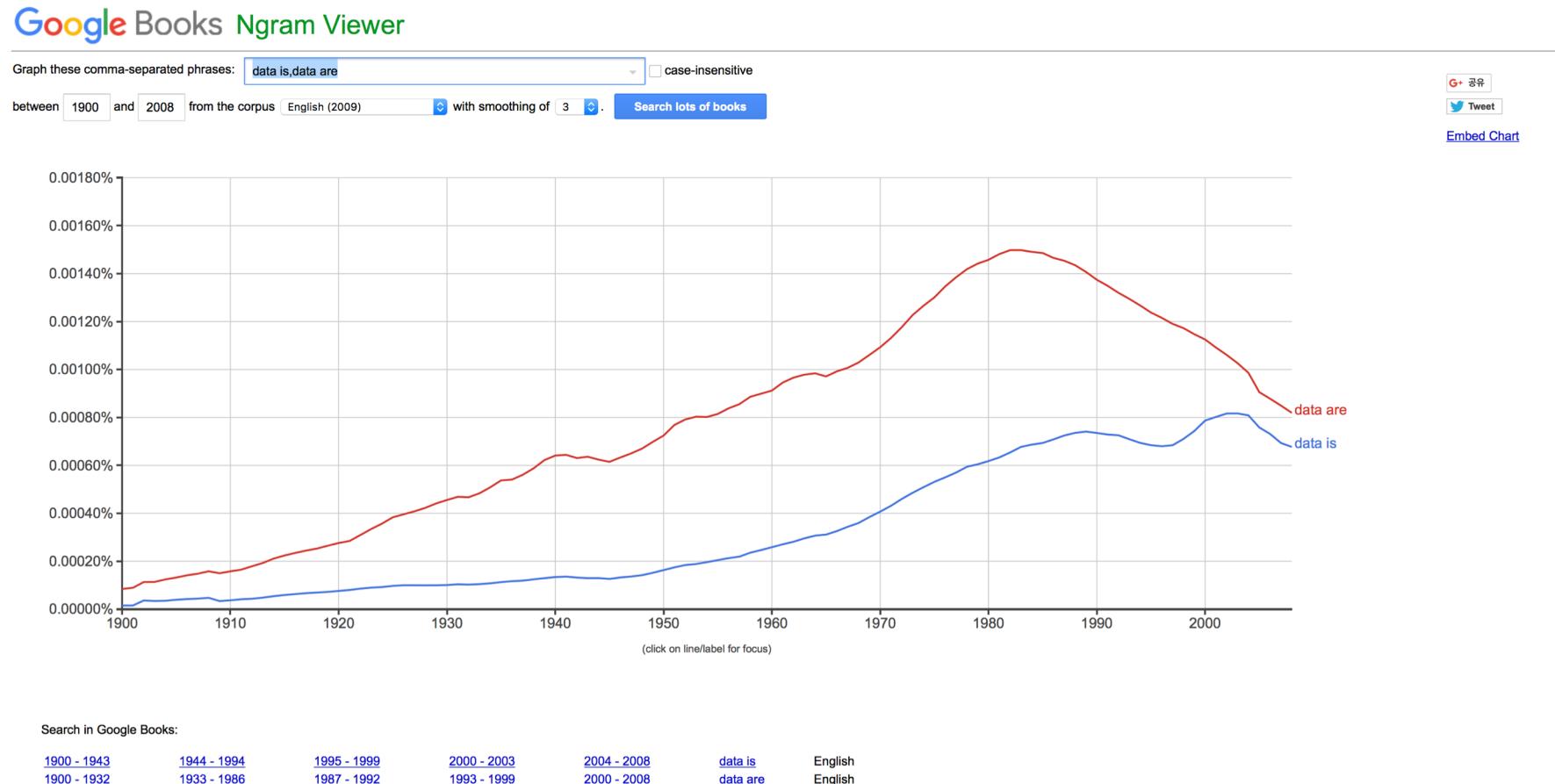
$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

More about backpropagation

- 칠판

Google n-gram viewer



<https://zariski.wordpress.com/2010/12/22/거대한-코퍼스로-놀기/>

———— How **not** to build a minimum viable product ——



———— How **to** build a minimum viable product ——

