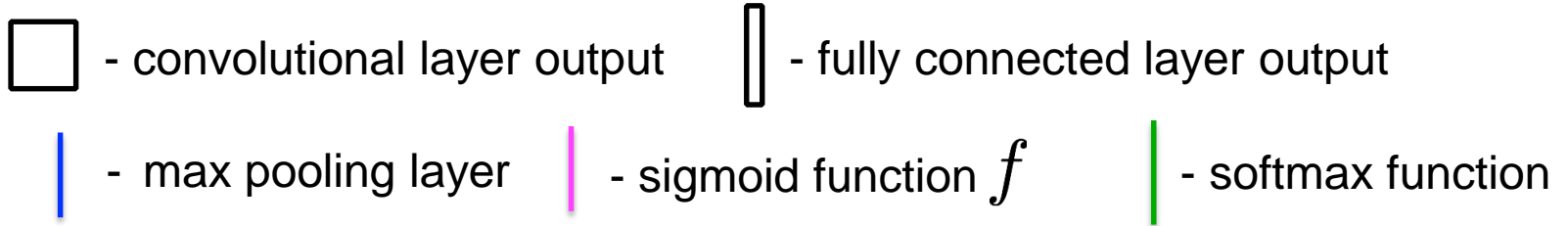


# Improving Learning in Neural Networks

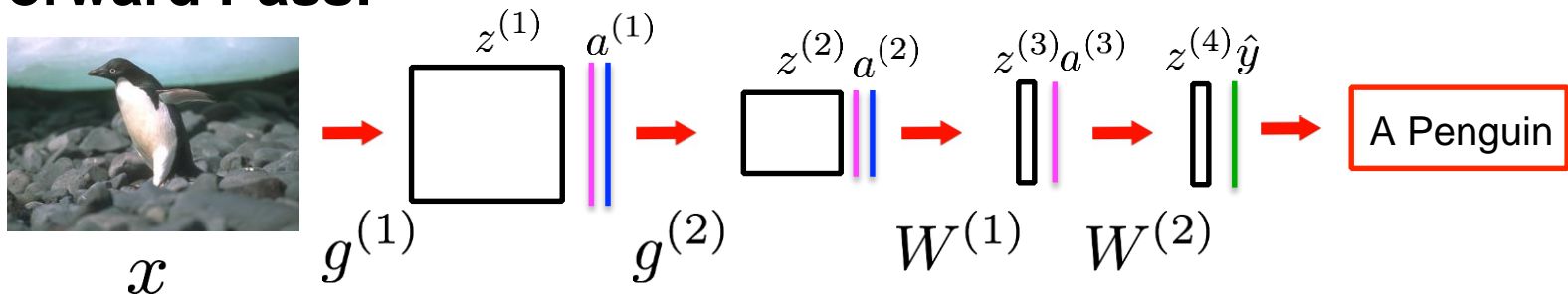
CIS 680

# Convolutional Networks

## Notation:

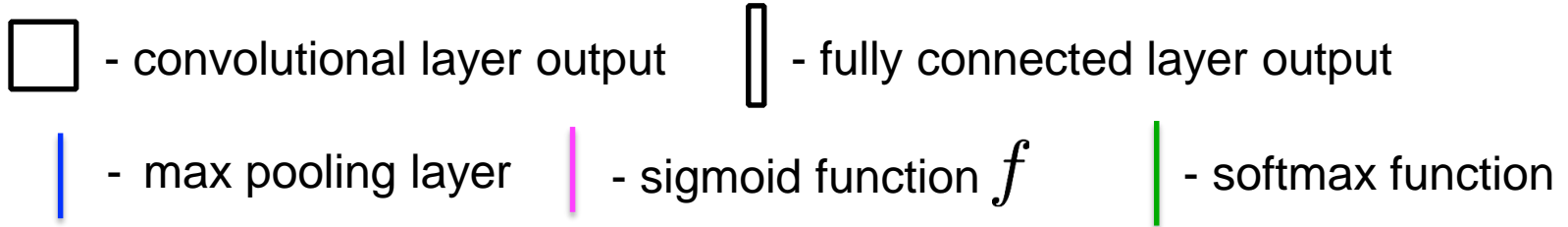


## Forward Pass:

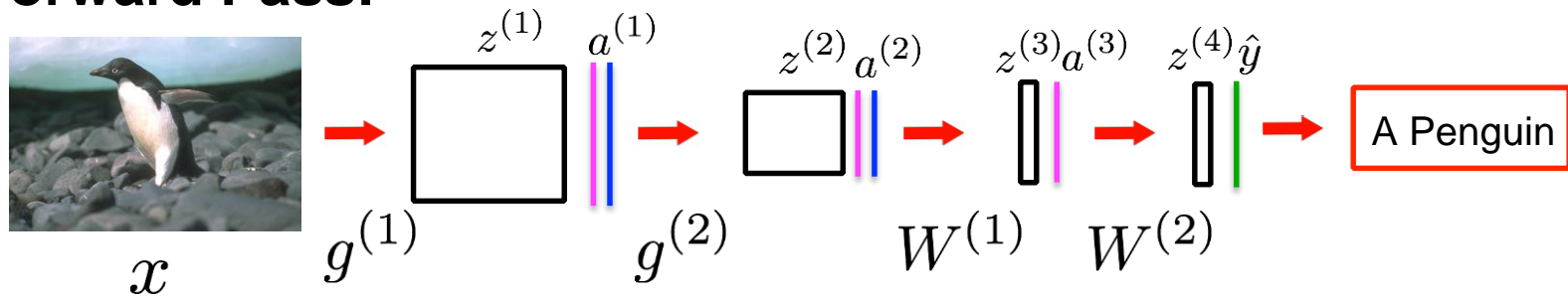


# Convolutional Networks

## Notation:



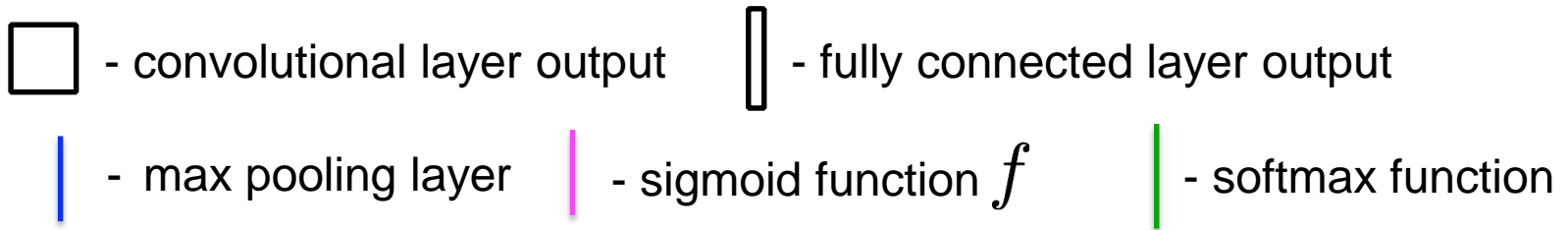
## Forward Pass:



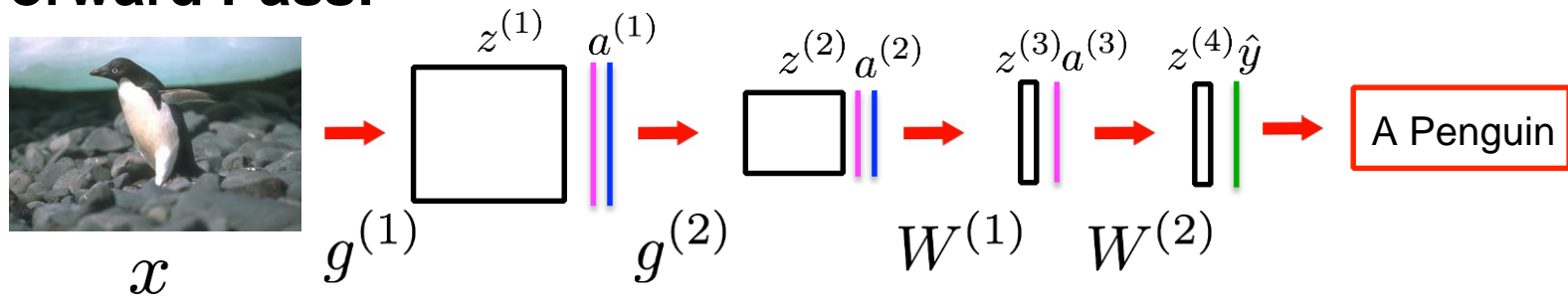
$$1. \quad a^{(1)} = \text{pool}(f(g^{(1)} * x))$$

# Convolutional Networks

## Notation:



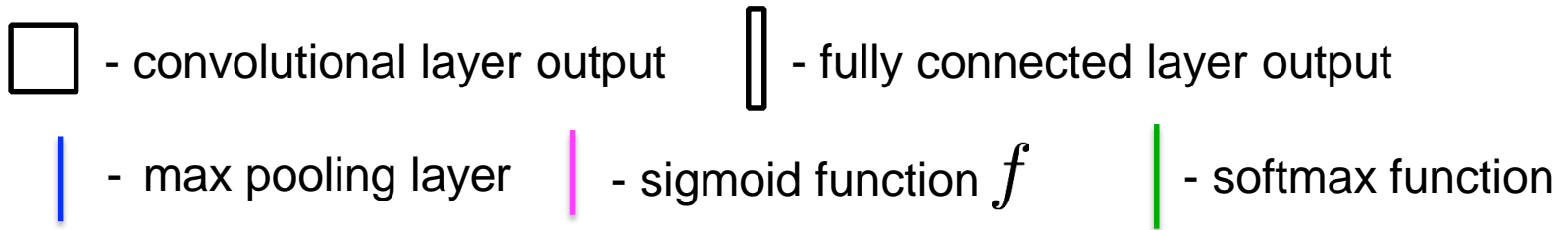
## Forward Pass:



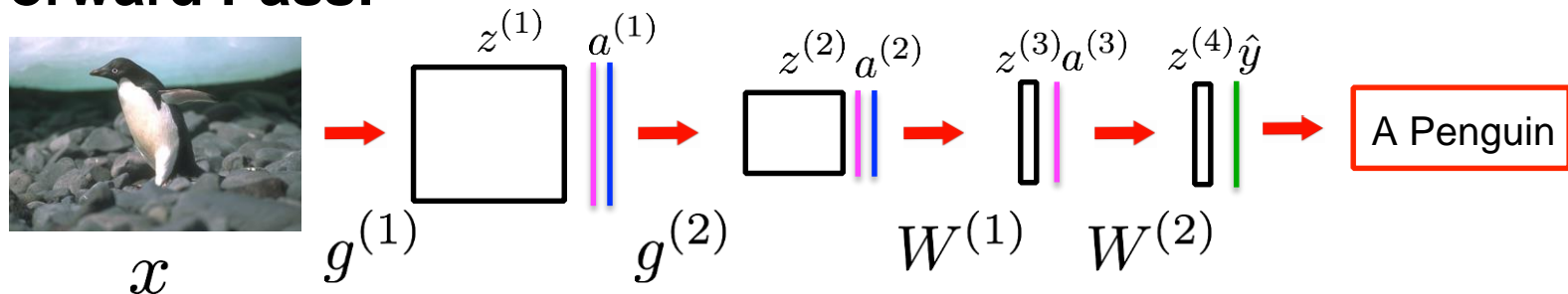
- $a^{(1)} = \text{pool}(f(g^{(1)} * x))$
- $a^{(2)} = \text{pool}(f(g^{(2)} * a^{(1)}))$

# Convolutional Networks

## Notation:



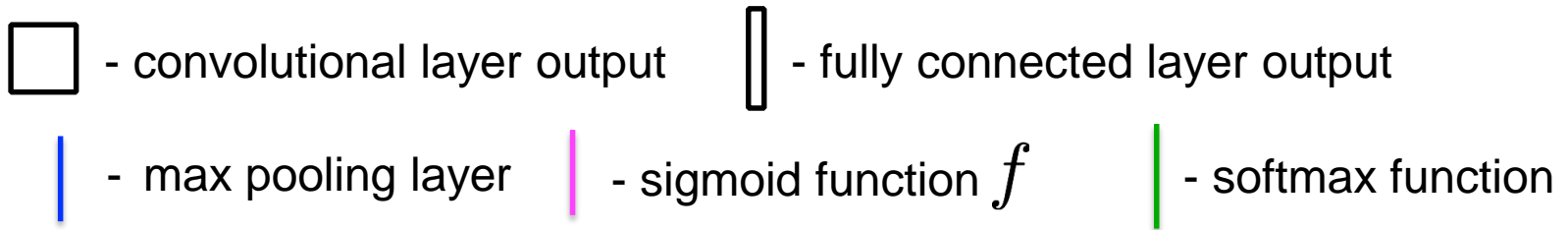
## Forward Pass:



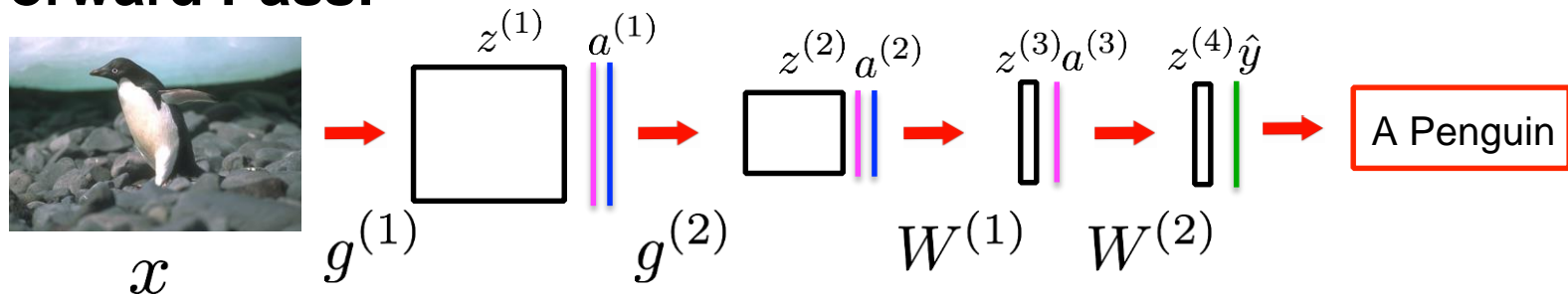
1.  $a^{(1)} = \text{pool}(f(g^{(1)} * x))$
2.  $a^{(2)} = \text{pool}(f(g^{(2)} * a^{(1)}))$
3.  $a^{(3)} = f(W^{(1)} a^{(2)})$

# Convolutional Networks

## Notation:



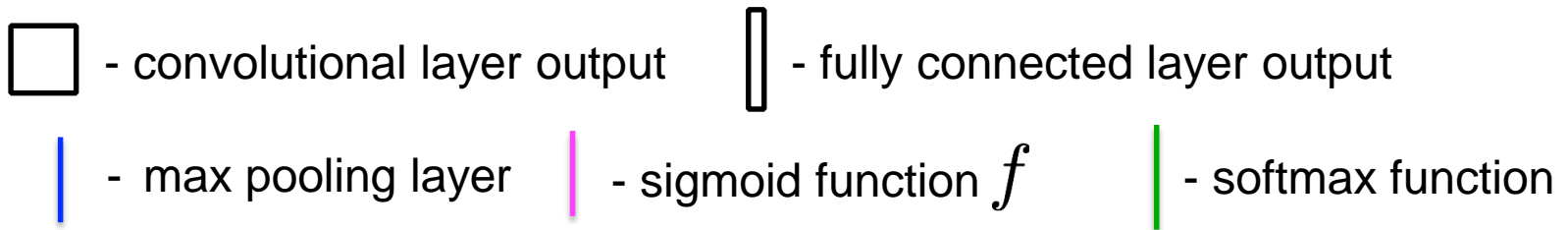
## Forward Pass:



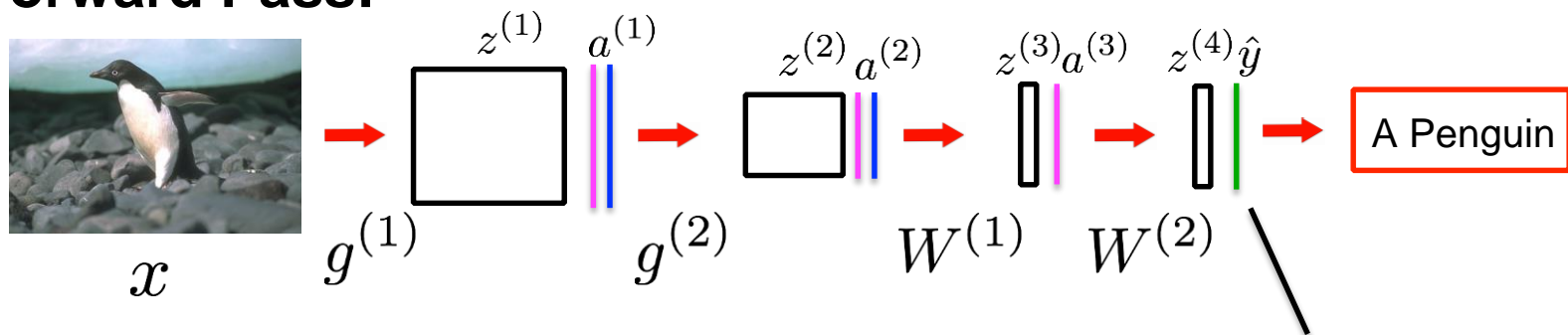
1.  $a^{(1)} = \text{pool}(f(g^{(1)} * x))$
2.  $a^{(2)} = \text{pool}(f(g^{(2)} * a^{(1)}))$
3.  $a^{(3)} = f(W^{(1)} a^{(2)})$
4.  $\hat{y} = \text{softmax}(W^{(2)} a^{(3)})$

# Convolutional Networks

## Notation:



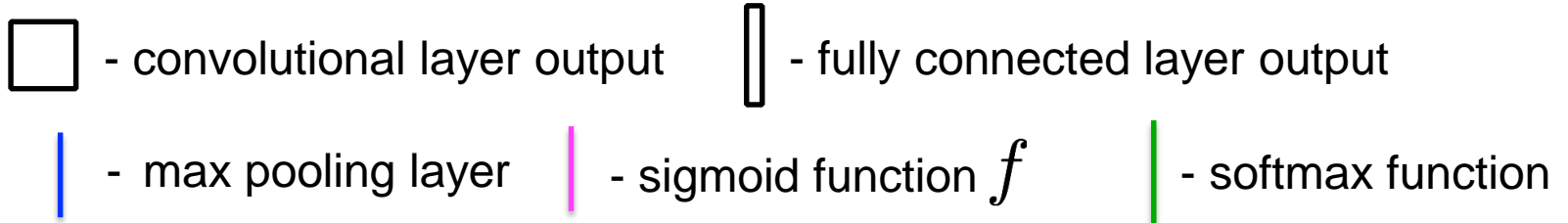
## Forward Pass:



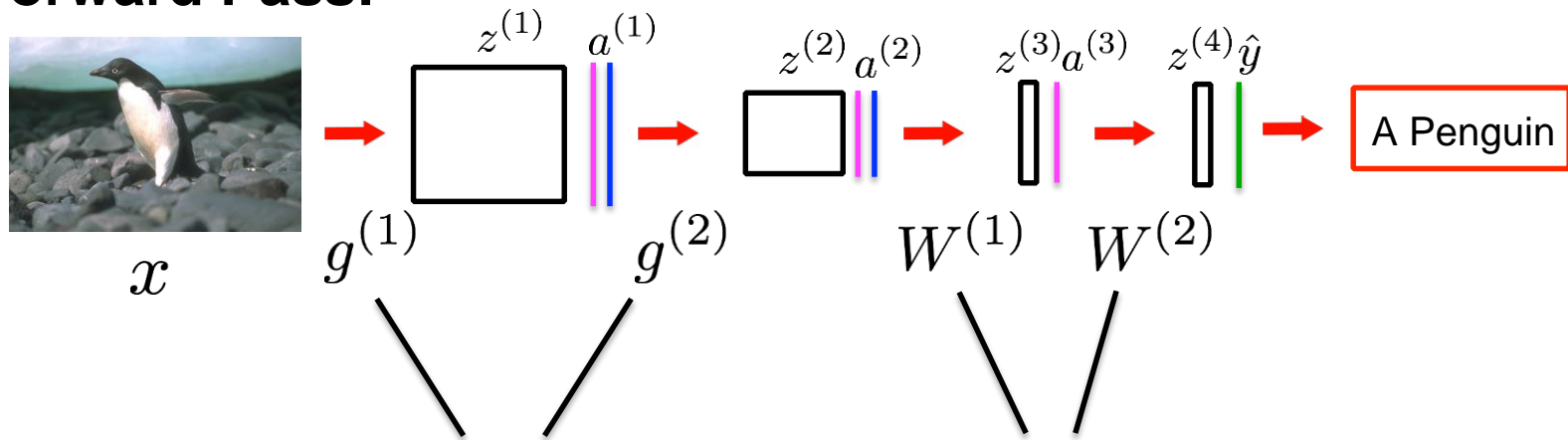
Predicted probabilities, which class this image belongs to

# Convolutional Networks

## Notation:



## Forward Pass:

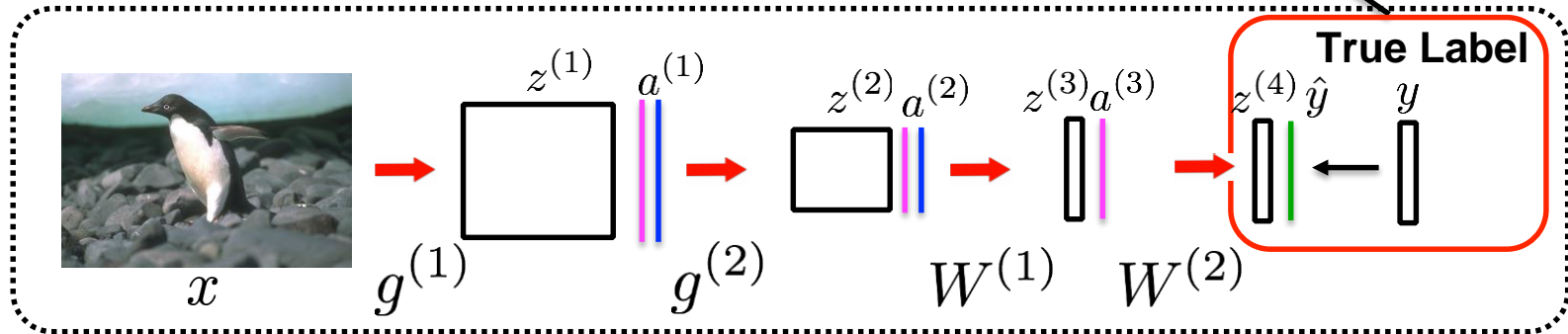


**How to learn the parameters from the data?**

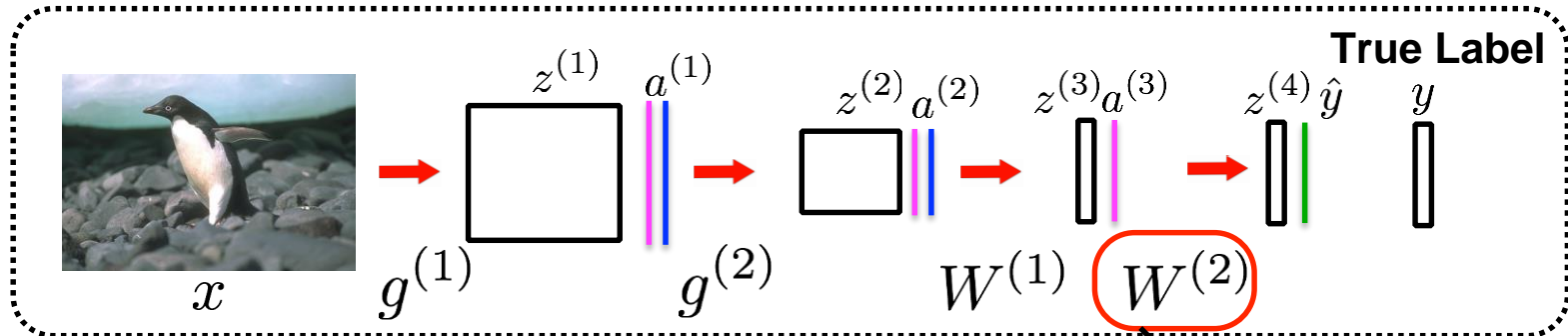


# Backpropagation

1. Compute the gradients of the overall loss and propagate it back:  $\frac{\partial L}{\partial z^{(4)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(4)}}$



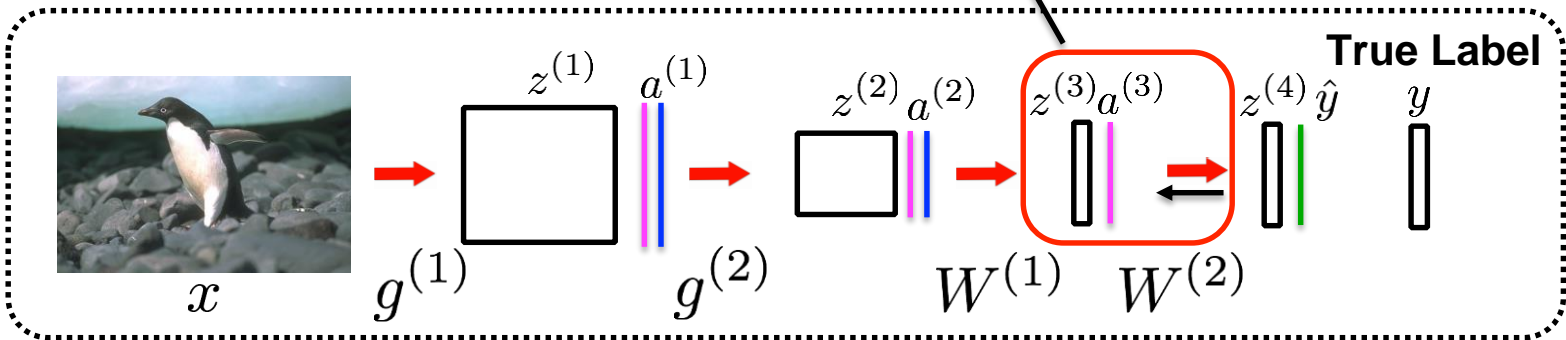
# Backpropagation



2. Compute the gradients to adjust the weights:  $\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial W^{(2)}}$  where  $z^{(4)} = W^{(2)} a^{(3)}$

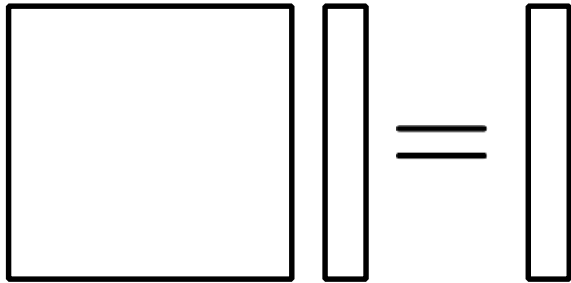
# Backpropagation

3. Backpropagate the gradients to previous layers:  $\frac{\partial L}{\partial z^{(3)}} = \frac{\partial L}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial f(z^{(3)})} \frac{\partial f(z^{(3)})}{\partial z^{(3)}}$  where  $z^{(4)} = W^{(2)} a^{(3)}$



## Fully Connected Layers:

Forward:

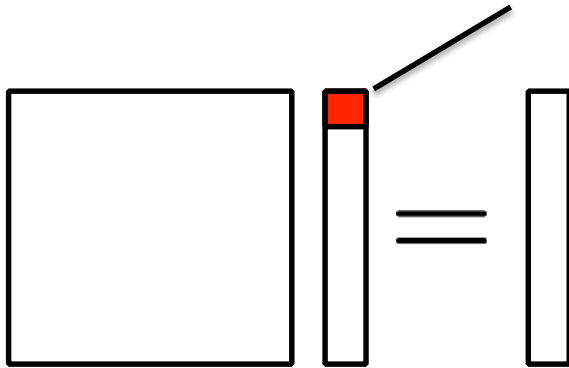


$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$

## Fully Connected Layers:

Forward:

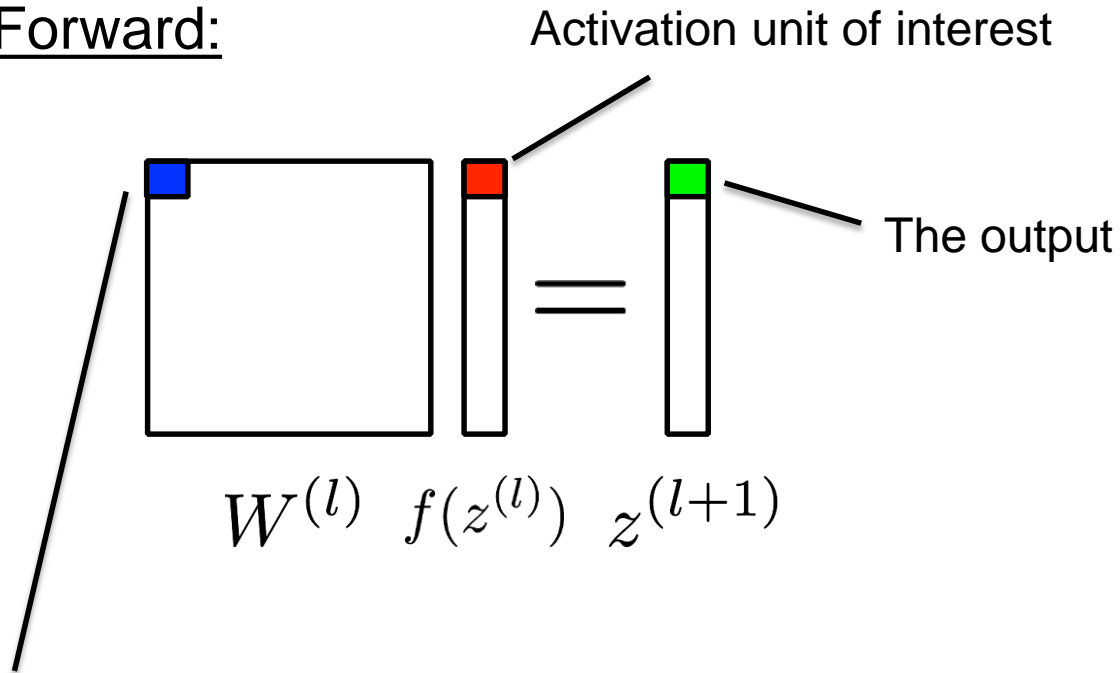
Activation unit of interest



$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$

# Fully Connected Layers:

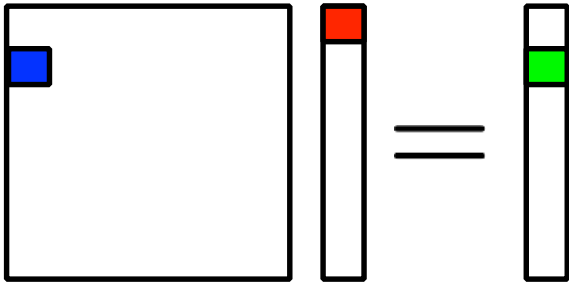
Forward:



The weight that is used in conjunction with the activation unit of interest

## Fully Connected Layers:

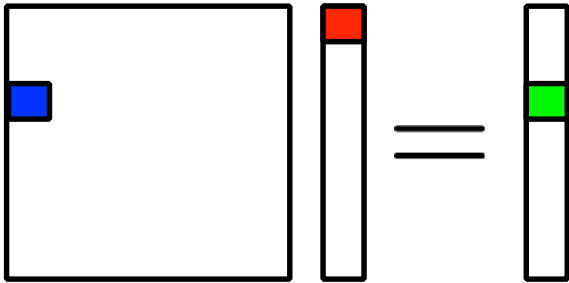
Forward:



$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$

## Fully Connected Layers:

Forward:

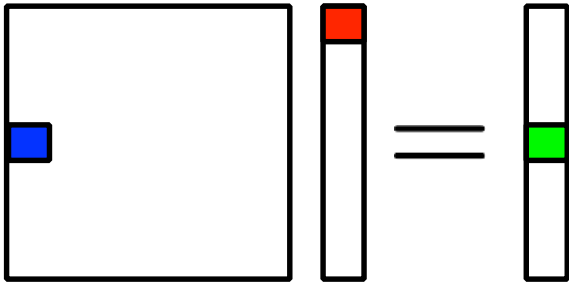


$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$



## Fully Connected Layers:

Forward:

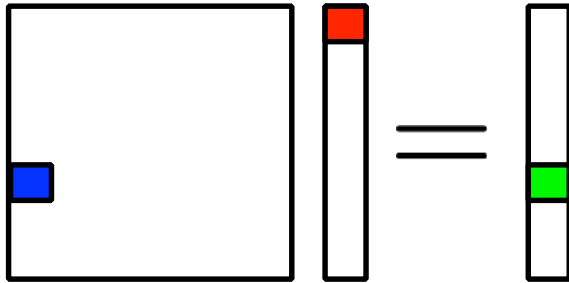


$$W^{(l)} f(z^{(l)}) = z^{(l+1)}$$

# Backpropagation

## Fully Connected Layers:

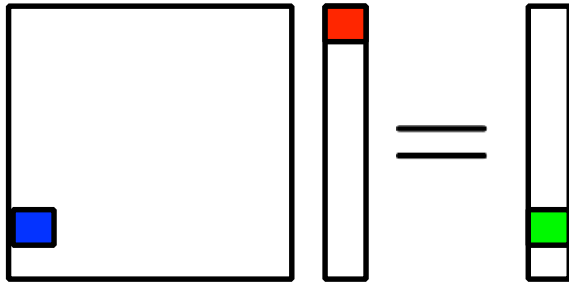
Forward:



$$W^{(l)} f(z^{(l)}) = z^{(l+1)}$$

## Fully Connected Layers:

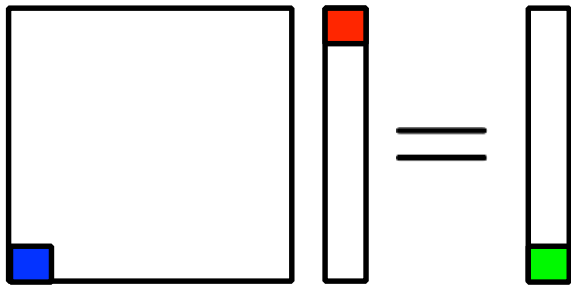
Forward:



$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$

## Fully Connected Layers:

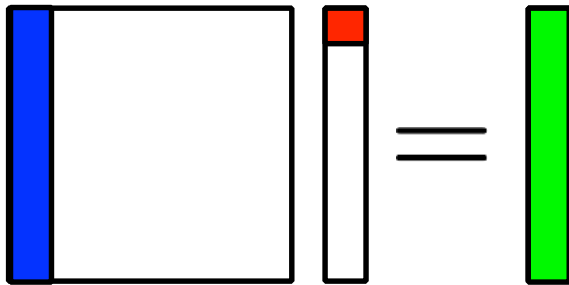
Forward:



$$W^{(l)} f(z^{(l)}) = z^{(l+1)}$$

## Fully Connected Layers:

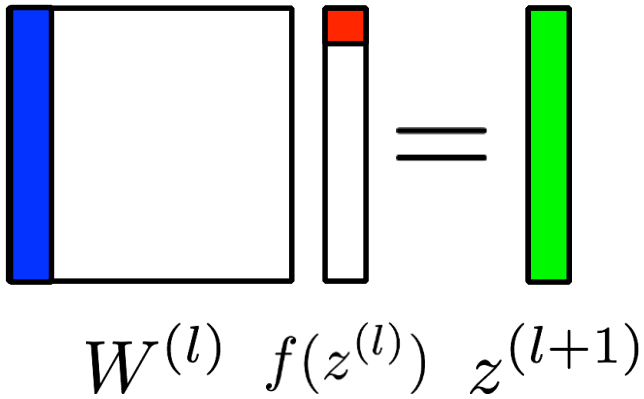
Forward:



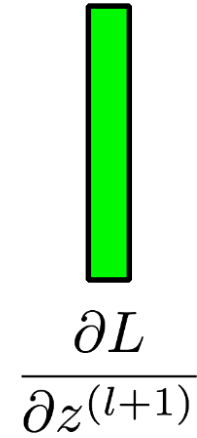
$$W^{(l)} \quad f(z^{(l)}) \quad z^{(l+1)}$$

## Fully Connected Layers:

Forward:

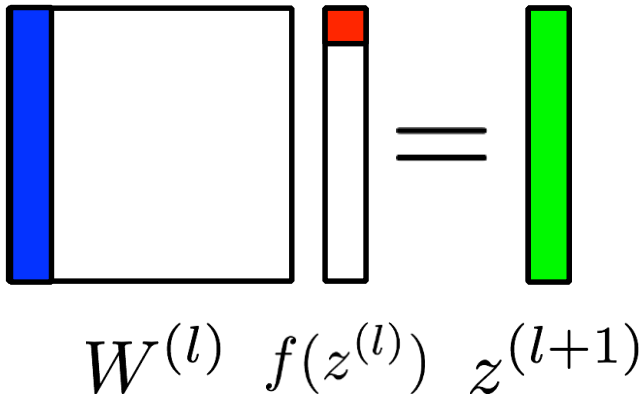


Backward:



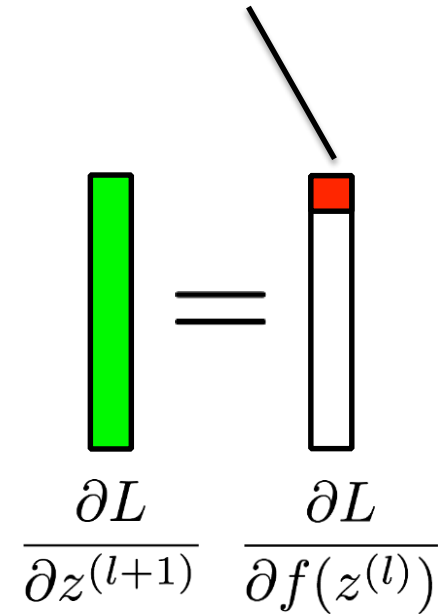
## Fully Connected Layers:

Forward:



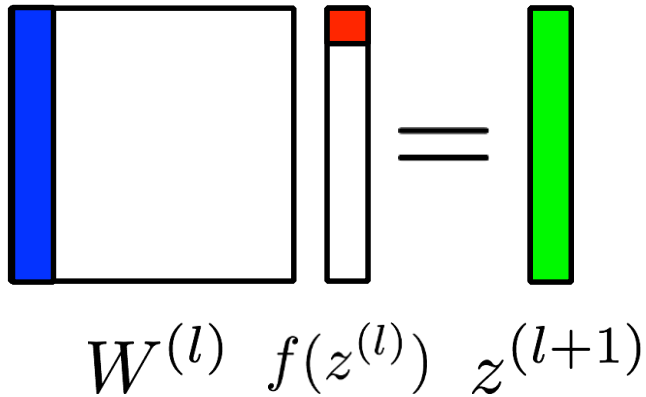
A measure how much an activation unit contributed to the loss

Backward:



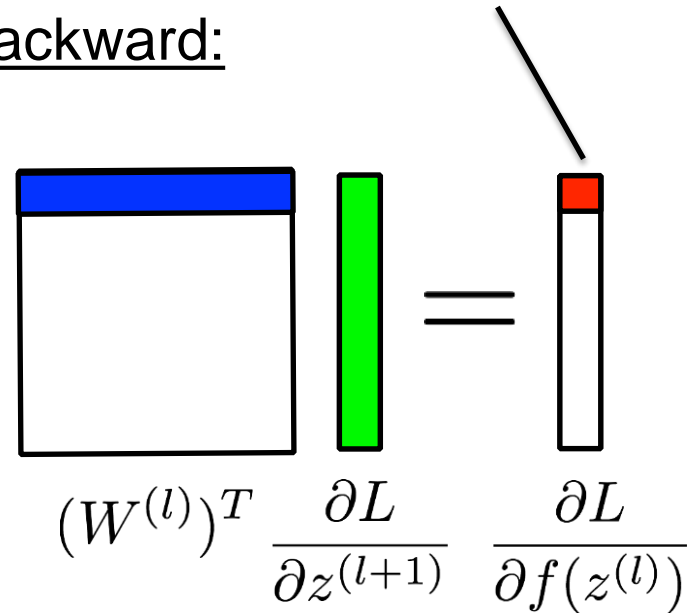
## Fully Connected Layers:

Forward:



A measure how much an activation unit contributed to the loss

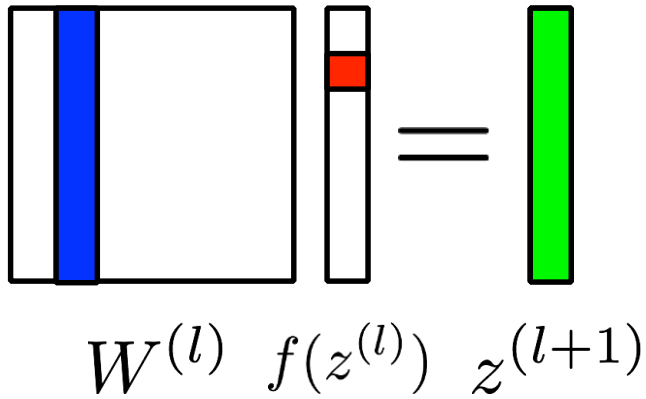
Backward:



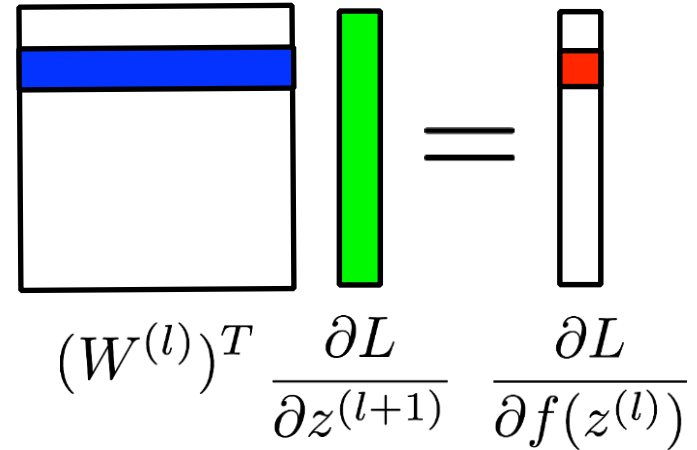


# Fully Connected Layers:

Forward:

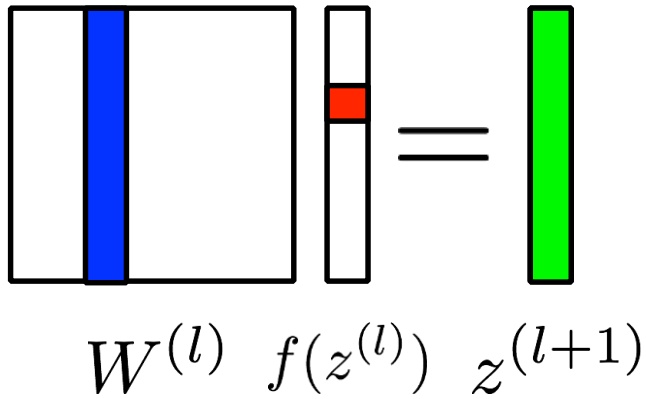


Backward:

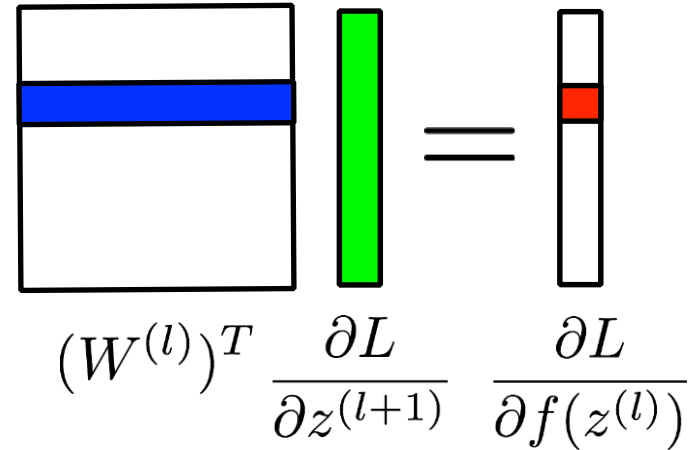


## Fully Connected Layers:

Forward:

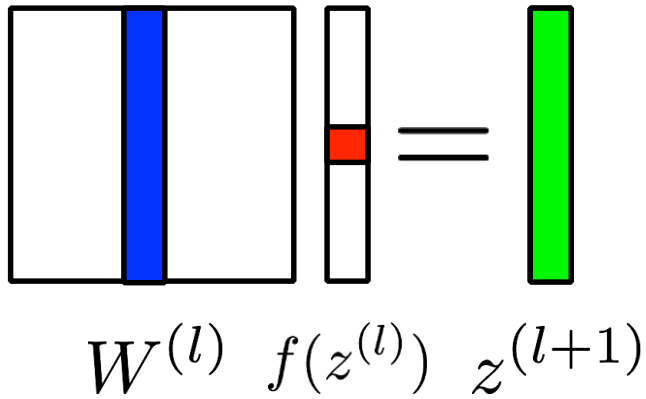


Backward:

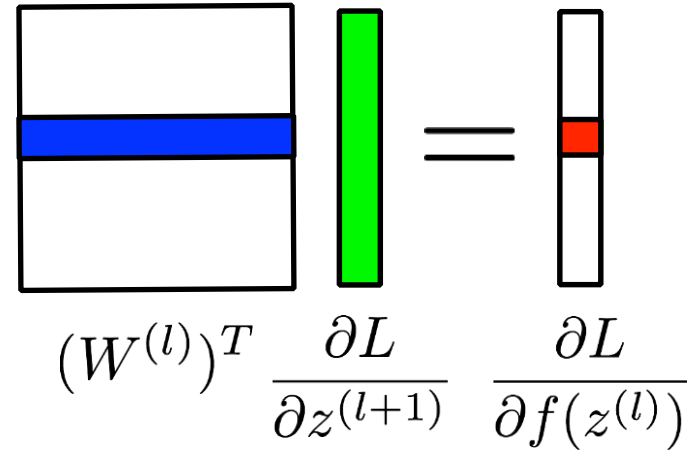


# Fully Connected Layers:

Forward:

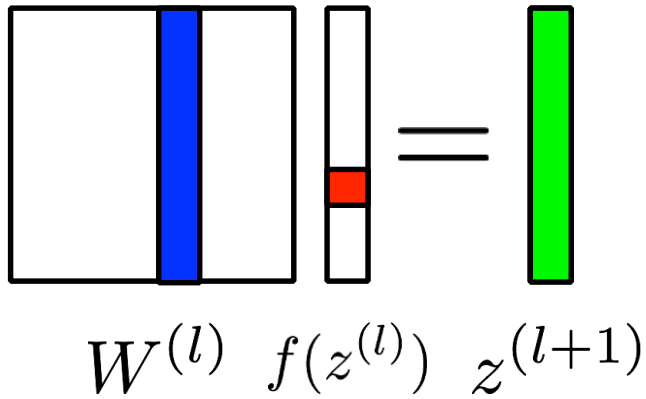


Backward:

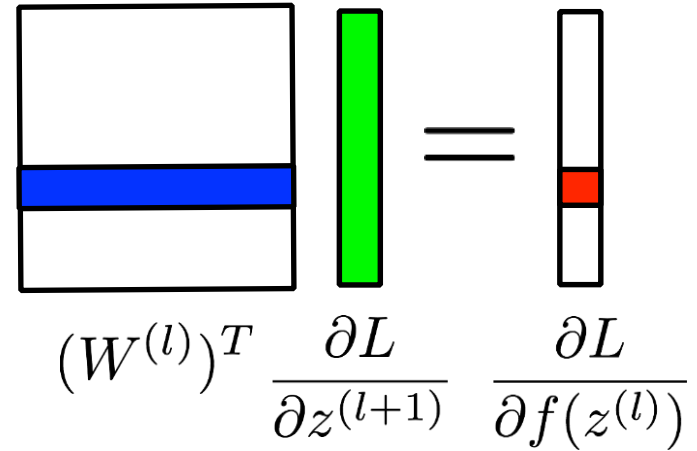


# Fully Connected Layers:

Forward:

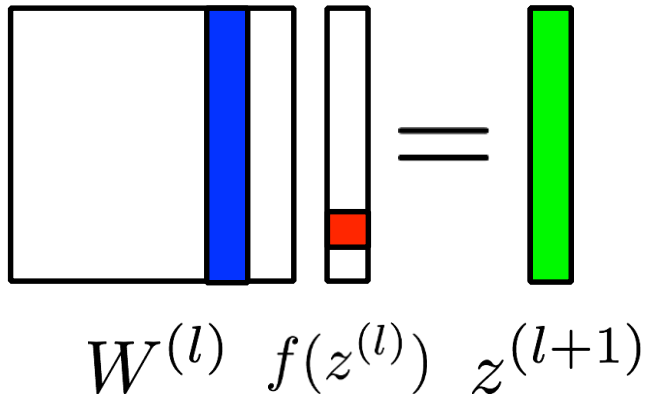


Backward:

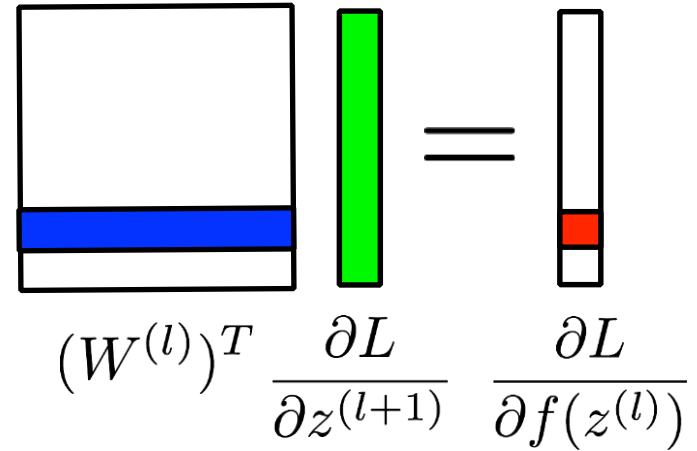


# Fully Connected Layers:

Forward:

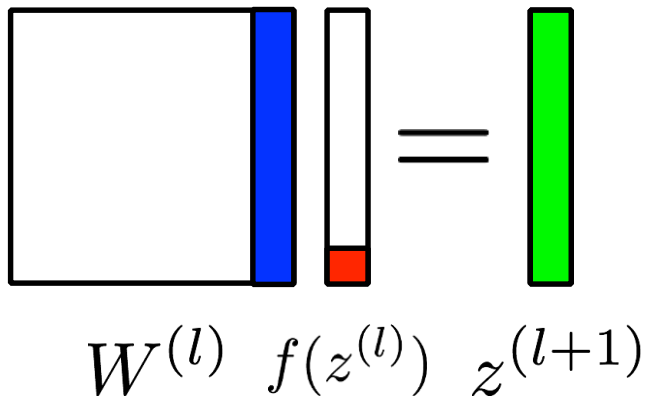


Backward:

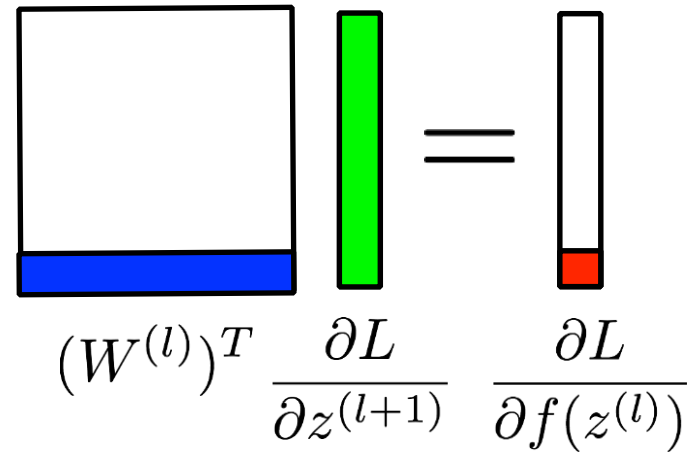


## Fully Connected Layers:

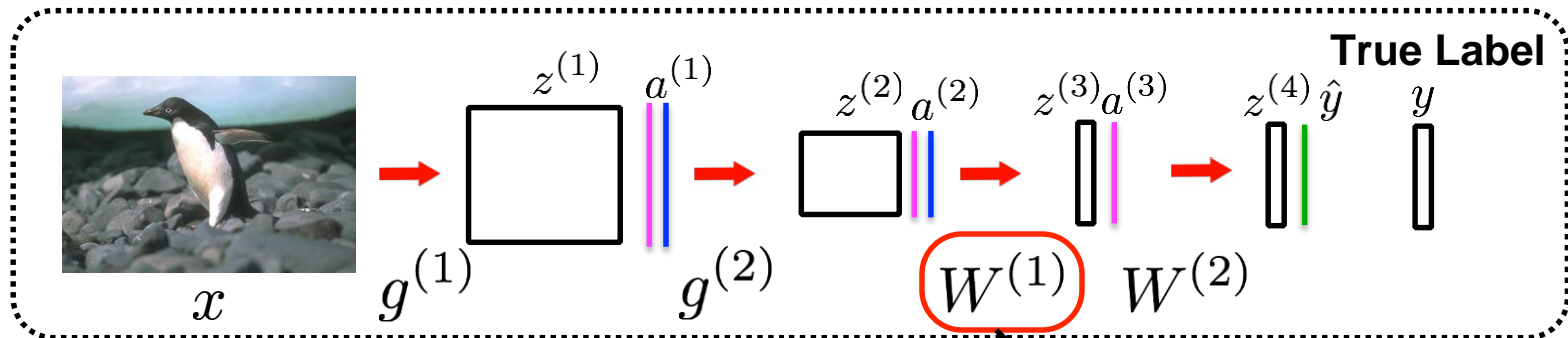
Forward:



Backward:



# Backpropagation

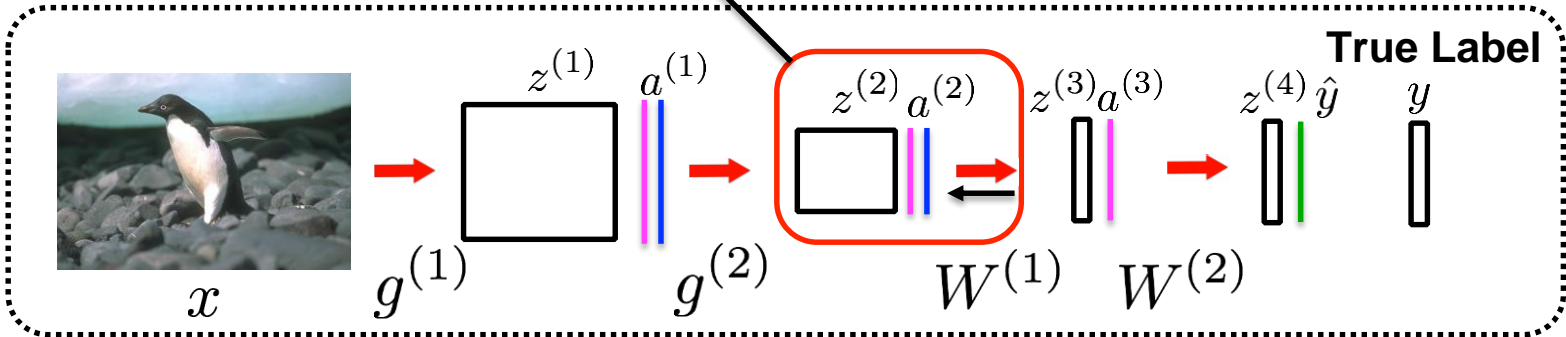


4. Compute the gradients to adjust the weights:  $\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial W^{(1)}}$  where  $z^{(3)} = W^{(1)} a^{(2)}$

# Backpropagation

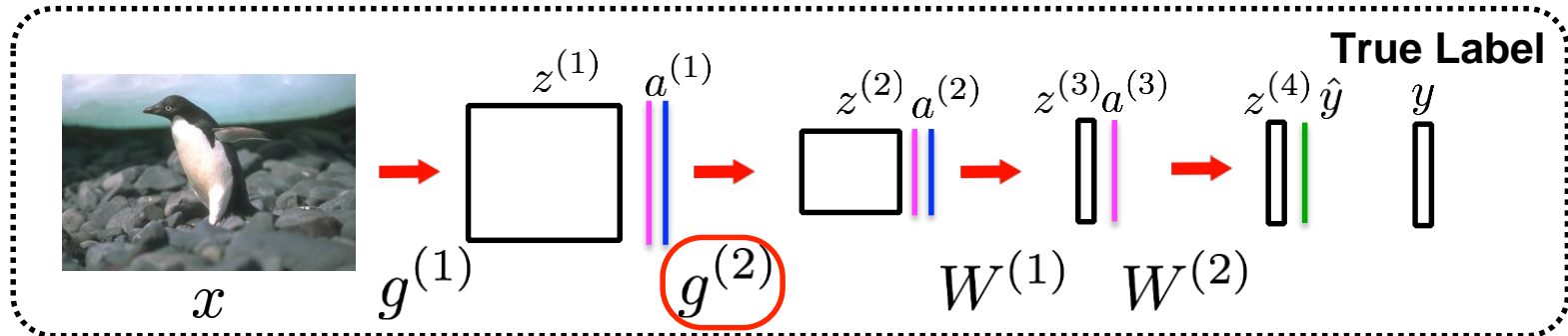
5. Backpropagate the gradients to previous layers:

$$\frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial f(z^{(2)})} \frac{\partial f(z^{(2)})}{\partial z^{(2)}} \quad \text{where } z^{(3)} = W^{(1)} a^{(2)}$$





# Backpropagation

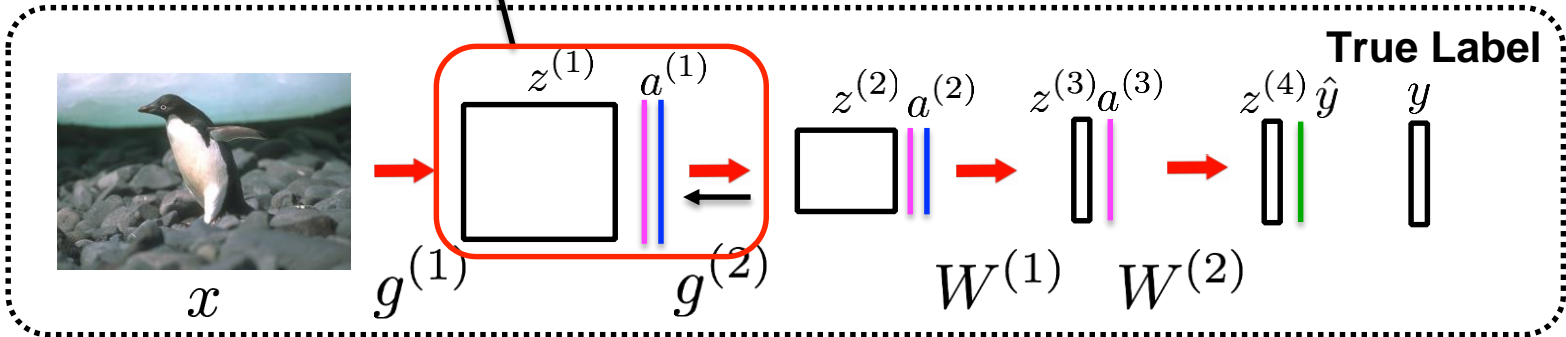


6. Compute the gradients to adjust the weights:  $\frac{\partial L}{\partial g^{(2)}} = \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial g^{(2)}}$  where  $z^{(2)} = g^{(2)} * a^{(1)}$

# Backpropagation

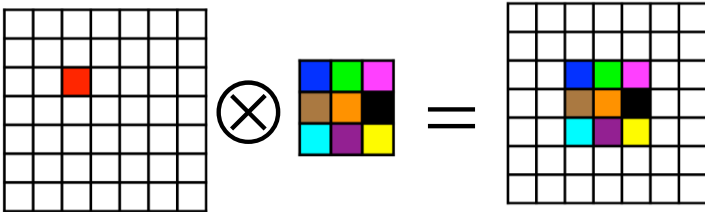
7. Backpropagate the gradients to previous layers:

$$\frac{\partial L}{\partial z^{(1)}} = \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial f(z^{(1)})} \frac{\partial f(z^{(1)})}{\partial z^{(1)}} \quad \text{where } z^{(2)} = g^{(2)} * a^{(1)}$$



# Convolutional Layers:

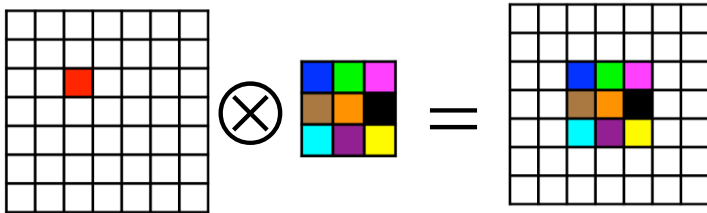
Forward:



$$a^{(\ell)} \otimes g^{(\ell)} = z^{(\ell+1)}$$

# Convolutional Layers:

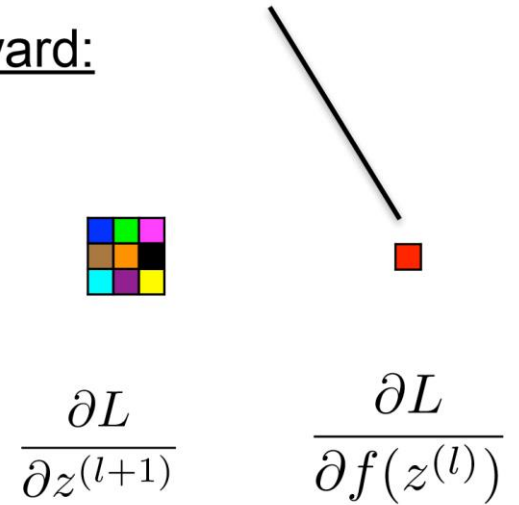
Forward:



$$a^{(l)} \otimes g^{(l)} = z^{(l+1)}$$

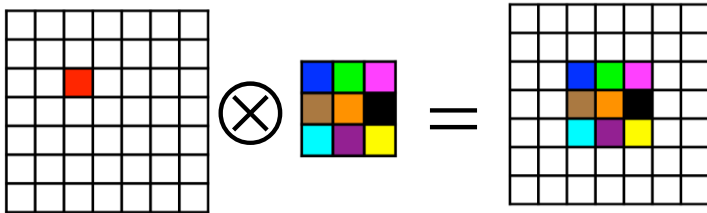
A measure how much an activation unit contributed to the loss

Backward:



# Convolutional Layers:

Forward:



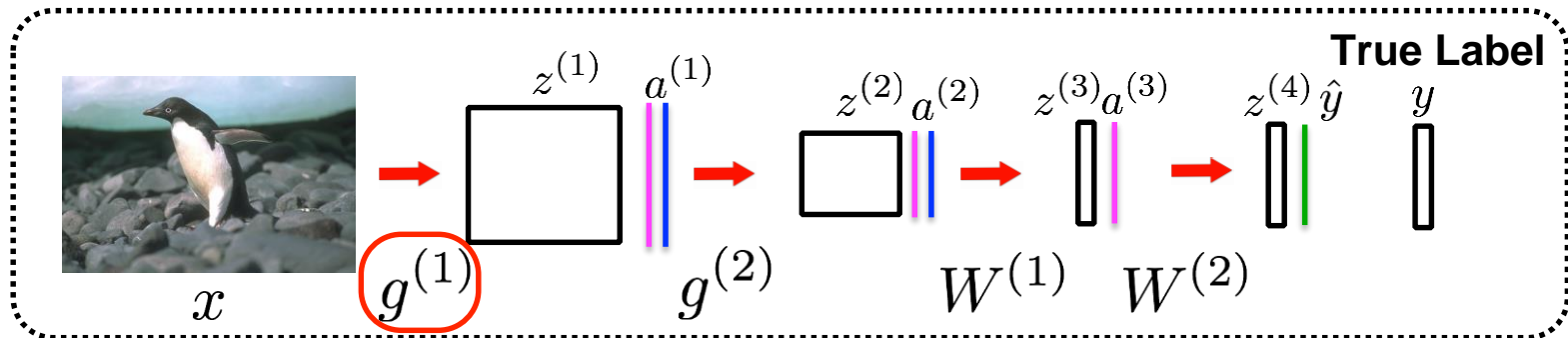
$$a^{(l)} \otimes g^{(l)} = z^{(l+1)}$$

Backward:



$$\text{sum} \left( g^{(l)} \odot \frac{\partial L}{\partial z^{(l+1)}} \right) = \frac{\partial L}{\partial f(z^{(l)})}$$

# Backpropagation



8. Compute the gradients to adjust the weights:  $\frac{\partial L}{\partial g^{(1)}} = \frac{\partial L}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial g^{(1)}}$  where  $z^{(1)} = g^{(1)} * x$

# Neural Network Learning

Google Deepmind DQN playing  
Atari Breakout

Setup:  
NVIDIA GTX 690  
i7-3770K - 16 GB RAM  
Ubuntu 16.04 LTS  
Google Deepmind DQN

# Neural Network Learning

Google Deepmind DQN playing  
Atari Breakout

Setup:  
NVIDIA GTX 690  
i7-3770K - 16 GB RAM  
Ubuntu 16.04 LTS  
Google Deepmind DQN

**How do we better understanding various properties  
behind the learning process in neural nets?**



# How to Improve NN Learning?

## Loss / Cost function:

- Loss function defines what we want the neural network to learn.

# How to Improve NN Learning?

## Loss / Cost function:

- Loss function defines what we want the neural network to learn.

**How do we select a good loss function?**

# How to Improve NN Learning?

## Loss / Cost function:

- Loss function defines what we want the neural network to learn.
- Humans learn best when they get feedback after being very wrong (e.g. a person learns to avoid scams after the first time he/she was scammed).

# How to Improve NN Learning?

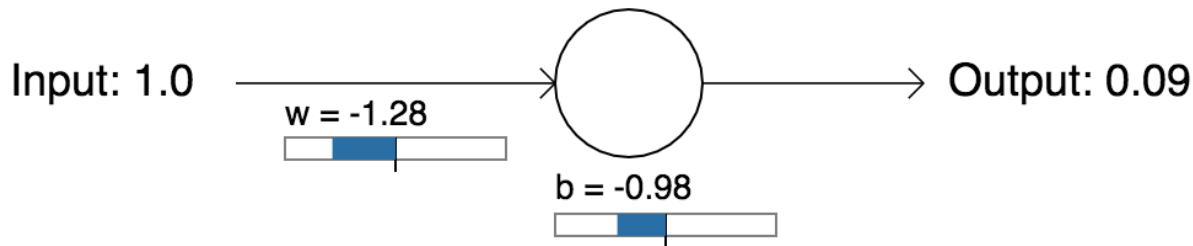
## Loss / Cost function:

- Loss function defines what we want the neural network to learn.
- Humans learn best when they get feedback after being very wrong (e.g. a person learns to avoid scams after the first time he/she was scammed).
- Does the same learning trend apply to neural networks? If not we want to design a loss function with such learning characteristics.

# How to Improve NN Learning?

## Loss / Cost function:

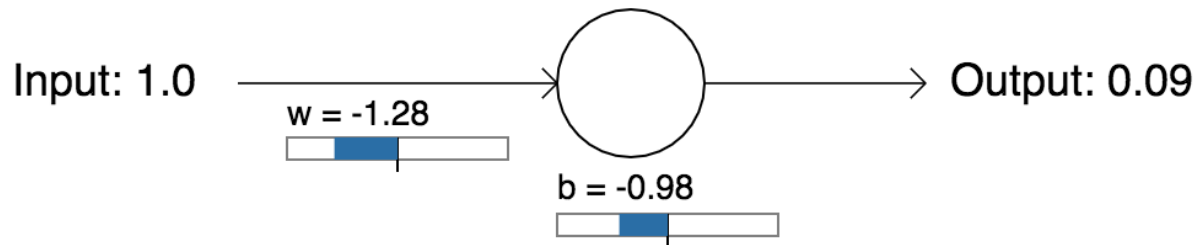
- Consider a neural network consisting of a single hidden neuron:



# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



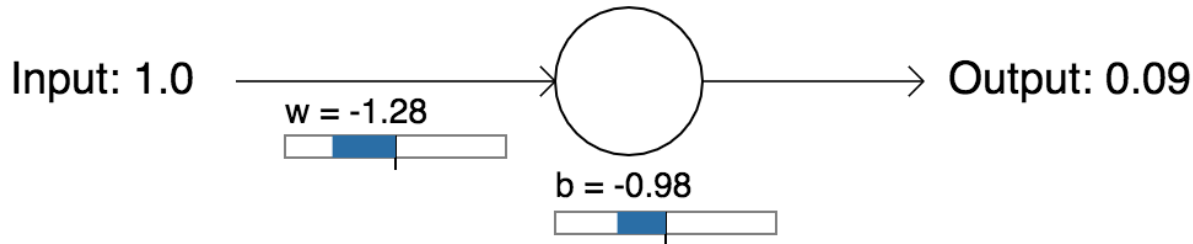
- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



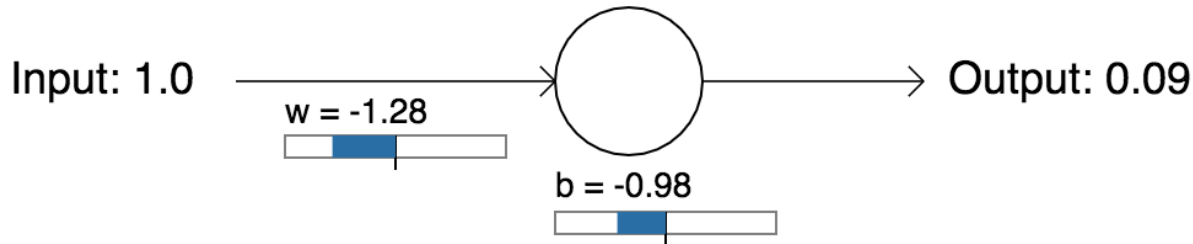
- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2} (\underbrace{\sigma(z)}_{\text{NN prediction}} - \underbrace{y}_{\text{ground truth}})^2$$

# How to Improve NN Learning?

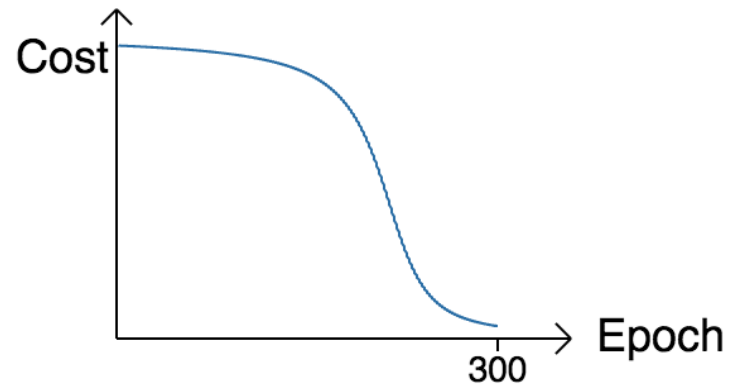
## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

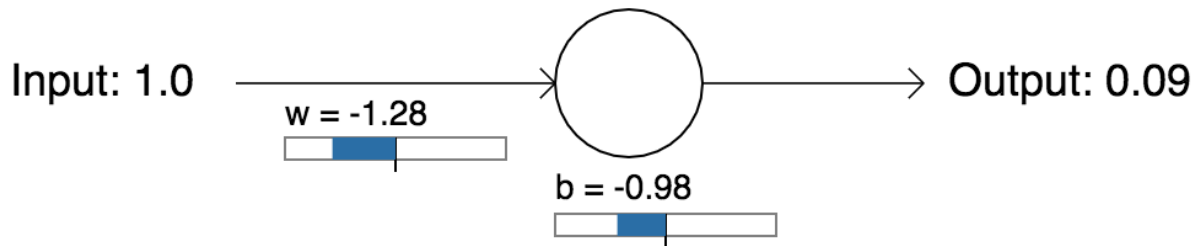




# How to Improve NN Learning?

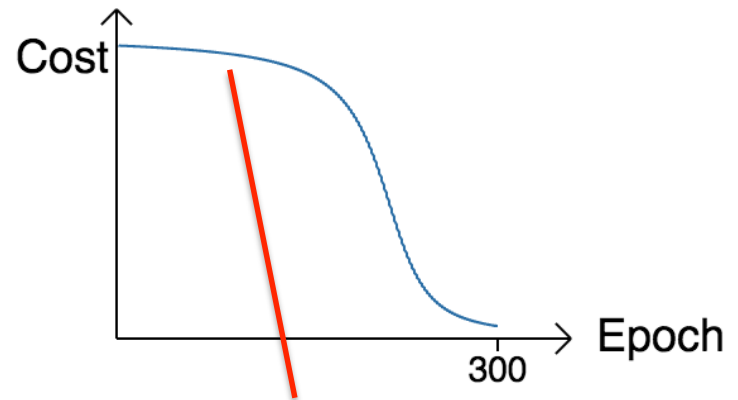
## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

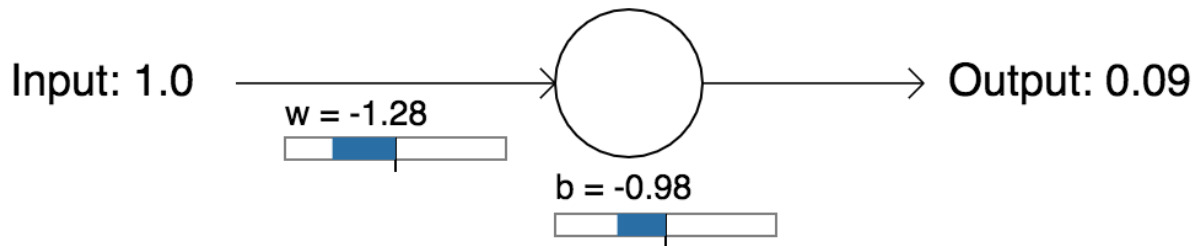


**why is the learning slow initially?**

# How to Improve NN Learning?

## Loss / Cost function:

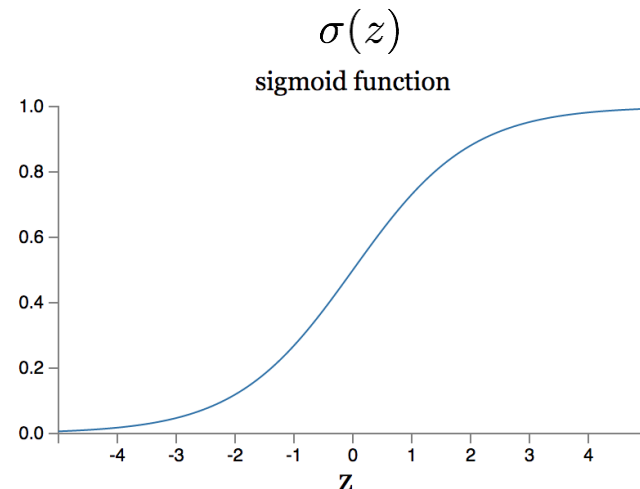
- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

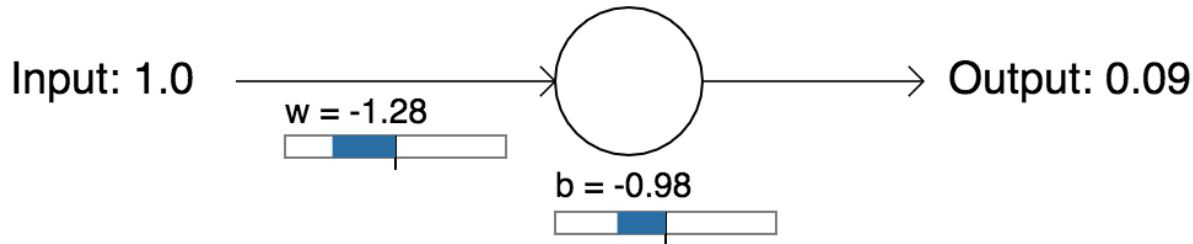
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)\sigma'(z)x$$



# How to Improve NN Learning?

## Loss / Cost function:

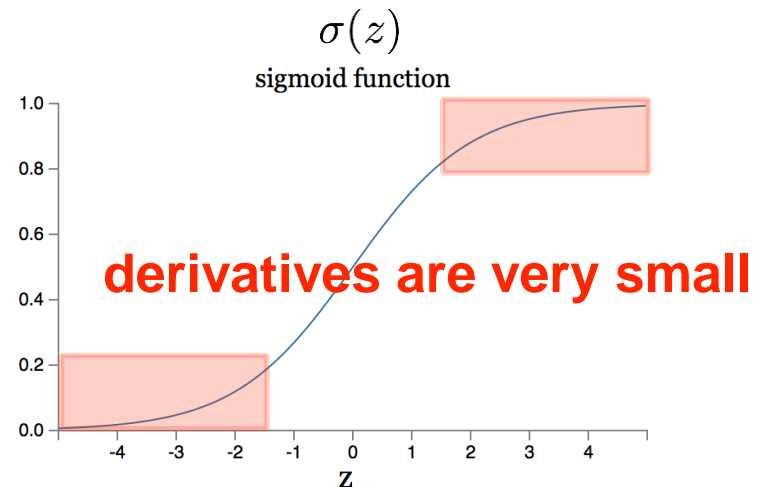
- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

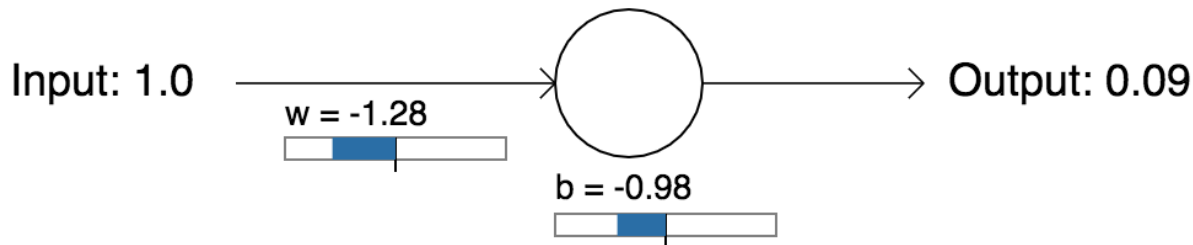
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)\sigma'(z)x$$



# How to Improve NN Learning?

## Loss / Cost function:

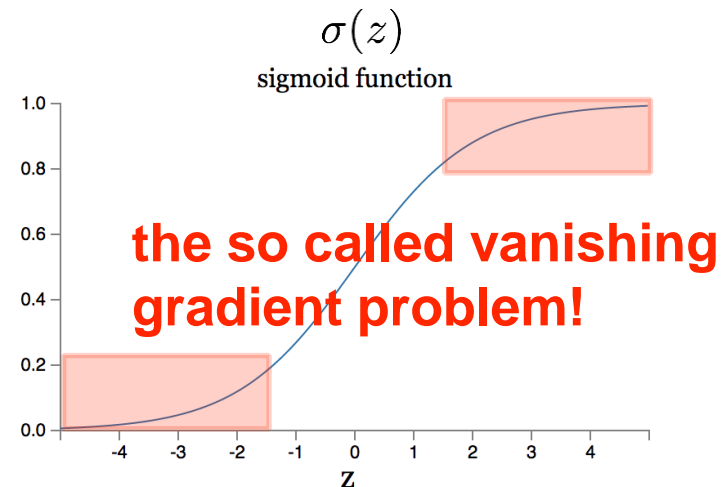
- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

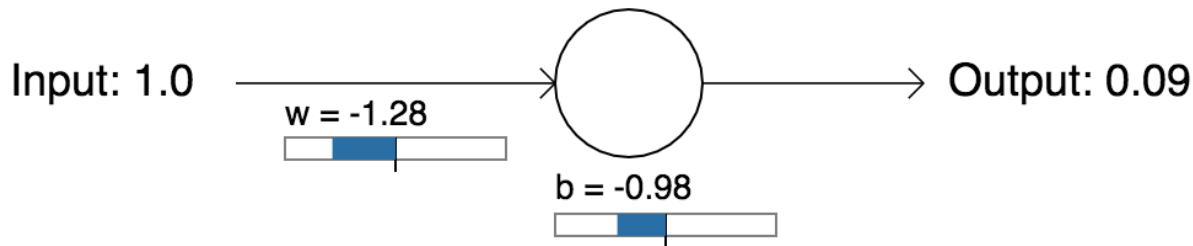
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)\sigma'(z)x$$



# How to Improve NN Learning?

## Loss / Cost function:

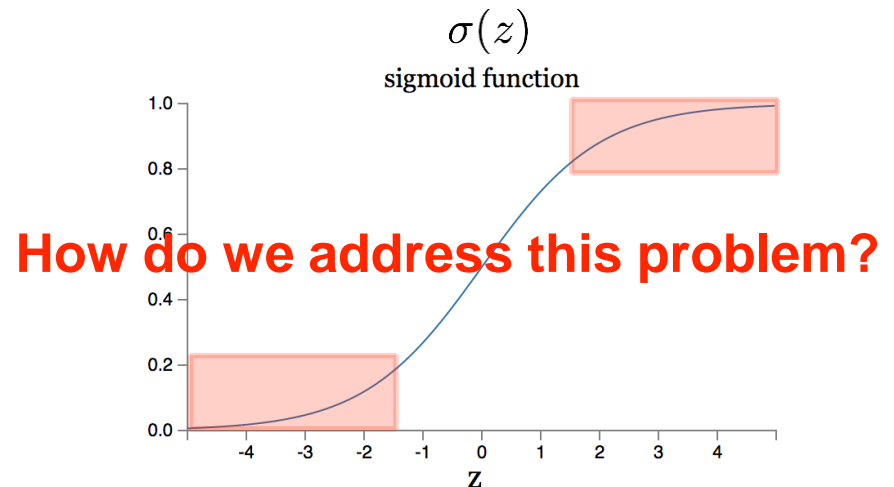
- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

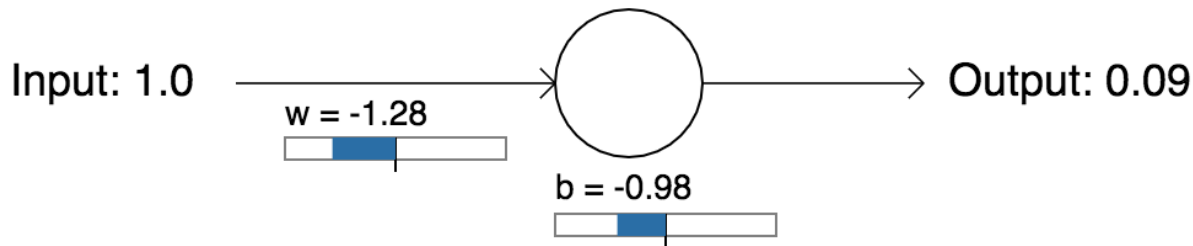
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)\sigma'(z)x$$



# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

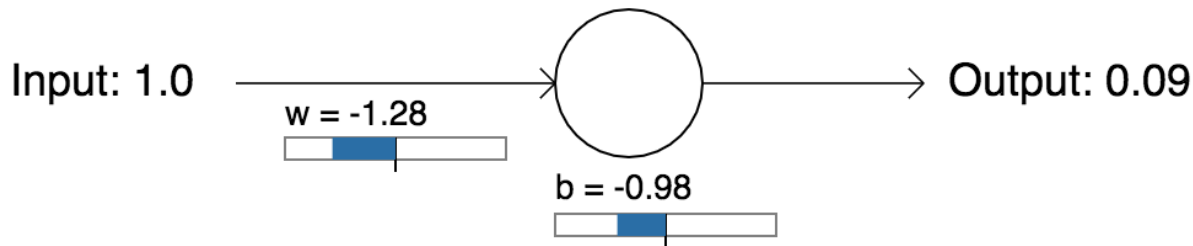
**We need a learning objective that wouldn't have a sigmoid derivative in the gradient.**

$$\frac{\partial L}{\partial w} = (\sigma(z) - y) \underline{\sigma'(z)} x$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- First, let's examine a commonly used L2 loss objective:

$$L = \frac{1}{2}(\sigma(z) - y)^2$$

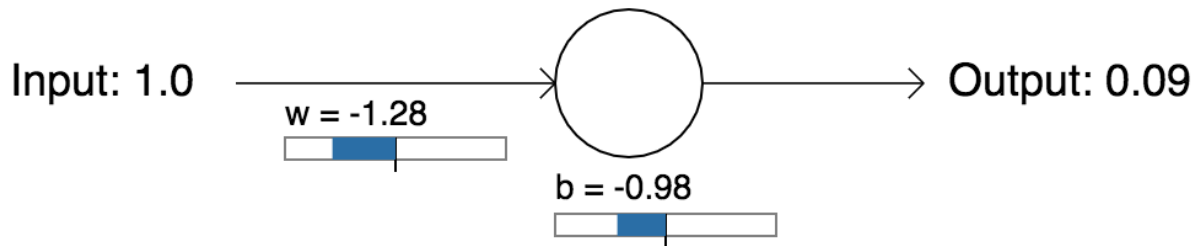
**We want the following gradient:**

$$\frac{\partial L}{\partial w} = (\sigma(z) - y)\sigma'(z)x \rightarrow \boxed{\frac{\partial L}{\partial w} = (\sigma(z) - y)x}$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



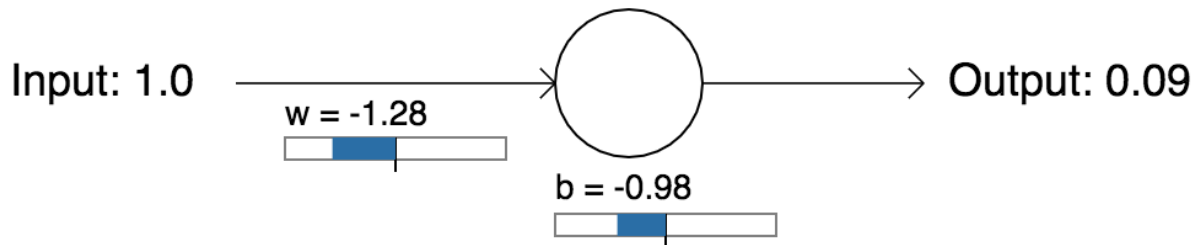
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$



# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



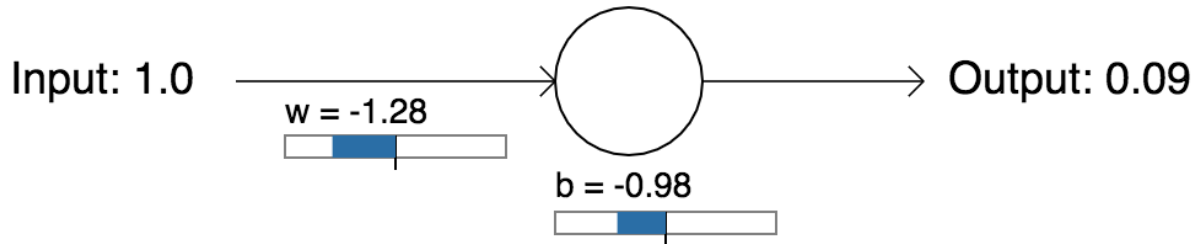
$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} \quad \text{where } \frac{\partial a}{\partial z} = a(1 - a) \quad \& \quad \frac{\partial z}{\partial w} = x$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$

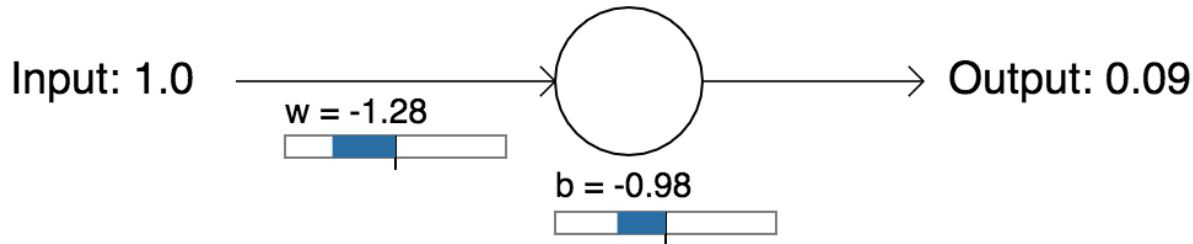
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} \quad \text{where } \frac{\partial a}{\partial z} = a(1 - a) \quad \& \quad \frac{\partial z}{\partial w} = x$$

$$(a - y)x = \frac{\partial L}{\partial a} a(1 - a)x$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$

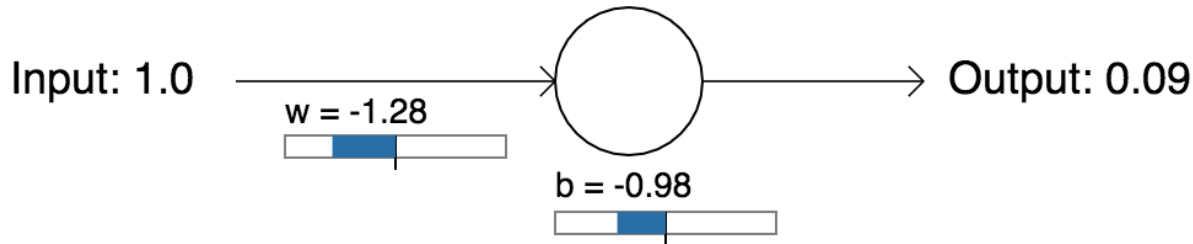
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} \quad \text{where } \frac{\partial a}{\partial z} = a(1 - a) \quad \& \quad \frac{\partial z}{\partial w} = x$$

$$(a - y)x = \frac{\partial L}{\partial a} a(1 - a)x \rightarrow \boxed{\frac{\partial L}{\partial a} = \frac{(a - y)}{a(1 - a)}}$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} \quad \text{where } \frac{\partial a}{\partial z} = a(1 - a) \quad \& \quad \frac{\partial z}{\partial w} = x$$

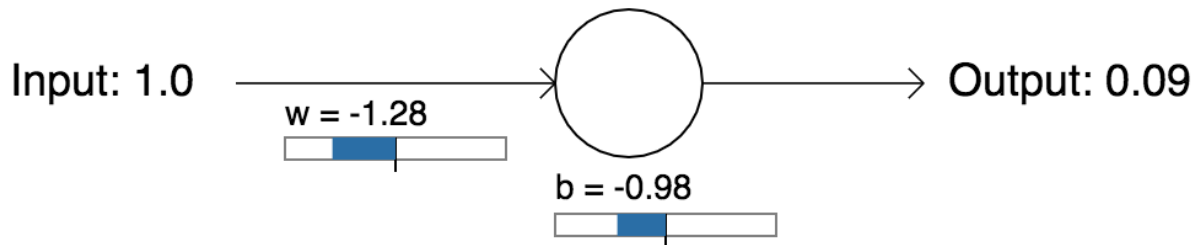
$$(a - y)x = \frac{\partial L}{\partial a} a(1 - a)x \rightarrow \frac{\partial L}{\partial a} = \frac{(a - y)}{a(1 - a)}$$

$$L = -(y \log a + (1 - y) \log (1 - a))$$

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



$$\frac{\partial L}{\partial w} = (\sigma(z) - y)x \quad \text{Let } a = \sigma(z)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} \quad \text{where } \frac{\partial a}{\partial z} = a(1 - a) \quad \& \quad \frac{\partial z}{\partial w} = x$$

$$(a - y)x = \frac{\partial L}{\partial a} a(1 - a)x \rightarrow \frac{\partial L}{\partial a} = \frac{(a - y)}{a(1 - a)}$$

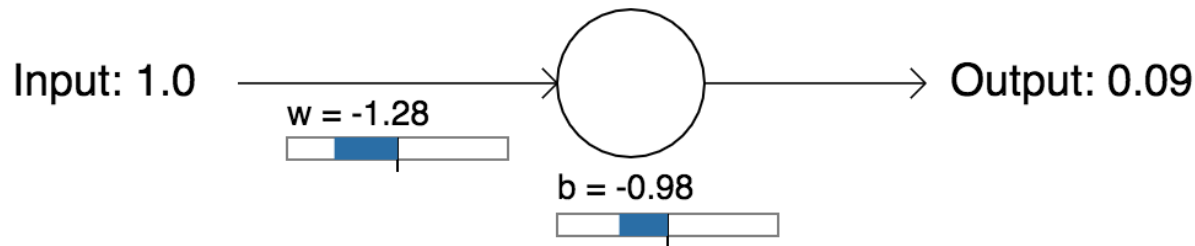
$$L = -(y \log a + (1 - y) \log (1 - a))$$

**We just derived a cross-entropy loss function!**

# How to Improve NN Learning?

## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



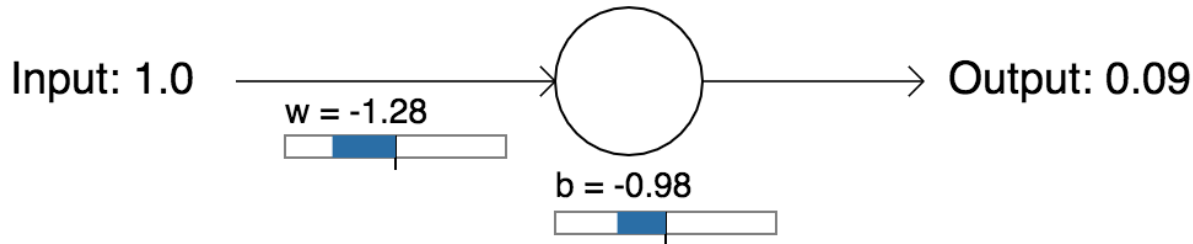
- Now let's examine a cross-entropy loss objective:

$$L = -(y \log a + (1 - y) \log (1 - a))$$

# How to Improve NN Learning?

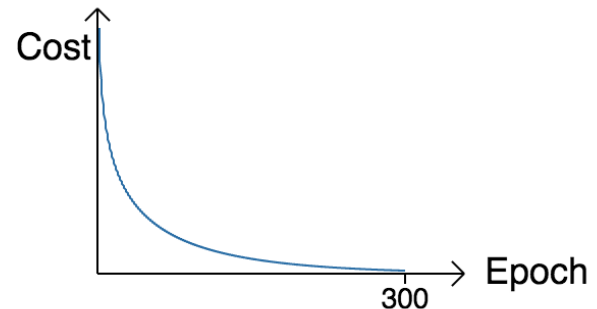
## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- Now let's examine a cross-entropy loss objective:

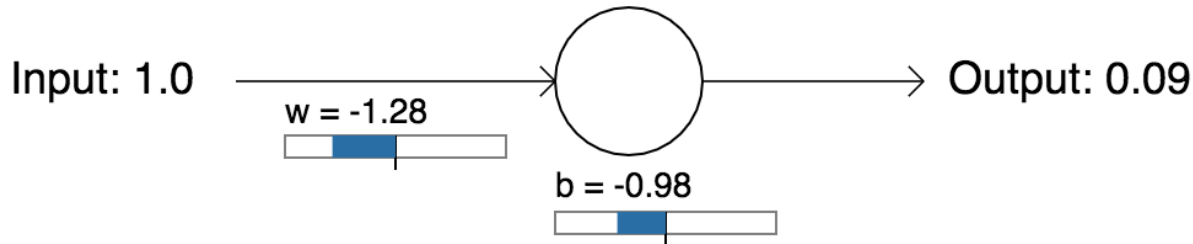
$$L = -(y \log a + (1 - y) \log (1 - a))$$



# How to Improve NN Learning?

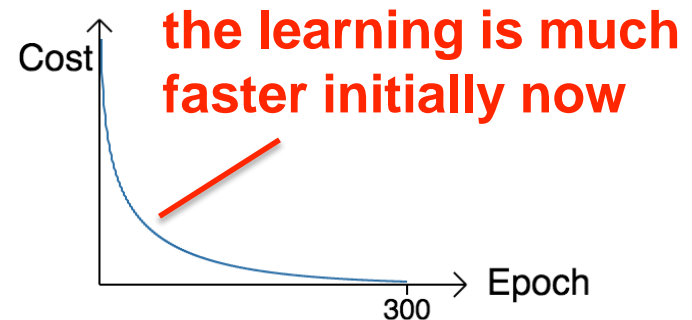
## Loss / Cost function:

- Consider a neural network consisting of a single hidden neuron:



- Now let's examine a cross-entropy loss objective:

$$L = -(y \log a + (1 - y) \log (1 - a))$$

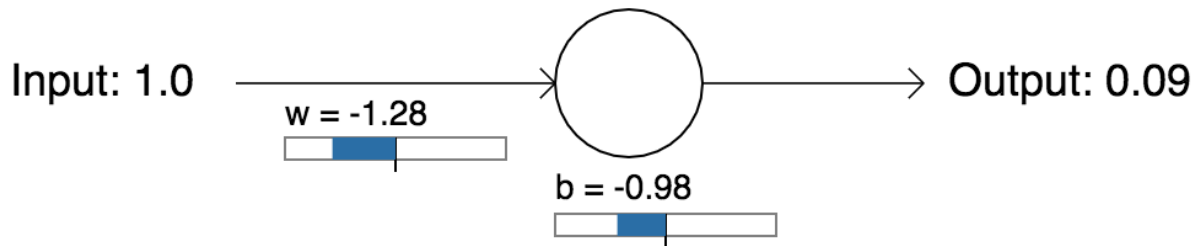




# How to Improve NN Learning?

## Loss / Cost function:

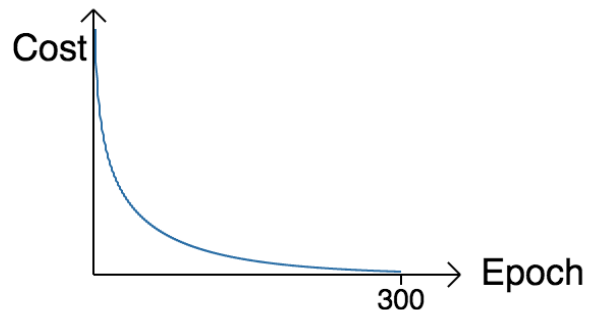
- Consider a neural network consisting of a single hidden neuron:



- Now let's examine a cross-entropy loss objective:

$$L = -(y \log a + (1 - y) \log (1 - a))$$

$$\frac{\partial L}{\partial w} = x(a - y)$$



**no sigmoid derivatives in the gradient equation**

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$L = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad \text{where} \quad \hat{y}_i = \frac{\exp(z_i^{(4)})}{\sum_{j=1}^K \exp(z_j^{(4)})}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$L = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad \text{where} \quad \hat{y}_i = \frac{\exp(z_i^{(4)})}{\sum_{j=1}^K \exp(z_j^{(4)})}$$

$$\frac{\partial L}{\partial z_i^{(4)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$L = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad \text{where} \quad \hat{y}_i = \frac{\exp(z_i^{(4)})}{\sum_{j=1}^K \exp(z_j^{(4)})}$$

$$\frac{\partial L}{\partial z_i^{(4)}} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}}$$

$$\frac{\partial L}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$L = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad \text{where} \quad \hat{y}_i = \frac{\exp(z_i^{(4)})}{\sum_{j=1}^K \exp(z_j^{(4)})}$$

$$\frac{\partial L}{\partial z_i^{(4)}} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}}$$

$$\frac{\partial L}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i}$$

$$\frac{\partial \hat{y}_i}{\partial z_j^{(4)}} = \begin{cases} \hat{y}_i(1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \hat{y}_j, & \text{if } i \neq j \end{cases}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$\frac{\partial L}{\partial z_i^{(4)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$\begin{aligned}\frac{\partial L}{\partial z_i^{(4)}} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}} \\ &= \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}} + \sum_{i \neq j} \frac{\partial L}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial z_i^{(4)}}\end{aligned}$$

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$\begin{aligned}\frac{\partial L}{\partial z_i^{(4)}} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}} \\ &= \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}} + \sum_{i \neq j} \frac{\partial L}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial z_i^{(4)}} \\ &= \hat{y}_i - y_i\end{aligned}$$



# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?

$$\begin{aligned}\frac{\partial L}{\partial z_i^{(4)}} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}} \\ &= \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}} + \sum_{i \neq j} \frac{\partial L}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial z_i^{(4)}} \\ &= \hat{y}_i - y_i\end{aligned}$$

- **No sigmoid derivatives in the gradient!**

# How to Improve NN Learning?

## Loss / Cost function:

- How about the softmax loss function?




$$\begin{aligned}\frac{\partial L}{\partial z_i^{(4)}} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i^{(4)}} \\ &= \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(4)}} + \sum_{i \neq j} \frac{\partial L}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial z_i^{(4)}} \\ &= \hat{y}_i - y_i\end{aligned}$$

- **No sigmoid derivatives in the gradient!**
- **Therefore, learning shouldn't be slowed down.**

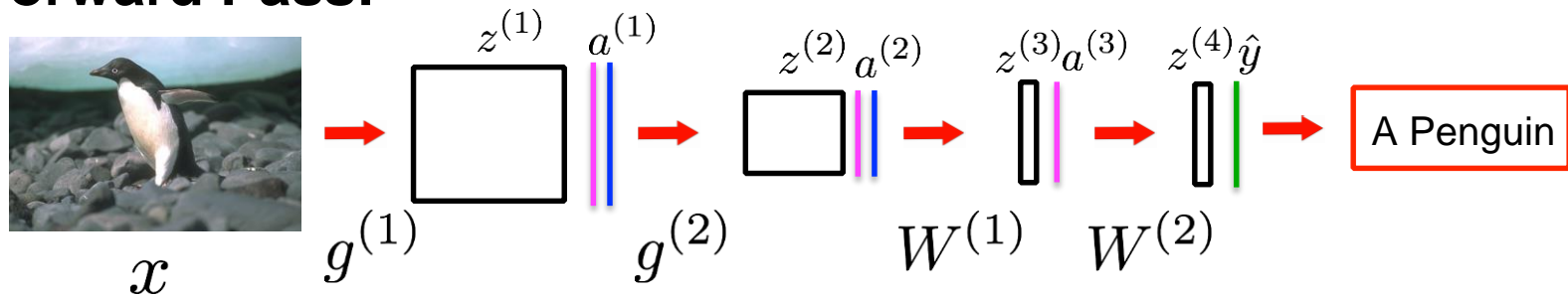
# Vanishing Gradients

## Notation:

 - convolutional layer output     - fully connected layer output

 - max pooling layer     - sigmoid function  $f$      - softmax function

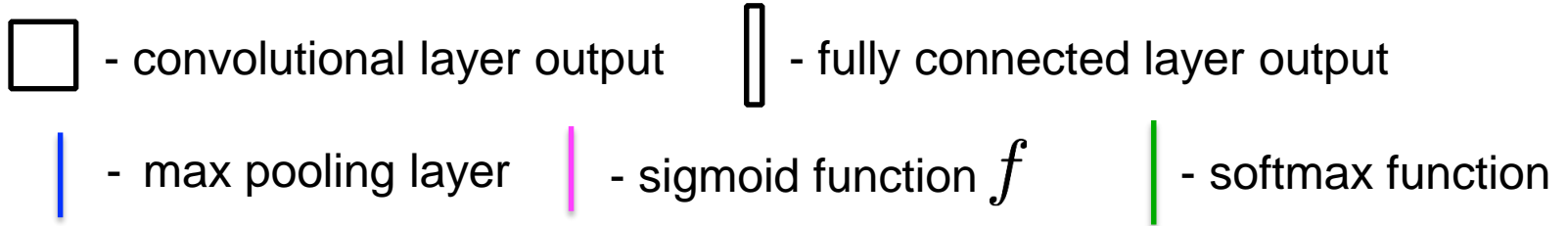
## Forward Pass:



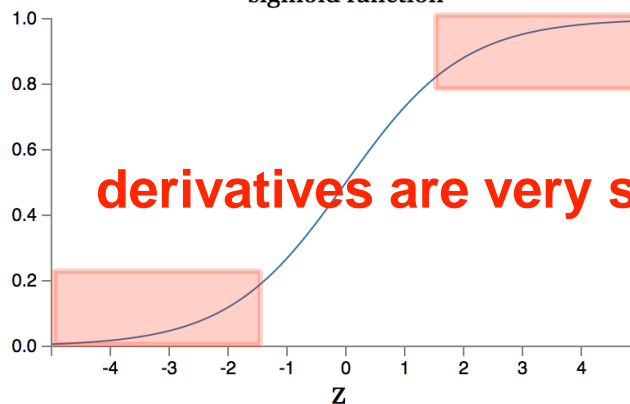
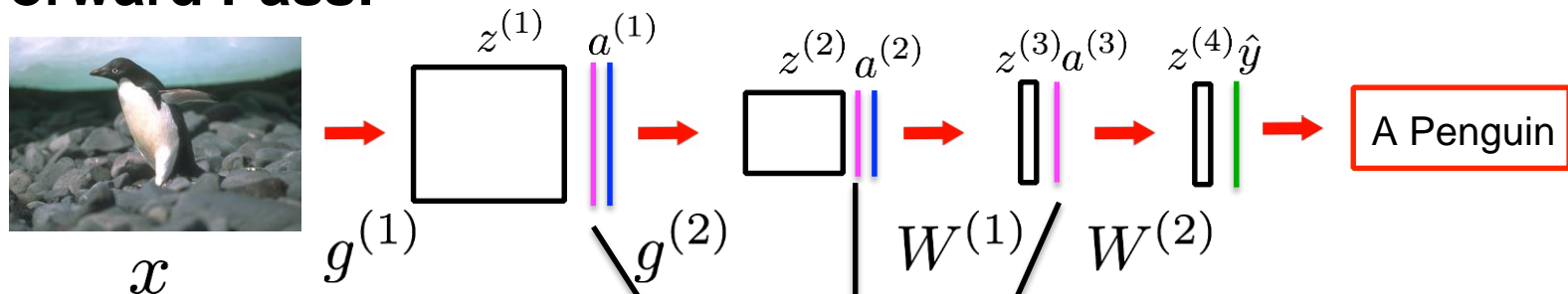
- Did we solve the vanishing gradient problem?

# Vanishing Gradients

## Notation:

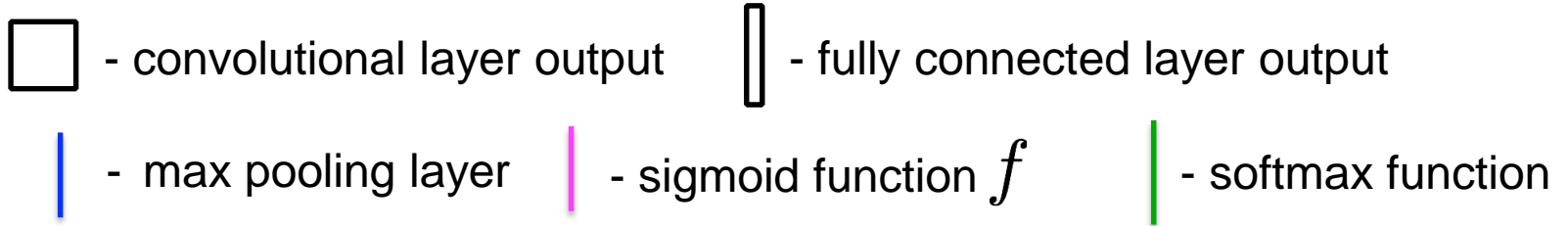


## Forward Pass:

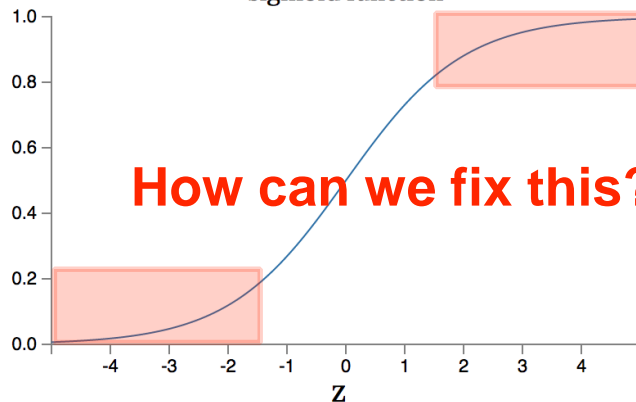
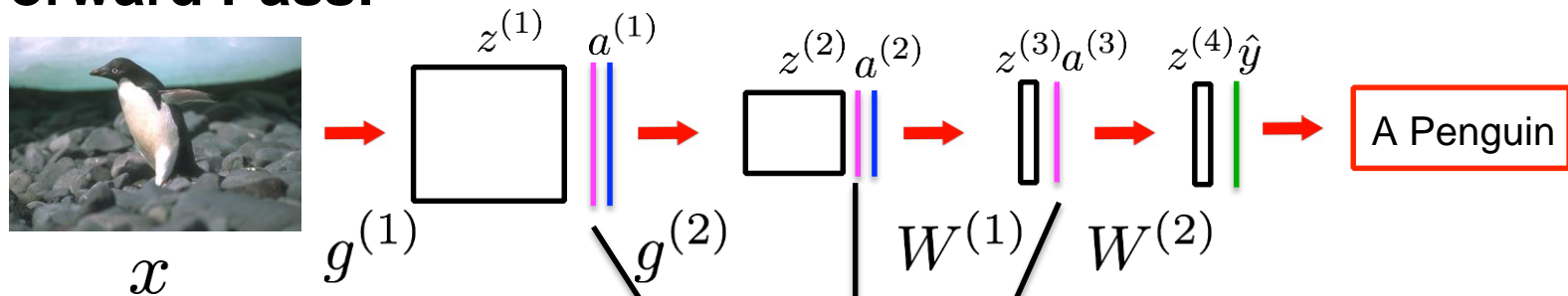


# Vanishing Gradients

## Notation:

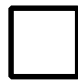






## Forward Pass:

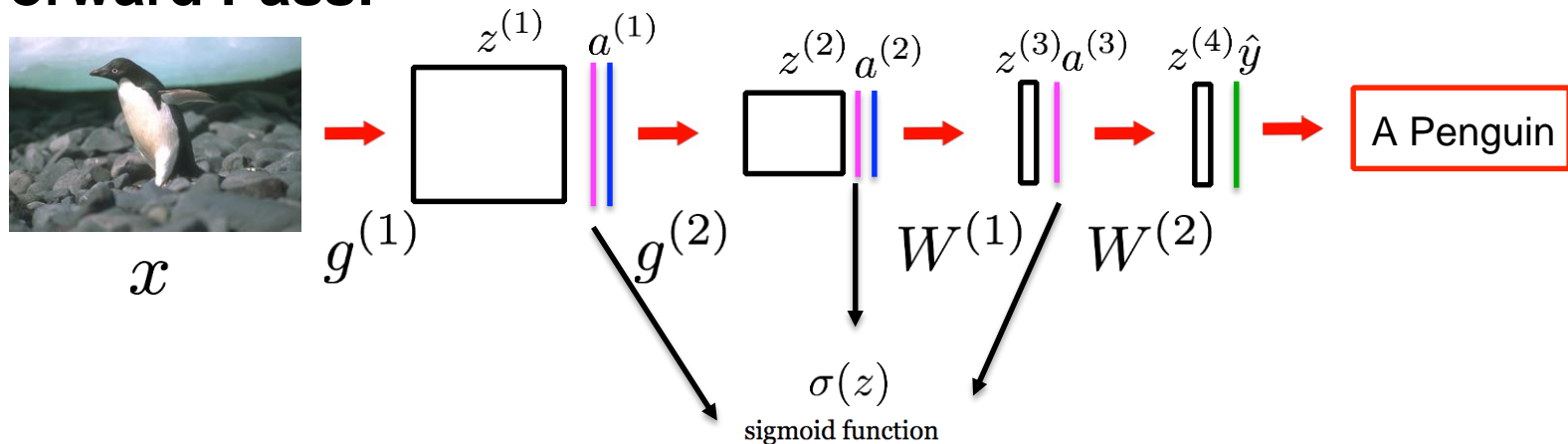


# Vanishing Gradients

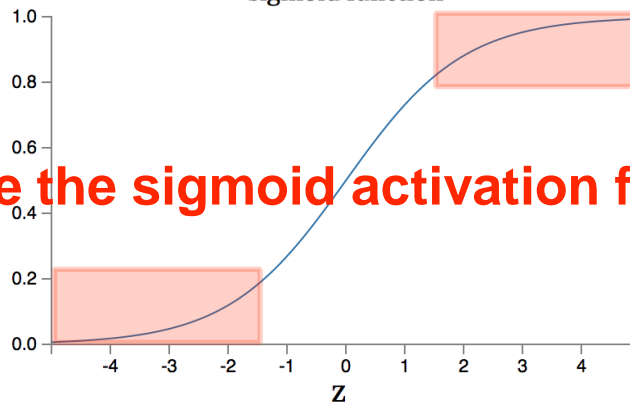
## Notation:

-  - convolutional layer output
-  - fully connected layer output
-  - max pooling layer
-  - sigmoid function  $f$
-  - softmax function

## Forward Pass:



**Replace the sigmoid activation function!**



# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- We want the activation function to be differentiable.
- We want an activation function that eliminates the vanishing gradient problem.

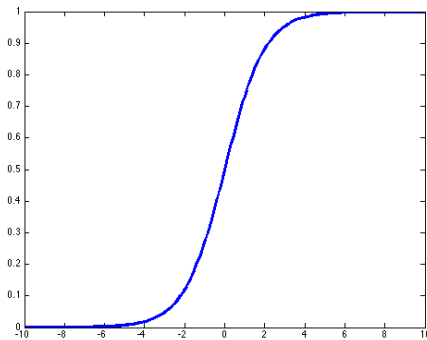
# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- We want the activation function to be differentiable.
- We want an activation function that eliminates the vanishing gradient problem.

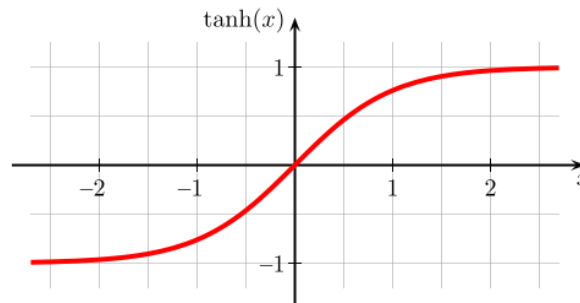
### Sigmoid Function

$$f(z) = \frac{1}{1 + \exp(-z)}$$



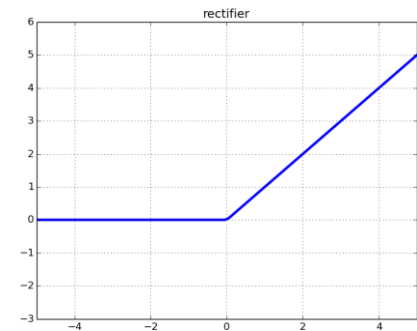
### Hyperbolic Tangent

$$f(z) = \tanh(z)$$



### RELU

$$f(z) = \max(0, z)$$





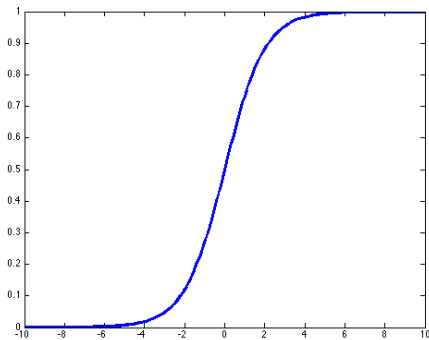
# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- We want the activation function to be differentiable.
- We want an activation function that eliminates the vanishing gradient problem.

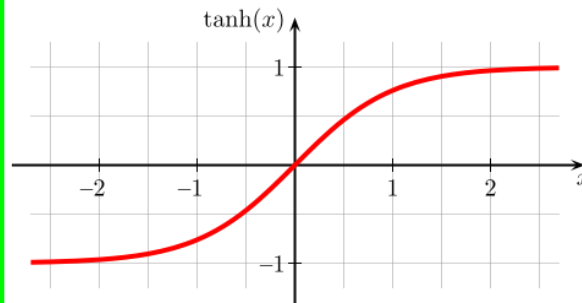
### Sigmoid Function

$$f(z) = \frac{1}{1 + \exp(-z)}$$



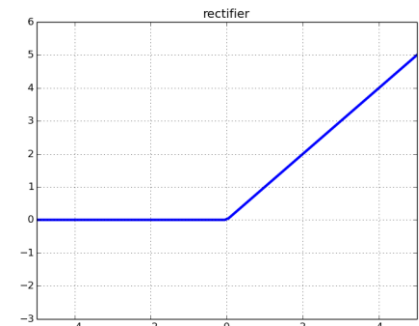
### Hyperbolic Tangent

$$f(z) = \tanh(z)$$



### RELU

$$f(z) = \max(0, z)$$



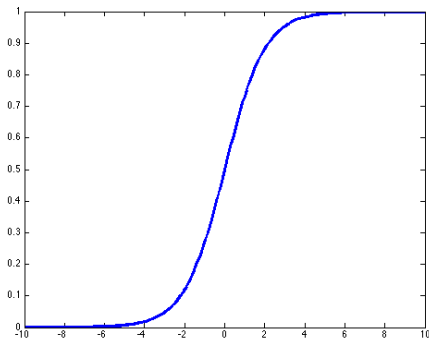
# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- **We want the activation function to be differentiable.**
- We want an activation function that eliminates the vanishing gradient problem.

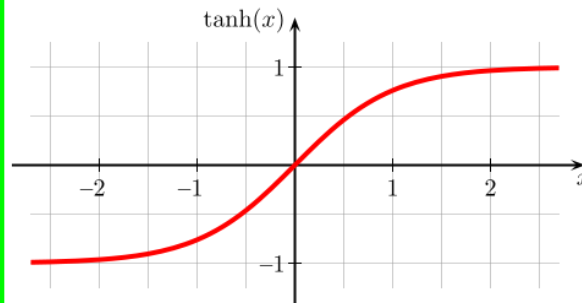
### Sigmoid Function

$$f(z) = \frac{1}{1 + \exp(-z)}$$



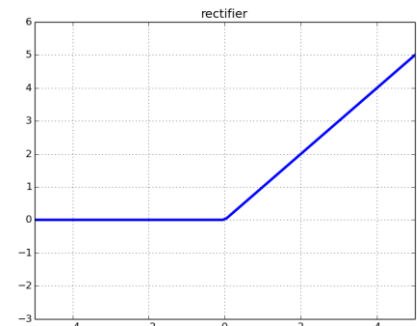
### Hyperbolic Tangent

$$f(z) = \tanh(z)$$



### RELU

$$f(z) = \max(0, z)$$



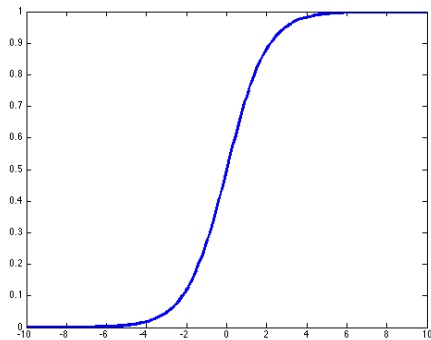
# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- We want the activation function to be differentiable.
- We want an activation function that eliminates the vanishing gradient problem.

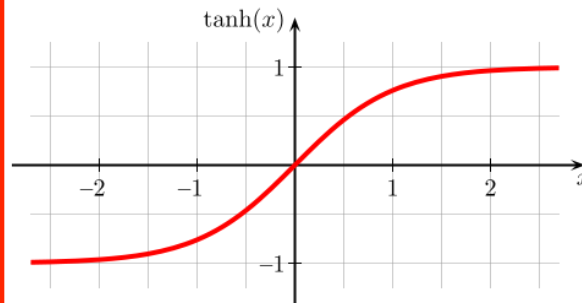
### Sigmoid Function

$$f(z) = \frac{1}{1 + \exp(-z)}$$



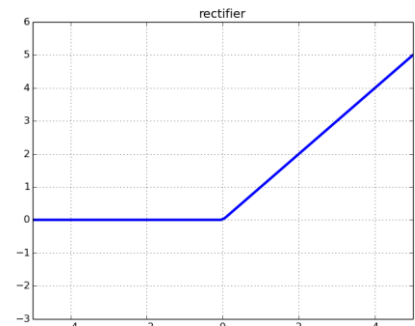
### Hyperbolic Tangent

$$f(z) = \tanh(z)$$



### RELU

$$f(z) = \max(0, z)$$



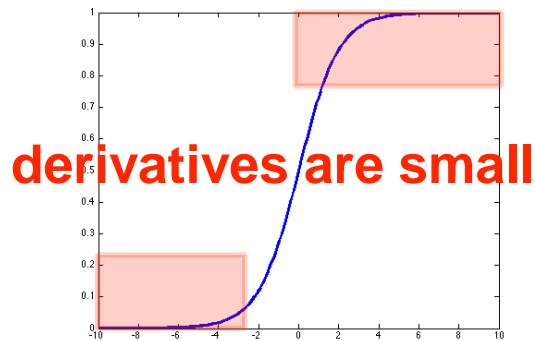
# Vanishing Gradients

## Activation function requirements:

- We want the activation function to be non-linear.
- We want the activation function to be differentiable.
- We want an activation function that eliminates the vanishing gradient problem.

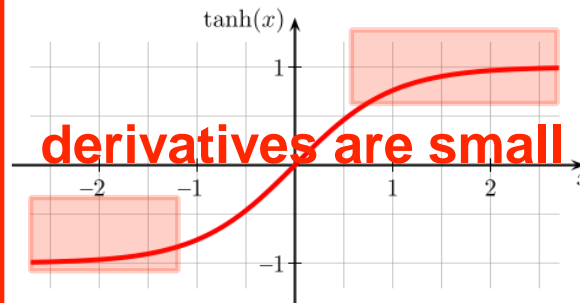
### Sigmoid Function

$$f(z) = \frac{1}{1 + \exp(-z)}$$



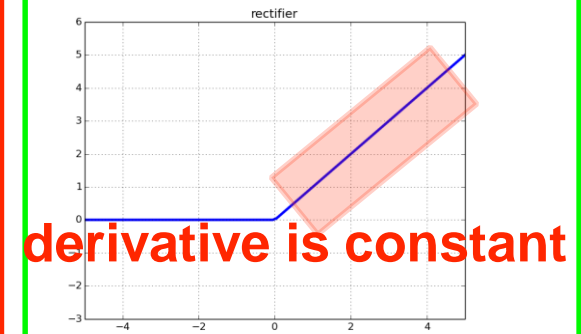
### Hyperbolic Tangent

$$f(z) = \tanh(z)$$



### RELU

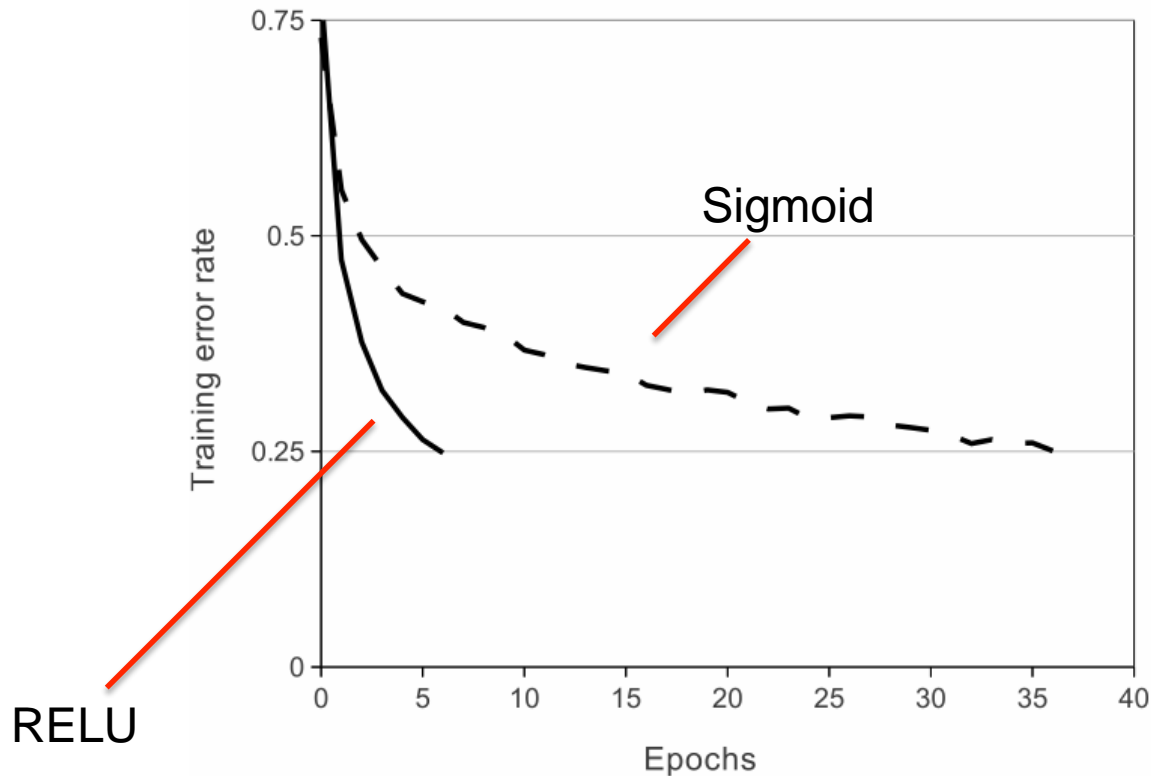
$$f(z) = \max(0, z)$$



# Vanishing Gradients

## Learning Speed:

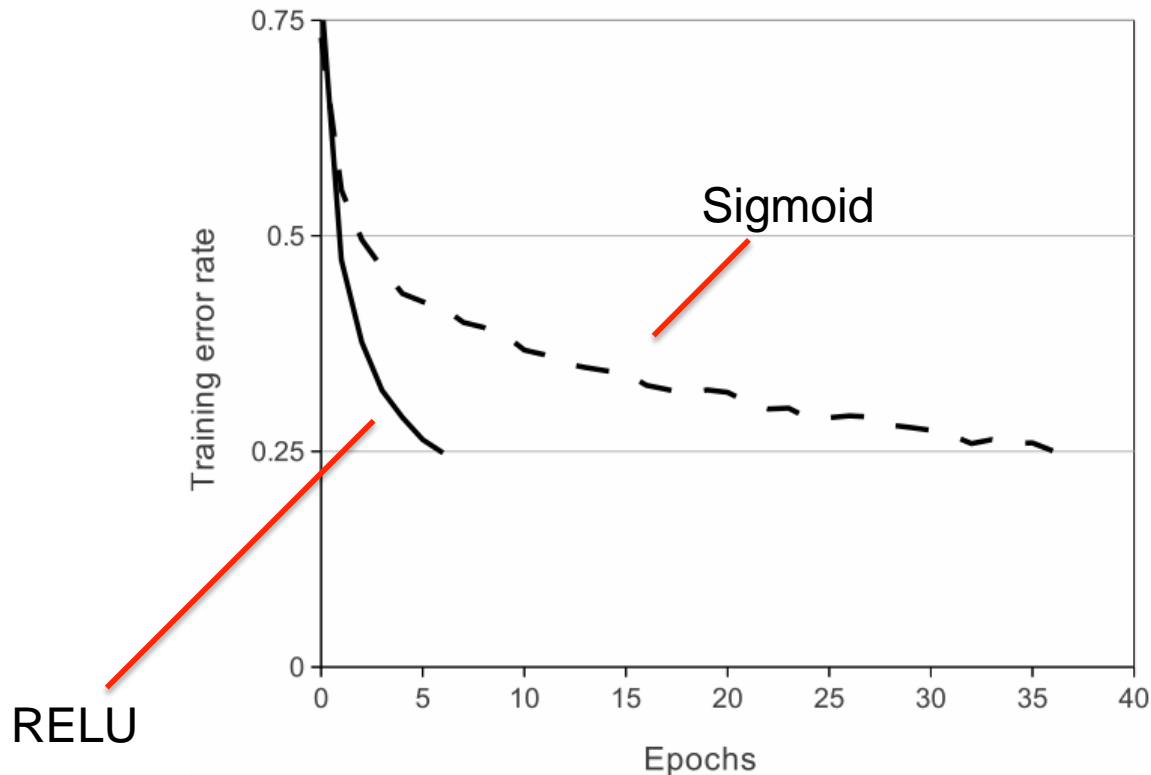
- The network that uses a RELU activation function learns significantly faster.



# Vanishing Gradients

## Learning Speed:

- The network that uses a RELU activation function learns significantly faster.

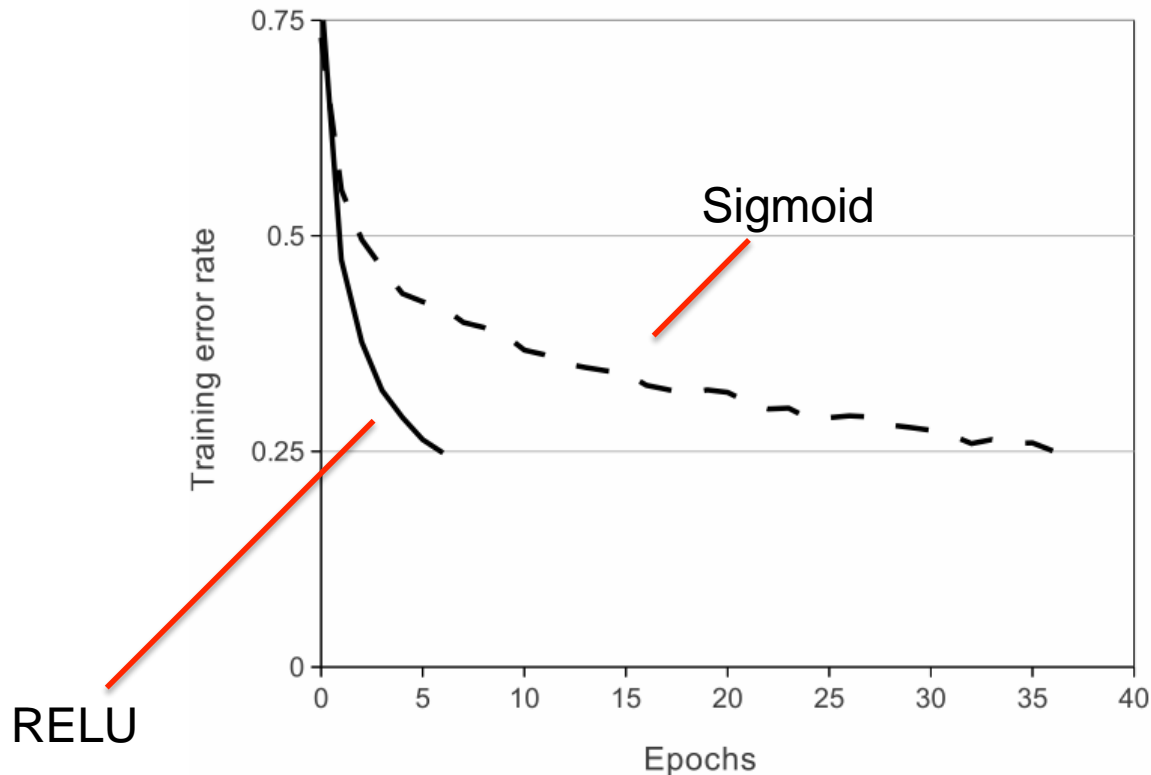


**Training is significantly faster when using the RELU function**

# Vanishing Gradients

## Learning Speed:

- The network that uses a RELU activation function learns significantly faster.



**Does this mean the problem of vanishing gradients is completely solved?**

# Vanishing Gradients

## Learning Speed:

- It turns out that even using RELU activation function doesn't completely eliminate the vanishing gradient problem.
- The key question is why the vanishing gradient problem still persists.
- To understand this we need to revisit the original back propagation algorithm.



# Backpropagation

1. Let  $\frac{\partial L}{\partial z_i^{(n)}} = \hat{y}_i - y_i$ , where  $n$  denotes the number of layers in the network.

2. For each **fully connected** layer  $l$ :

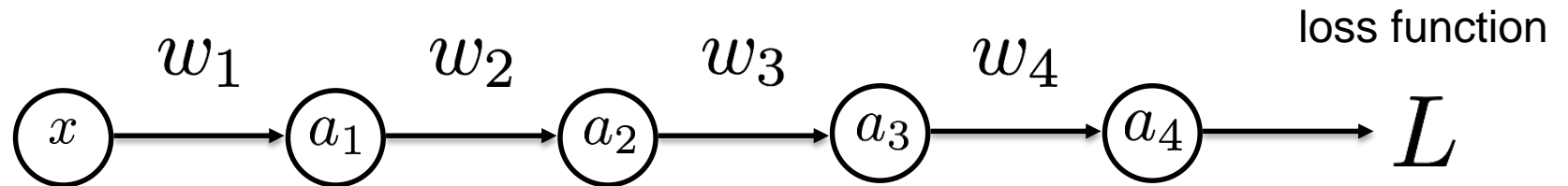
- For each node  $i$  in layer  $l$  set:

$$\frac{\partial L}{\partial z_i^{(l)}} = \left( \sum_{j=1}^{s^{l+1}} W_{ji}^{(l)} \frac{\partial L}{\partial z_j^{(l+1)}} \right) \frac{\partial f(z_i^{(l)})}{\partial z_i^{(l)}}$$

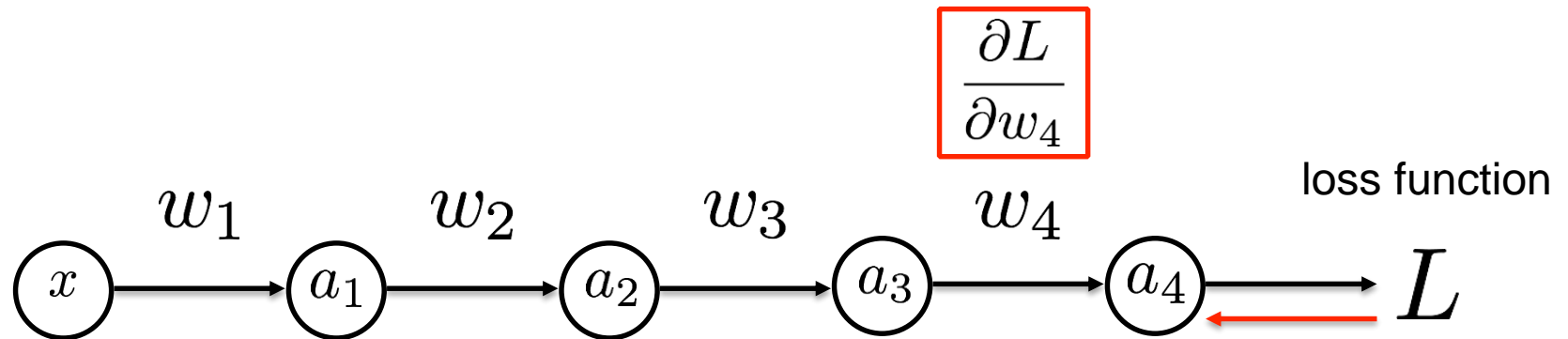
- Compute partial derivatives:  $\frac{\partial L}{\partial W_{ij}^{(l)}} = f'(z_j^{(l)}) \frac{\partial L}{\partial z_i^{(l+1)}}$

- Update the parameters:  $W_{ij}^{(l)} = W_{ij}^{(l)} - \alpha \frac{\partial L}{\partial W_{ij}^{(l)}}$

# Vanishing Gradients

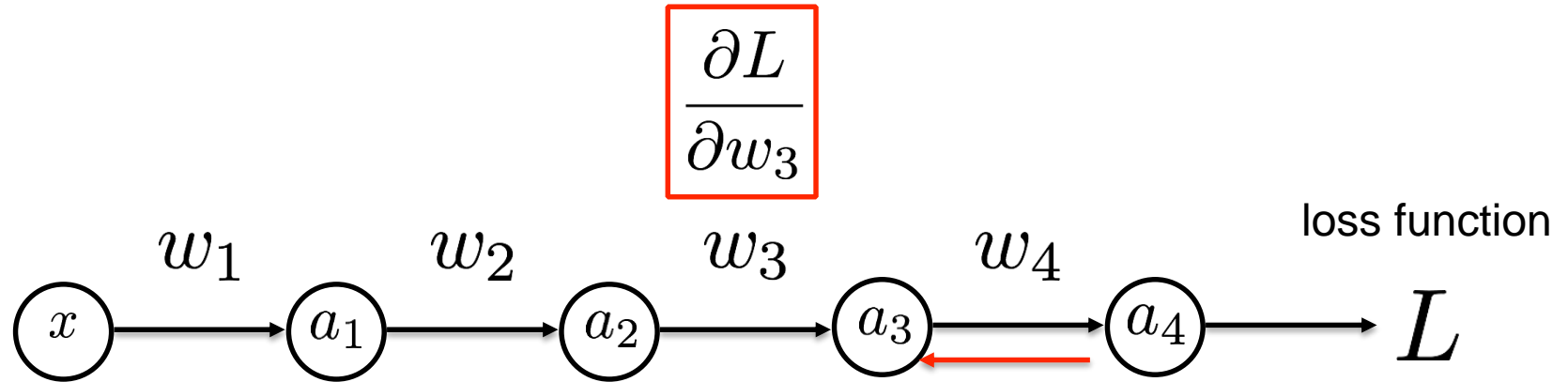


# Vanishing Gradients



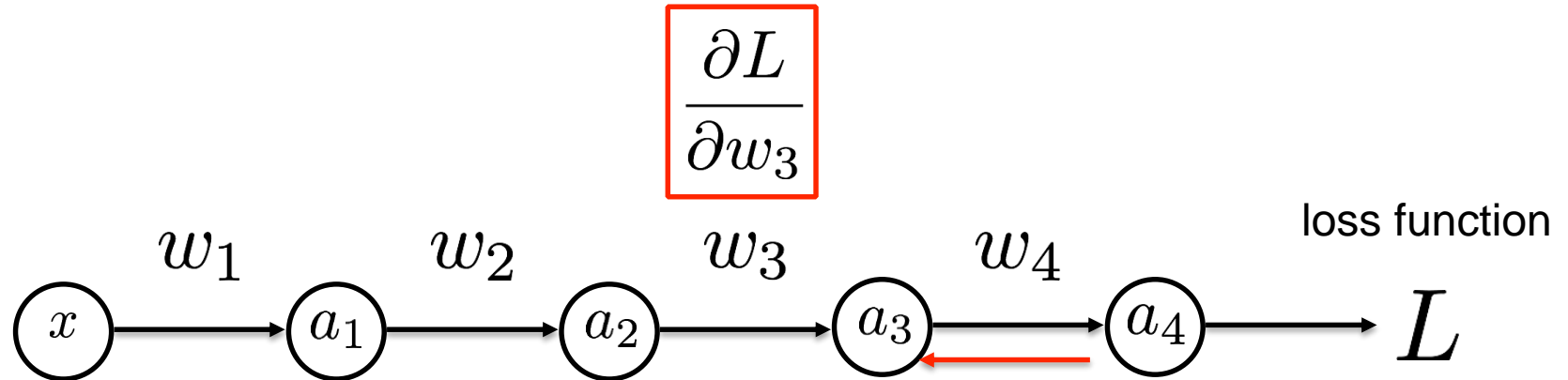
$$\frac{\partial L}{\partial w_4} = \frac{\partial L}{\partial a_4} \sigma'(z_4) a_3$$

# Vanishing Gradients



$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial a_3} \sigma'(z_3) a_2$$

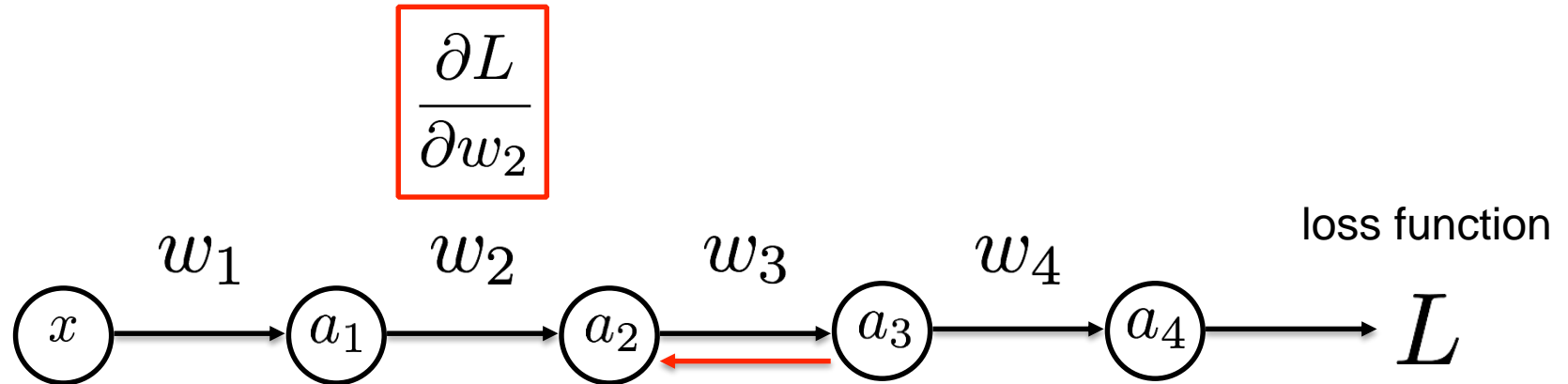
# Vanishing Gradients



$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial a_3} \sigma'(z_3) a_2$$

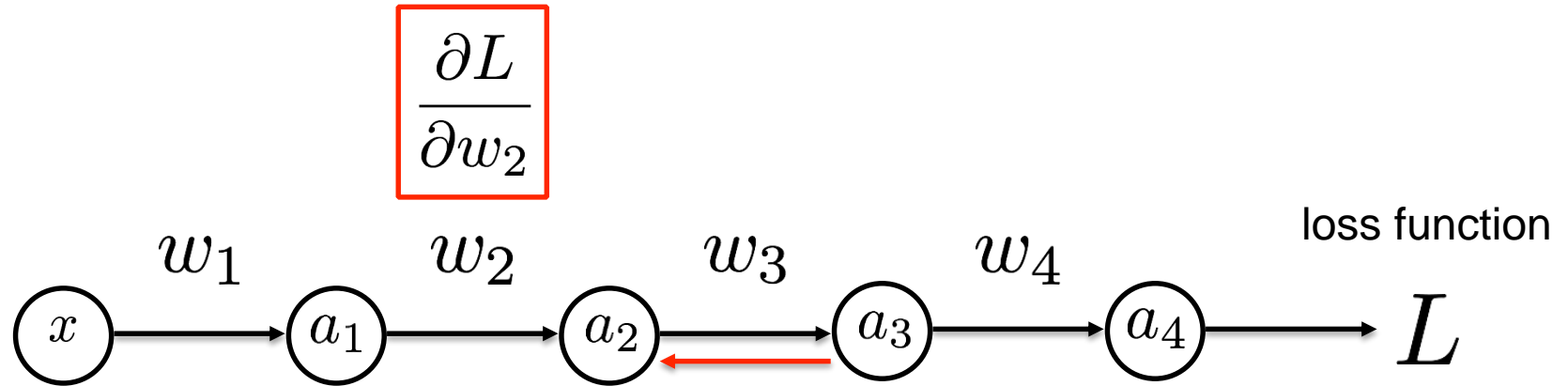
$$= \frac{\partial L}{\partial a_4} \sigma'(z_4) w_4 \sigma'(z_3) a_2$$

# Vanishing Gradients



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_2} \sigma'(z_2) a_1$$

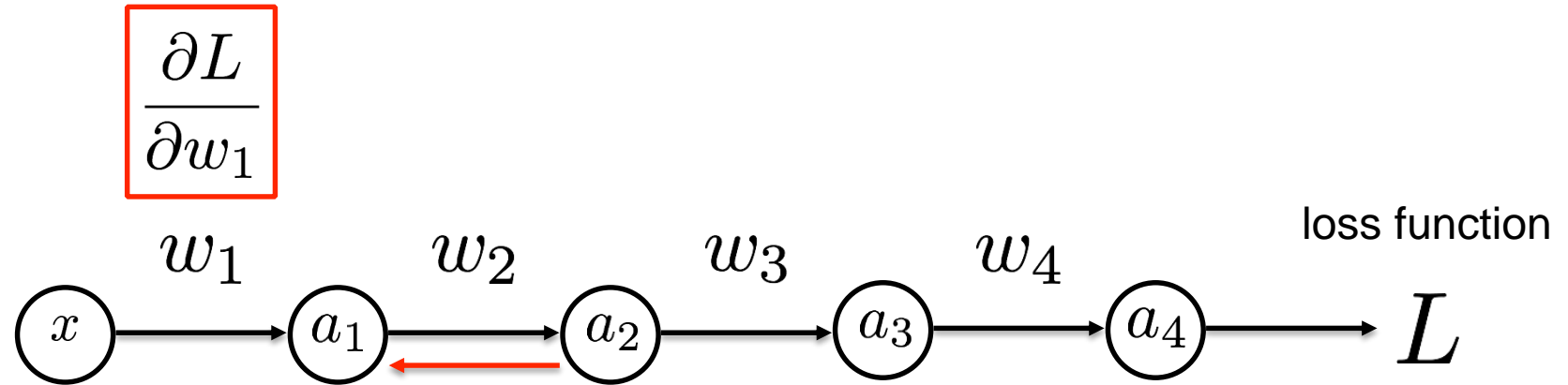
# Vanishing Gradients



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_2} \sigma'(z_2) a_1$$

$$= \frac{\partial L}{\partial a_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) a_1$$

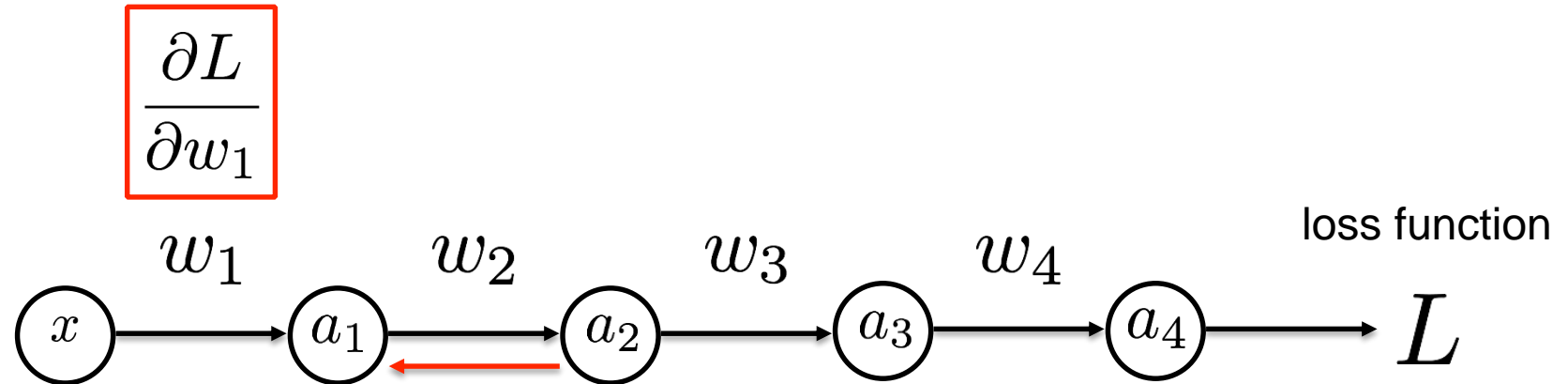
# Vanishing Gradients



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$



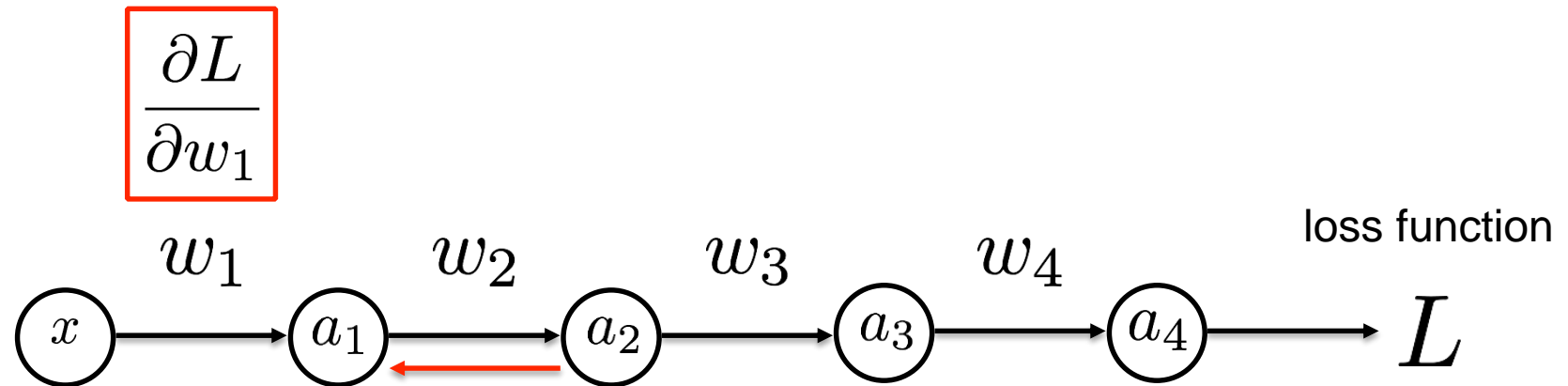
# Vanishing Gradients



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) x$$

# Vanishing Gradients

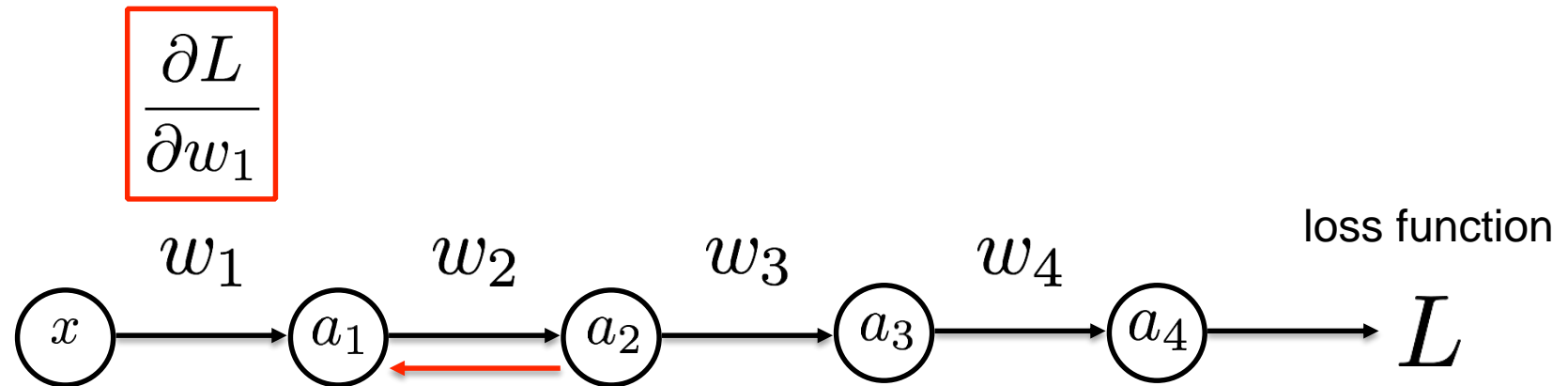


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \underbrace{\sigma'(z_4)}_{\sim |0.01|} w_4 \underbrace{\sigma'(z_3)}_{\sim |0.01|} w_3 \underbrace{\sigma'(z_2)}_{\sim |0.01|} w_2 \sigma'(z_1) x$$

- **Weights are typically initialized to small values (e.g. gaussian distribution with 0 mean and 0.01 std dev)**

# Vanishing Gradients

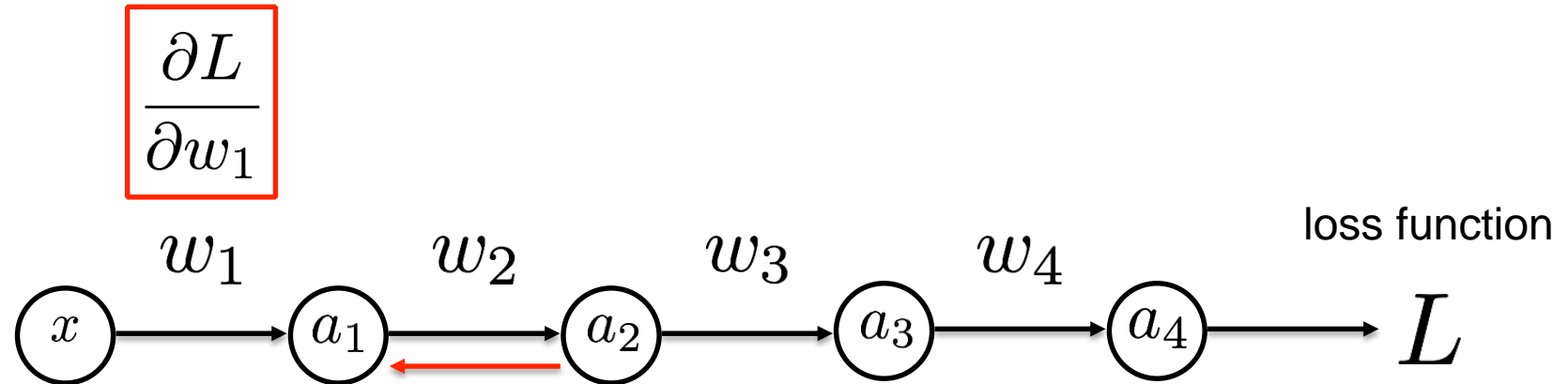


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \underbrace{\sigma'(z_4)}_{\sim |0.01|} w_4 \underbrace{\sigma'(z_3)}_{\sim |0.01|} w_3 \underbrace{\sigma'(z_2)}_{\sim |0.01|} w_2 \sigma'(z_1) x$$

- Exponential decrease in the gradient as we move towards the early hidden layers.

# Vanishing Gradients



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \underbrace{\sigma'(z_4)}_{\sim |0.01|} w_4 \underbrace{\sigma'(z_3)}_{\sim |0.01|} w_3 \underbrace{\sigma'(z_2)}_{\sim |0.01|} w_2 \sigma'(z_1) x$$

- As a result, the deepest hidden layers learn significantly faster, relative to the early layers that may not learn much at all.

# Deep Supervision

- We can reduce the vanishing gradient problem and make learning more effective via deep supervision.

# Deep Supervision

- We can reduce the vanishing gradient problem and make learning more effective via deep supervision.
- Deep supervision refers to a concept of adding learning objectives / loss functions to the intermediate hidden layers.

# Deep Supervision

- We can reduce the vanishing gradient problem and make learning more effective via deep supervision.
- Deep supervision refers to a concept of adding learning objectives / loss functions to the intermediate hidden layers.
- Backpropagation proceeds as usual, but now the gradients are propagated not from one but from multiple loss layers.

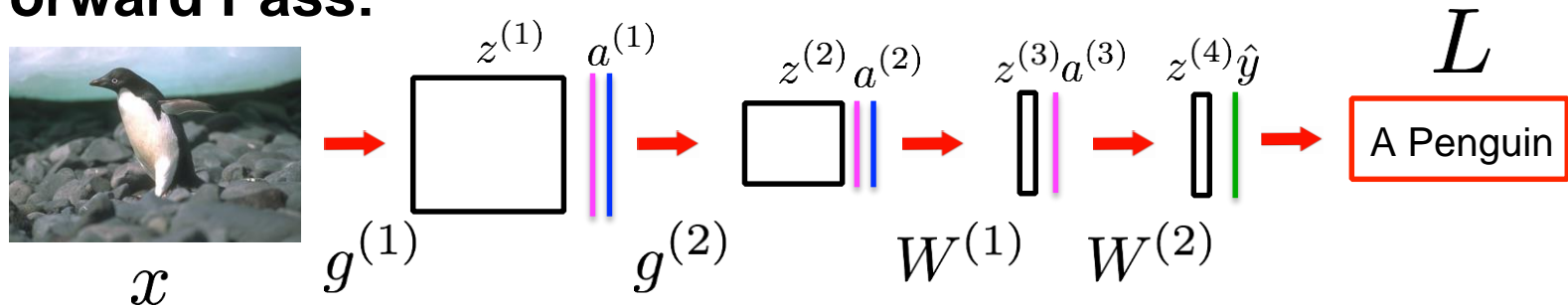
# Deep Supervision

## Standard CNN:

### Notation:

$\square$  - convolutional layer output     $\parallel$  - fully connected layer output  
 $|$  - max pooling layer     $|$  - sigmoid function  $f$      $|$  - softmax function

### Forward Pass:





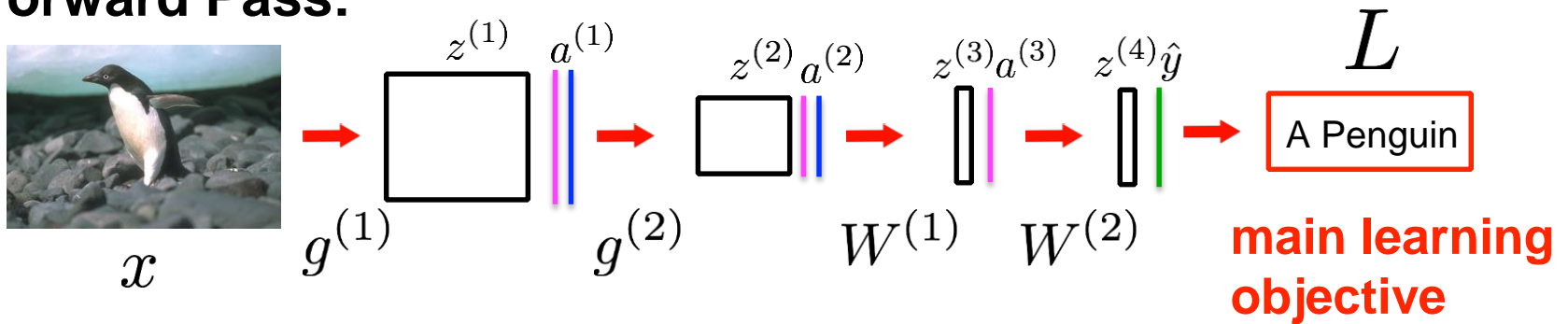
# Deep Supervision

## Standard CNN:

### Notation:

$\square$  - convolutional layer output     $\parallel$  - fully connected layer output  
 $|$  - max pooling layer     $|$  - sigmoid function  $f$      $|$  - softmax function

### Forward Pass:



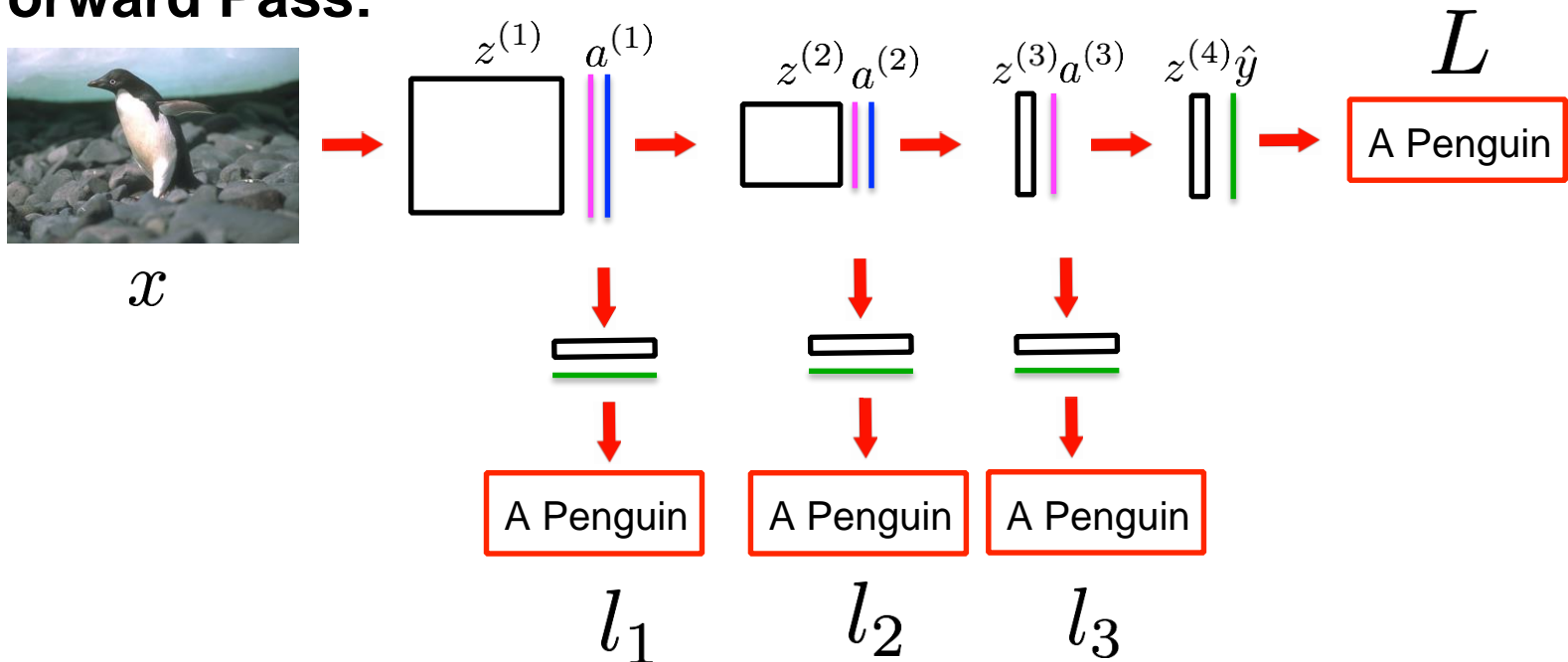
# Deep Supervision

## Deeply Supervised CNN:

### Notation:

$\square$  - convolutional layer output     $\parallel$  - fully connected layer output  
 $|$  - max pooling layer     $|$  - sigmoid function  $f$      $|$  - softmax function

### Forward Pass:



# Deep Supervision

## Deeply Supervised CNN:

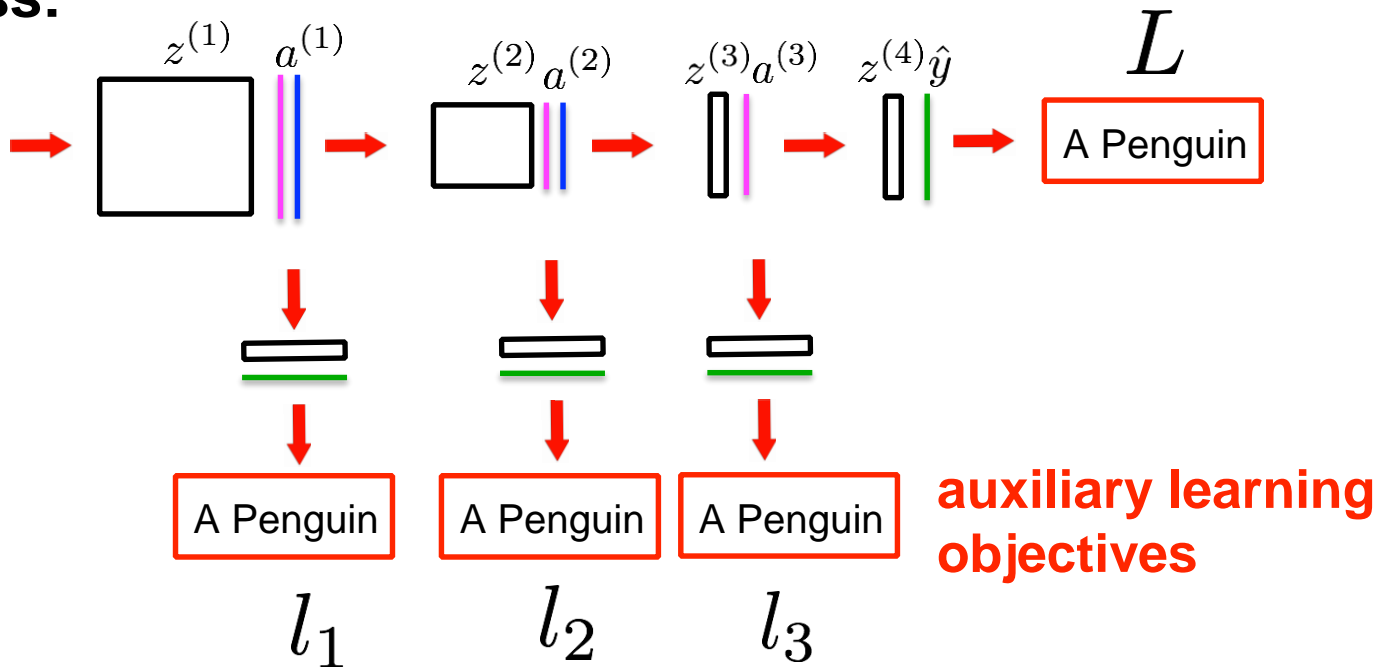
### Notation:

$\square$  - convolutional layer output     $\parallel$  - fully connected layer output  
 $|$  - max pooling layer     $|$  - sigmoid function  $f$      $|$  - softmax function

### Forward Pass:



$x$



# Deep Supervision

## Deeply Supervised CNN:

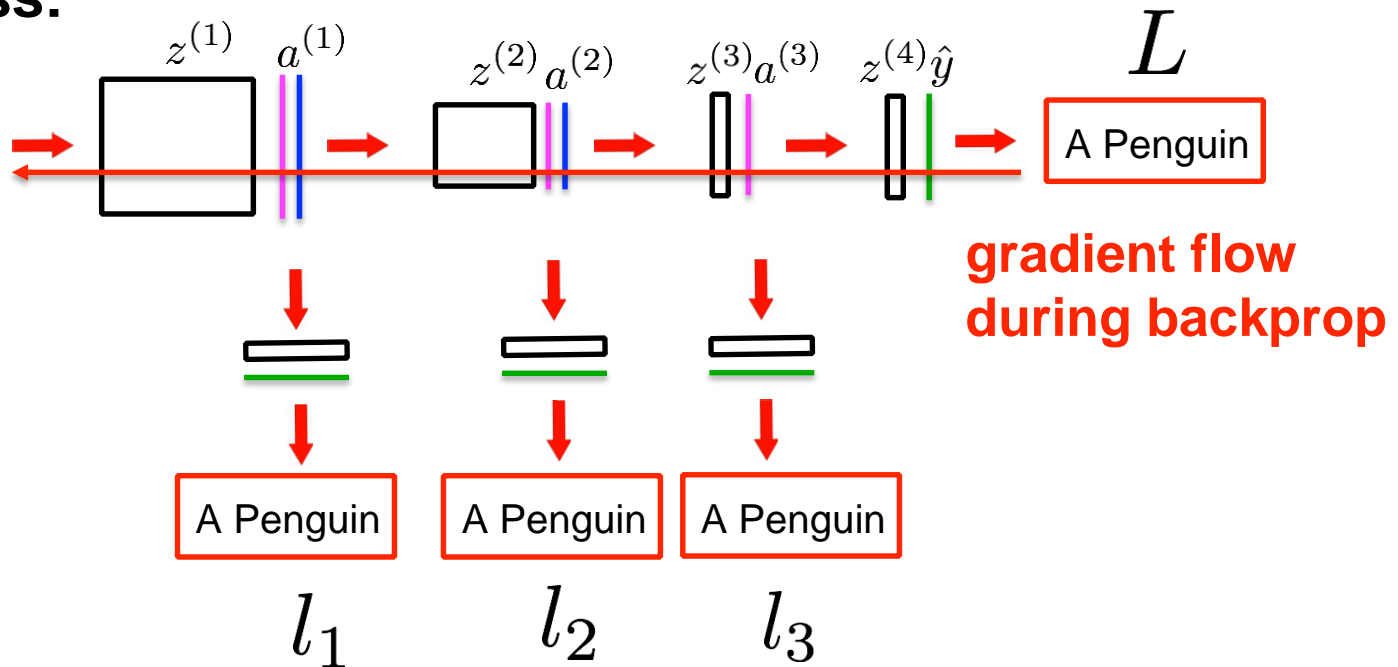
### Notation:

$\square$  - convolutional layer output     $\parallel$  - fully connected layer output  
 $|$  - max pooling layer     $|$  - sigmoid function  $f$      $|$  - softmax function

### Forward Pass:



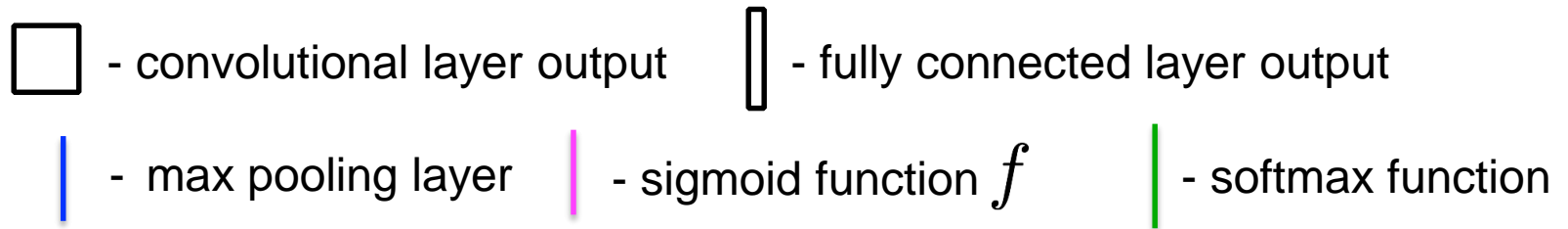
$x$



# Deep Supervision

## Deeply Supervised CNN:

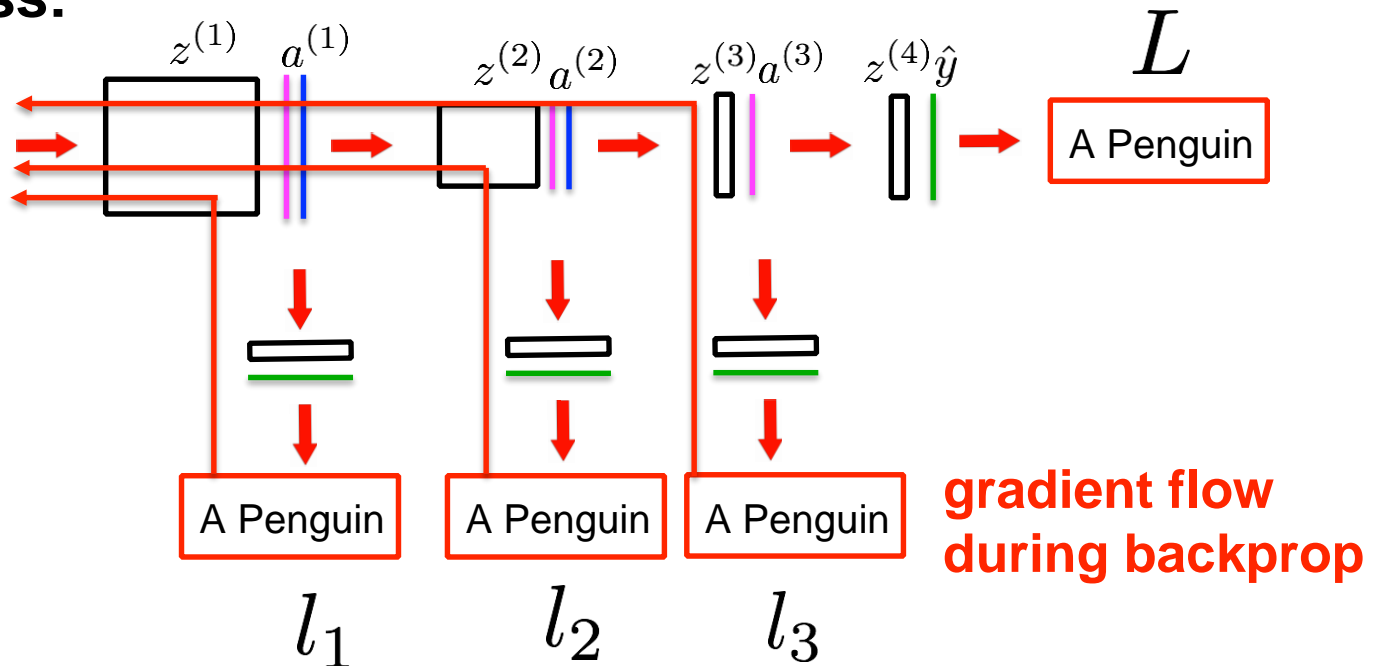
### Notation:



### Forward Pass:



$x$



# Deep Supervision

- Assume that we are given a **labeled** training dataset

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$$

- Typically we would employ the following loss function:

$$L = -(y \log f(x) + (1 - y) \log (1 - f(x)))$$

# Deep Supervision

- Assume that we are given a **labeled** training dataset

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$$

- Typically we would employ the following loss function:

$$L = -(y \log f(x) + (1 - y) \log (1 - f(x)))$$

- Under deeply supervised networks, we will use:

$$L_{new} = L + \sum_j \alpha_j l_j$$

# Deep Supervision

- Assume that we are given a **labeled** training dataset

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$$

- Typically we would employ the following loss function:

$$L = -(y \log f(x) + (1 - y) \log (1 - f(x)))$$

- Under deeply supervised networks, we will use:

$$L_{new} = L + \sum_j \alpha_j l_j$$

**auxiliary learning objectives**



# Deep Supervision

- Each auxiliary learning objective can be written as:

$$l_j = -(y \log f_j(x) + (1 - y) \log (1 - f_j(x)))$$

- $f_j$  refers to the network's output after a certain hidden layer  $j$ .

# Deep Supervision

- Each auxiliary learning objective can be written as:

$$l_j = -(y \log f_j(x) + (1 - y) \log (1 - f_j(x)))$$

- $f_j$  refers to the network's output after a certain hidden layer  $j$ .
- The differentiation can be done just as it's done with the main learning objective.

# Deep Supervision

- Each auxiliary learning objective can be written as:

$$l_j = -(y \log f_j(x) + (1 - y) \log (1 - f_j(x)))$$

- $f_j$  refers to the network's output after a certain hidden layer  $j$ .
- The differentiation can be done just as it's done with the main learning objective.
- The gradients from different learning objectives are summed in the hidden layers during the back propagation.

# Deep Supervision

- Each auxiliary learning objective can be written as:

$$l_j = -(y \log f_j(x) + (1 - y) \log (1 - f_j(x)))$$

- $f_j$  refers to the network's output after a certain hidden layer  $j$ .
- The differentiation can be done just as it's done with the main learning objective
- The gradients from different learning objectives are summed in the hidden layers during the back propagation.
  - **Helps to learn more discriminative features**
  - **Alleviates the vanishing gradient problem**

# Deep Supervision

- Standard CNN

```
1 name: "LeNet"
2 layers {
3   name: "mnist"
4   type: DATA
5   top: "data"
6   top: "label"
7   data_param {
8     source: "mnist-train-leveldb"
9     scale: 0.00390625
10    batch_size: 64
11  }
12 }
13 layers {
14   name: "conv1"
15   type: CONVOLUTION
16   bottom: "data"
17   top: "conv1"
18   blobs_lr: 1
19   blobs_lr: 2
20   convolution_param {
21     num_output: 20
22     kernel_size: 5
23     stride: 1
24     weight_filler {
25       type: "xavier"
26     }
27     bias_filler {
28       type: "constant"
29     }
30  }
31 }
32 layers {
33   name: "pool1"
34   type: POOLING
35   bottom: "conv1"
36   top: "pool1"
37   pooling_param {
38     pool: MAX
39     kernel_size: 2
40     stride: 2
41  }
42 }
```

```
43 layers {
44   name: "conv2"
45   type: CONVOLUTION
46   bottom: "pool1"
47   top: "conv2"
48   blobs_lr: 1
49   blobs_lr: 2
50   convolution_param {
51     num_output: 50
52     kernel_size: 5
53     stride: 1
54     weight_filler {
55       type: "xavier"
56     }
57     bias_filler {
58       type: "constant"
59     }
60  }
61 }
62 layers {
63   name: "pool2"
64   type: POOLING
65   bottom: "conv2"
66   top: "pool2"
67   pooling_param {
68     pool: MAX
69     kernel_size: 2
70     stride: 2
71  }
72 }
73 layers {
74   name: "ip1"
75   type: INNER_PRODUCT
76   bottom: "pool2"
77   top: "ip1"
78   blobs_lr: 1
79   blobs_lr: 2
80   inner_product_param {
81     num_output: 500
82     weight_filler {
83       type: "xavier"
84     }
85     bias_filler {
86       type: "constant"
87     }
88  }
89 }
```

# Deep Supervision

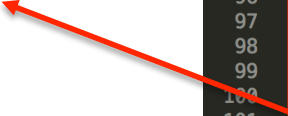
- Standard CNN

```
73 layers {
74   name: "ip1"
75   type: INNER_PRODUCT
76   bottom: "pool2"
77   top: "ip1"
78   blobs_lr: 1
79   blobs_lr: 2
80   inner_product_param {
81     num_output: 500
82     weight_filler {
83       type: "xavier"
84     }
85     bias_filler {
86       type: "constant"
87     }
88   }
89 }
90 layers {
91   name: "relu1"
92   type: RELU
93   bottom: "ip1"
94   top: "ip1"
95 }
96 layers {
97   name: "ip2"
98   type: INNER_PRODUCT
99   bottom: "ip1"
100  top: "ip2"
101  blobs_lr: 1
102  blobs_lr: 2
103  inner_product_param {
104    num_output: 10
105    weight_filler {
106      type: "xavier"
107    }
108    bias_filler {
109      type: "constant"
110    }
111  }
112 }
113 layers {
114   name: "loss"
115   type: SOFTMAX_LOSS
116   bottom: "ip2"
117   bottom: "label"
118 }
119 }
```

# Deep Supervision

- Standard CNN

Layer that  
produces  
predictions



```
73 layers {
74   name: "ip1"
75   type: INNER_PRODUCT
76   bottom: "pool2"
77   top: "ip1"
78   blobs_lr: 1
79   blobs_lr: 2
80   inner_product_param {
81     num_output: 500
82     weight_filler {
83       type: "xavier"
84     }
85     bias_filler {
86       type: "constant"
87     }
88   }
89 }
90 layers {
91   name: "relu1"
92   type: RELU
93   bottom: "ip1"
94   top: "ip1"
95 }
96 layers {
97   name: "ip2"
98   type: INNER_PRODUCT
99   bottom: "ip1"
100  top: "ip2"
101  blobs_lr: 1
102  blobs_lr: 2
103  inner_product_param {
104    num_output: 10
105    weight_filler {
106      type: "xavier"
107    }
108    bias_filler {
109      type: "constant"
110    }
111  }
112 }
113 layers {
114   name: "loss"
115   type: SOFTMAX_LOSS
116   bottom: "ip2"
117   bottom: "label"
118 }
119 }
```

# Deep Supervision

- Standard CNN

```
73 layers {
74   name: "ip1"
75   type: INNER_PRODUCT
76   bottom: "pool2"
77   top: "ip1"
78   blobs_lr: 1
79   blobs_lr: 2
80   inner_product_param {
81     num_output: 500
82     weight_filler {
83       type: "xavier"
84     }
85     bias_filler {
86       type: "constant"
87     }
88   }
89 }
90 layers {
91   name: "relu1"
92   type: RELU
93   bottom: "ip1"
94   top: "ip1"
95 }
96 layers {
97   name: "ip2"
98   type: INNER_PRODUCT
99   bottom: "ip1"
100  top: "ip2"
101  blobs_lr: 1
102  blobs_lr: 2
103  inner_product_param {
104    num_output: 10
105    weight_filler {
106      type: "xavier"
107    }
108    bias_filler {
109      type: "constant"
110    }
111  }
112 }
113 layers {
114   name: "loss"
115   type: SOFTMAX_LOSS
116   bottom: "ip2"
117   bottom: "label"
118 }
119 }
```

A softmax loss  
layer





# Deep Supervision

- Deeply Supervised CNN

```
1 name: "LeNet"
2 layers {
3   name: "mnist"
4   type: DATA
5   top: "data"
6   top: "label"
7   data_param {
8     source: "mnist-train-leveldb"
9     scale: 0.00390625
10    batch_size: 64
11  }
12 }
13 layers {
14   name: "conv1"
15   type: CONVOLUTION
16   bottom: "data"
17   top: "conv1"
18   blobs_lr: 1
19   blobs_lr: 2
20   convolution_param {
21     num_output: 20
22     kernel_size: 5
23     stride: 1
24     weight_filler {
25       type: "xavier"
26     }
27     bias_filler {
28       type: "constant"
29     }
30   }
31 }
32 layers {
33   name: "pool1"
34   type: POOLING
35   bottom: "conv1"
36   top: "pool1"
37   pooling_param {
38     pool: MAX
39     kernel_size: 2
40     stride: 2
41   }
42 }
```

```
44 layers {
45   name: "ip1_side"
46   type: INNER_PRODUCT
47   bottom: "pool1"
48   top: "ip1_side"
49   blobs_lr: 1
50   blobs_lr: 2
51   inner_product_param {
52     num_output: 10
53     weight_filler {
54       type: "xavier"
55     }
56     bias_filler {
57       type: "constant"
58     }
59   }
60 }
61 layers {
62   name: "l1"
63   type: SOFTMAX_LOSS
64   bottom: "ip1_side"
65   bottom: "label"
66 }
67 }
```

# Deep Supervision

- Deeply Supervised CNN

A side layer that produces predictions

```
1 name: "LeNet"
2 layers {
3   name: "mnist"
4   type: DATA
5   top: "data"
6   top: "label"
7   data_param {
8     source: "mnist-train-leveldb"
9     scale: 0.00390625
10    batch_size: 64
11  }
12 }
13 layers {
14   name: "conv1"
15   type: CONVOLUTION
16   bottom: "data"
17   top: "conv1"
18   blobs_lr: 1
19   blobs_lr: 2
20   convolution_param {
21     num_output: 20
22     kernel_size: 5
23     stride: 1
24     weight_filler {
25       type: "xavier"
26     }
27     bias_filler {
28       type: "constant"
29     }
30   }
31 }
32 layers {
33   name: "pool1"
34   type: POOLING
35   bottom: "conv1"
36   top: "pool1"
37   pooling_param {
38     pool: MAX
39     kernel_size: 2
40     stride: 2
41   }
42 }
```

```
44 layers {
45   name: "ip1_side"
46   type: INNER_PRODUCT
47   bottom: "pool1"
48   top: "ip1_side"
49   blobs_lr: 1
50   blobs_lr: 2
51   inner_product_param {
52     num_output: 10
53     weight_filler {
54       type: "xavier"
55     }
56     bias_filler {
57       type: "constant"
58     }
59   }
60 }
61 layers {
62   name: "l1"
63   type: SOFTMAX_LOSS
64   bottom: "ip1_side"
65   bottom: "label"
66 }
67 }
```

# Deep Supervision

- Deeply Supervised CNN

```
1 name: "LeNet"
2 layers {
3   name: "mnist"
4   type: DATA
5   top: "data"
6   top: "label"
7   data_param {
8     source: "mnist-train-leveldb"
9     scale: 0.00390625
10    batch_size: 64
11  }
12 }
13 layers {
14   name: "conv1"
15   type: CONVOLUTION
16   bottom: "data"
17   top: "conv1"
18   blobs_lr: 1
19   blobs_lr: 2
20   convolution_param {
21     num_output: 20
22     kernel_size: 5
23     stride: 1
24     weight_filler {
25       type: "xavier"
26     }
27     bias_filler {
28       type: "constant"
29     }
30   }
31 }
32 layers {
33   name: "pool1"
34   type: POOLING
35   bottom: "conv1"
36   top: "pool1"
37   pooling_param {
38     pool: MAX
39     kernel_size: 2
40     stride: 2
41   }
42 }
```

```
44 layers {
45   name: "ip1_side"
46   type: INNER_PRODUCT
47   bottom: "pool1"
48   top: "ip1_side"
49   blobs_lr: 1
50   blobs_lr: 2
51   inner_product_param {
52     num_output: 10
53     weight_filler {
54       type: "xavier"
55     }
56     bias_filler {
57       type: "constant"
58     }
59   }
60 }
61 layers {
62   name: "l1"
63   type: SOFTMAX_LOSS
64   bottom: "ip1_side"
65   bottom: "label"
66 }
67 }
```



**An auxiliary  
loss function**

# Deep Supervision

## Some interesting results:

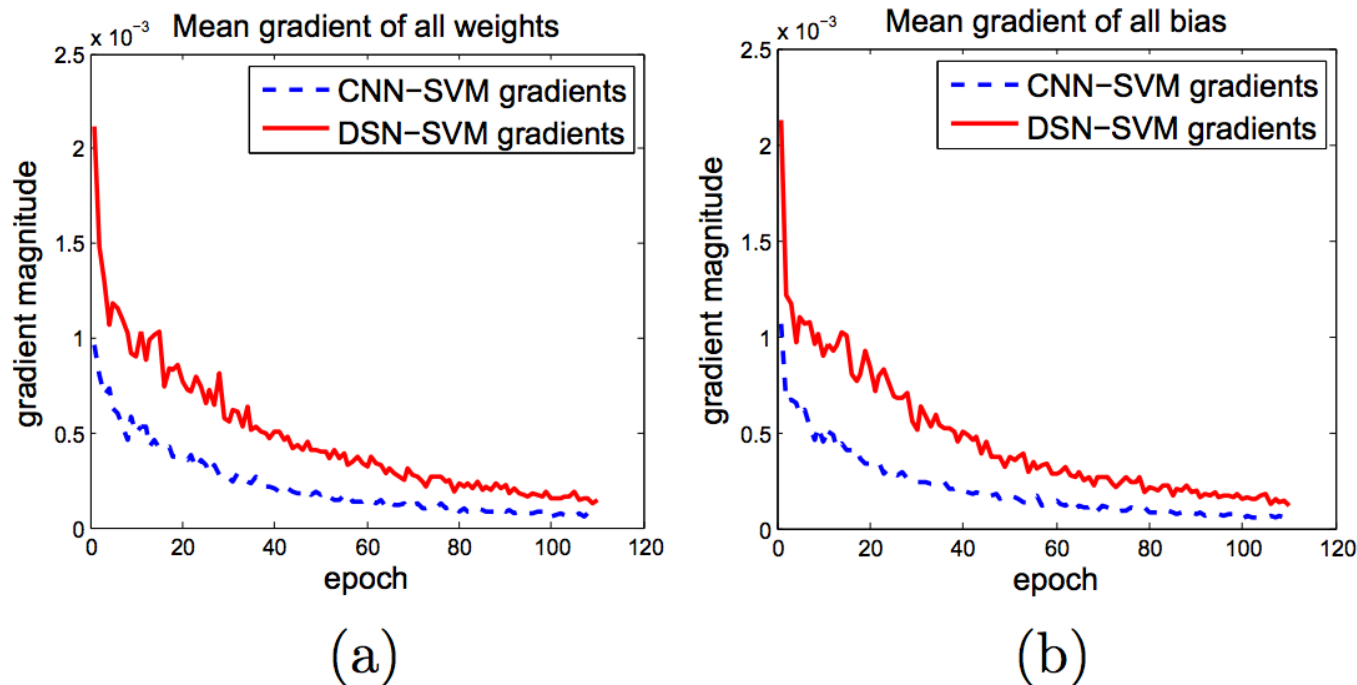


Figure 3: Average absolute value of gradient matrices during the training on MNIST for (a) weights; (b) biases.

# Deep Supervision

## Some interesting results:

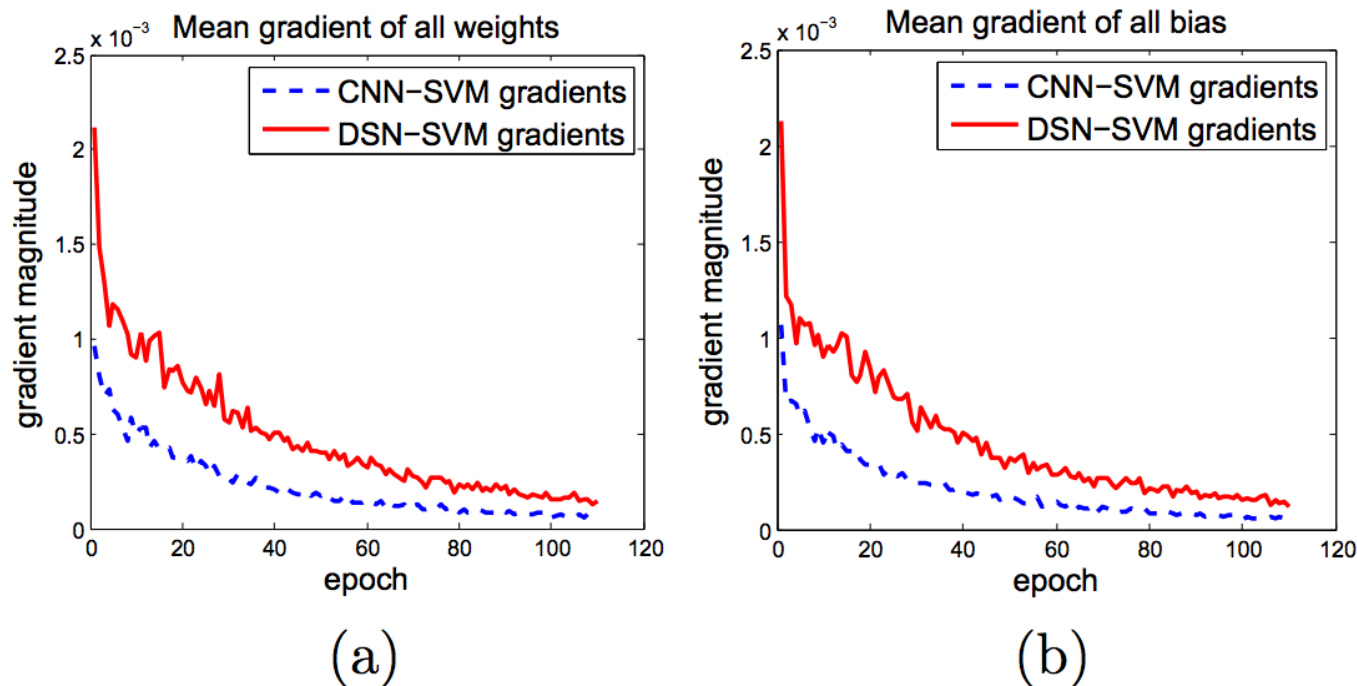
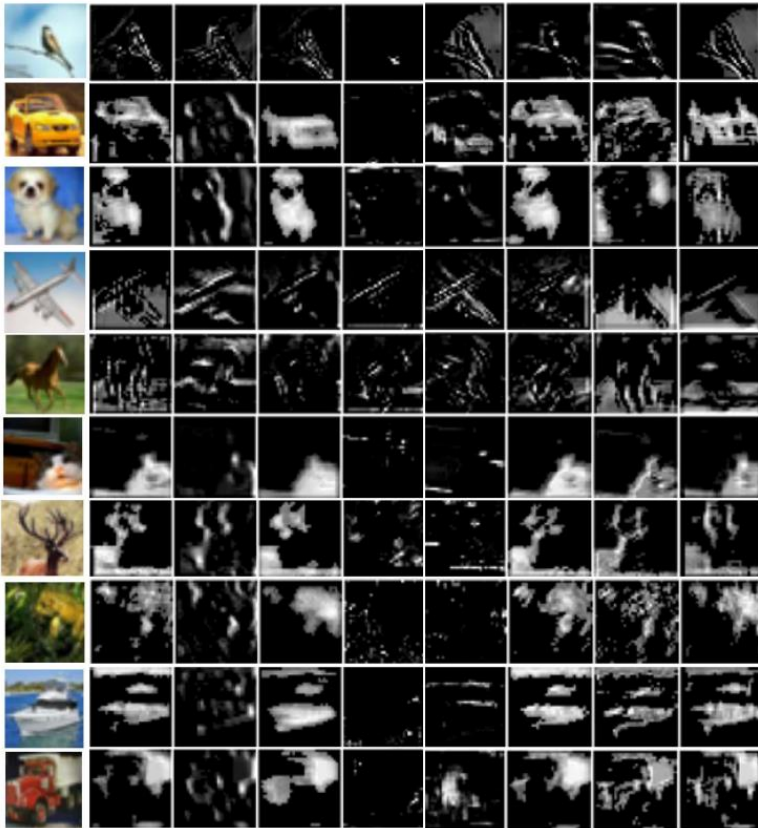


Figure 3: Average absolute value of gradient matrices during the training on MNIST for (a) weights; (b) biases.

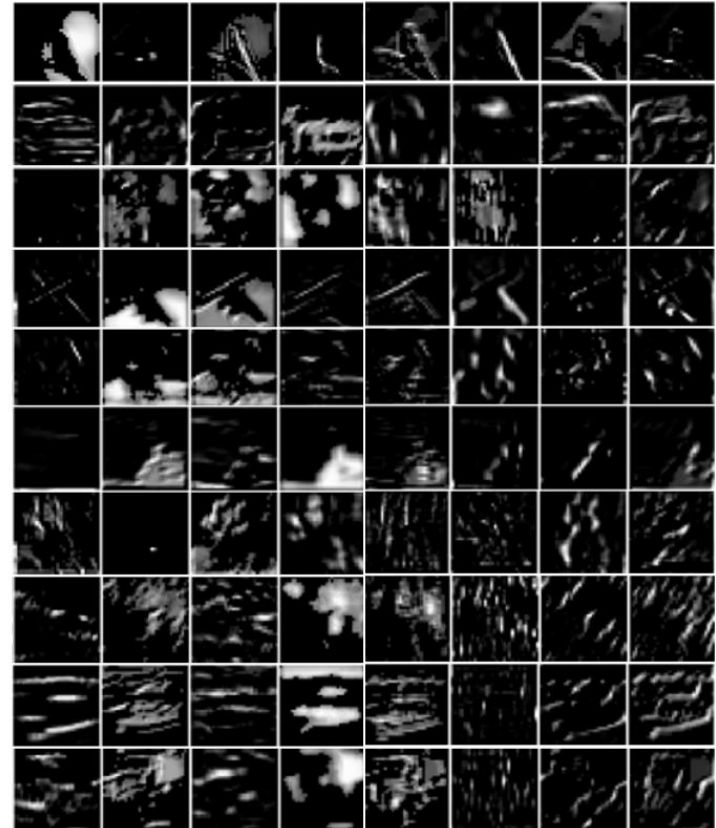
- **Deep supervision helps to reduce the vanishing gradient problem!**

# Deep Supervision

## Some interesting results:



(a) by DSN

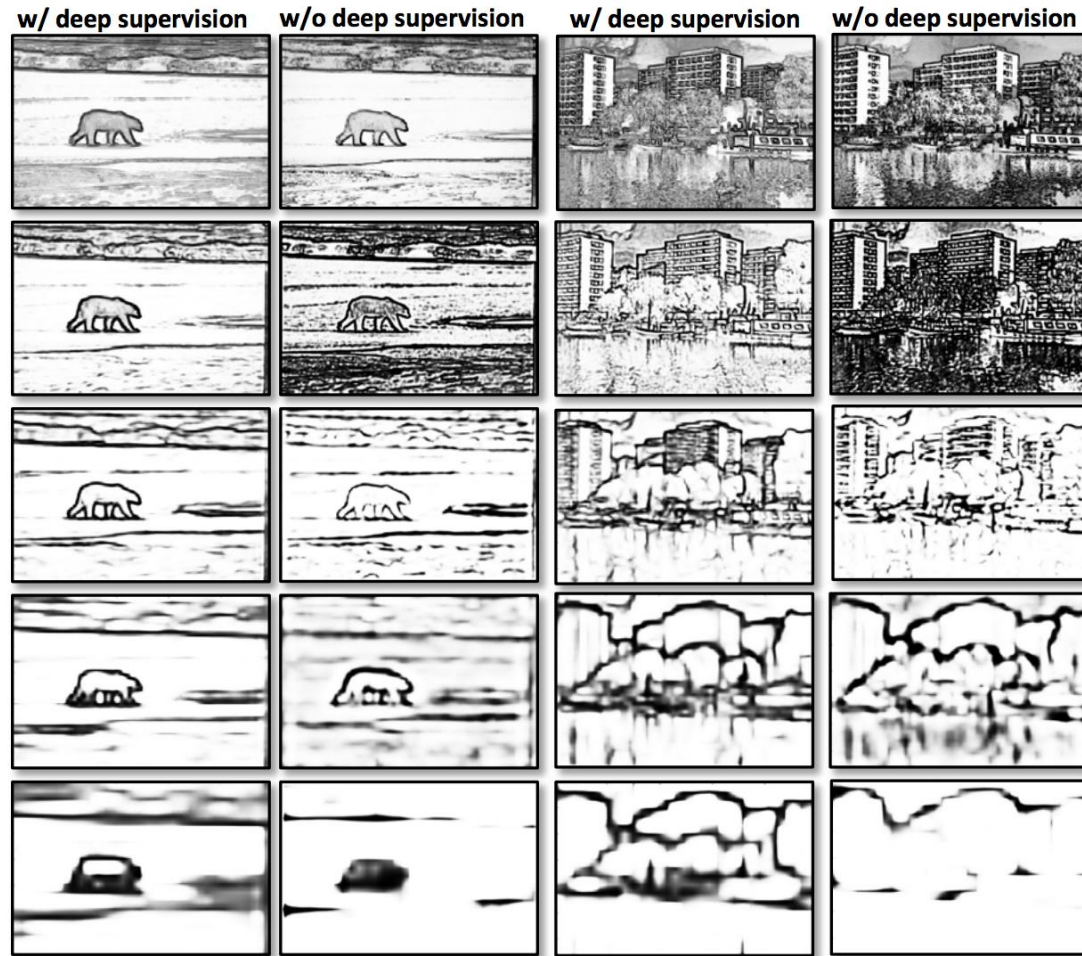


(b) by CNN

- **Learned features are more discriminative**

# Deep Supervision

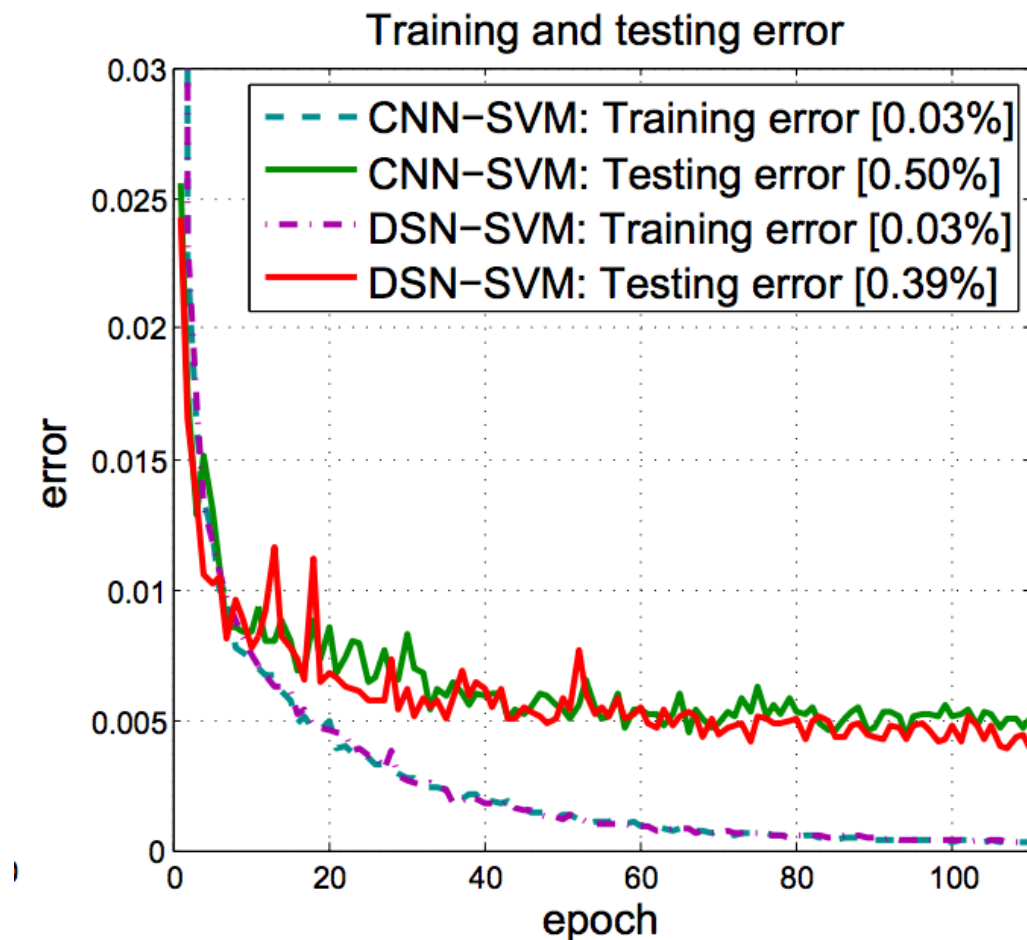
## Some interesting results:



- Learned features are more discriminative

# Deep Supervision

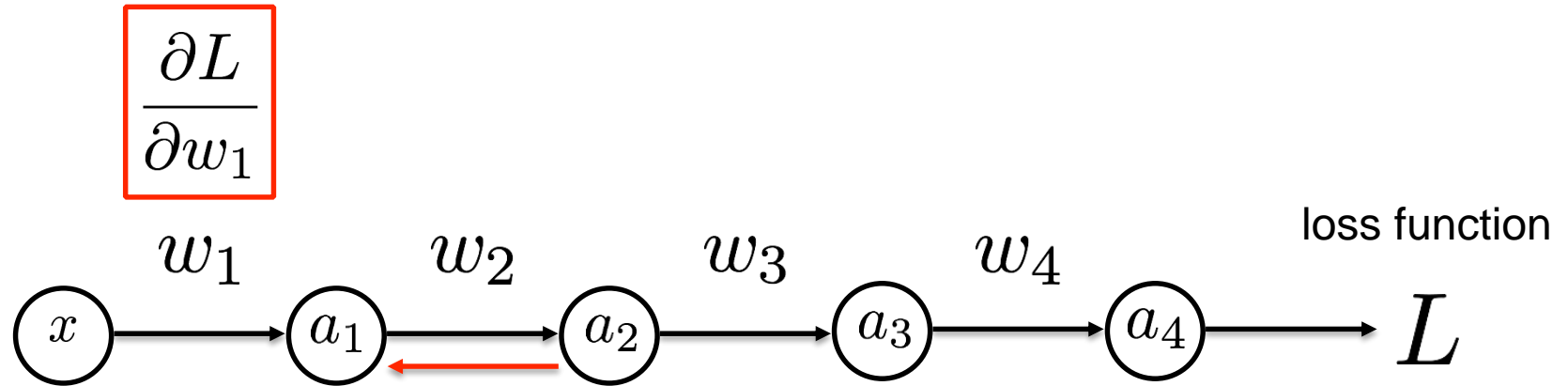
## Some interesting results:



- **Deep supervision reduces testing error without overfitting the training data.**



# Exploding Gradients

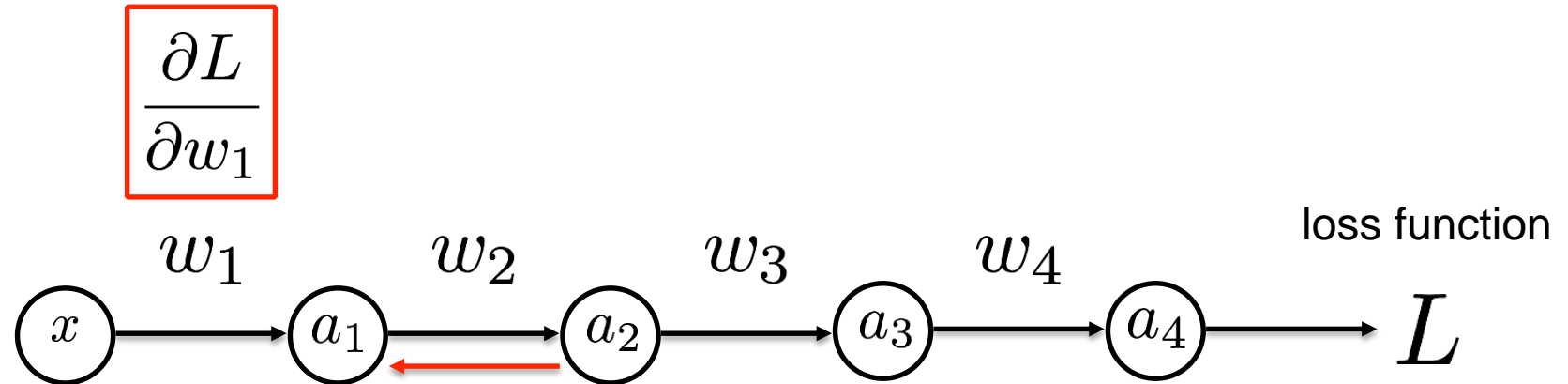


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) x$$

$>|1|$                        $>|1|$                        $>|1|$

# Exploding Gradients



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \sigma'(z_1) x$$

$$= \frac{\partial L}{\partial a_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) x$$

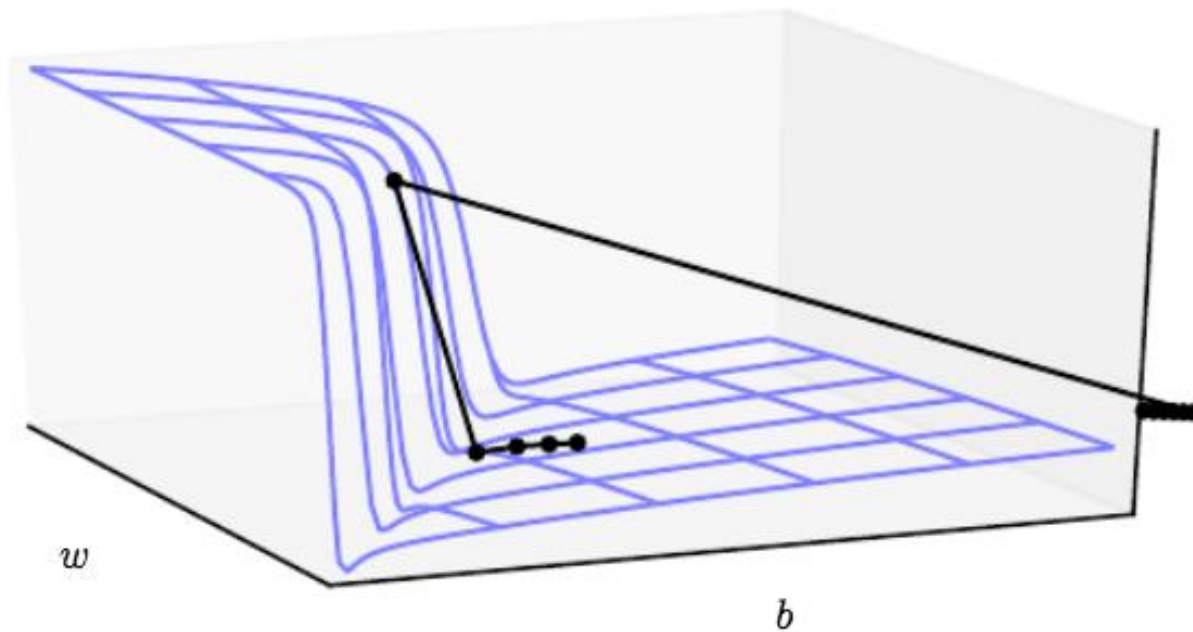
$>|1|$                        $>|1|$                        $>|1|$

- An instance of exploding gradients

# Exploding Gradients

## Exploding gradients:

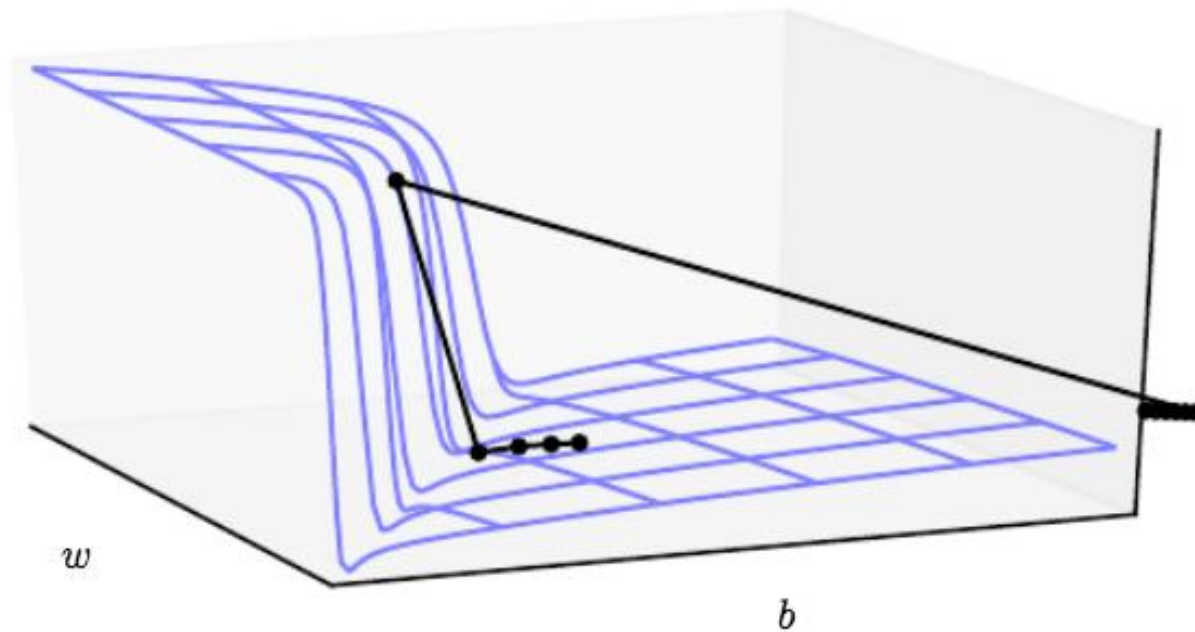
- One of the most often occurring learning problems.
- Due to large jumps parameter update becomes extremely unstable.



# Exploding Gradients

## Exploding gradients:

- One of the most often occurring learning problems.
- Due to large jumps parameter update becomes extremely unstable.

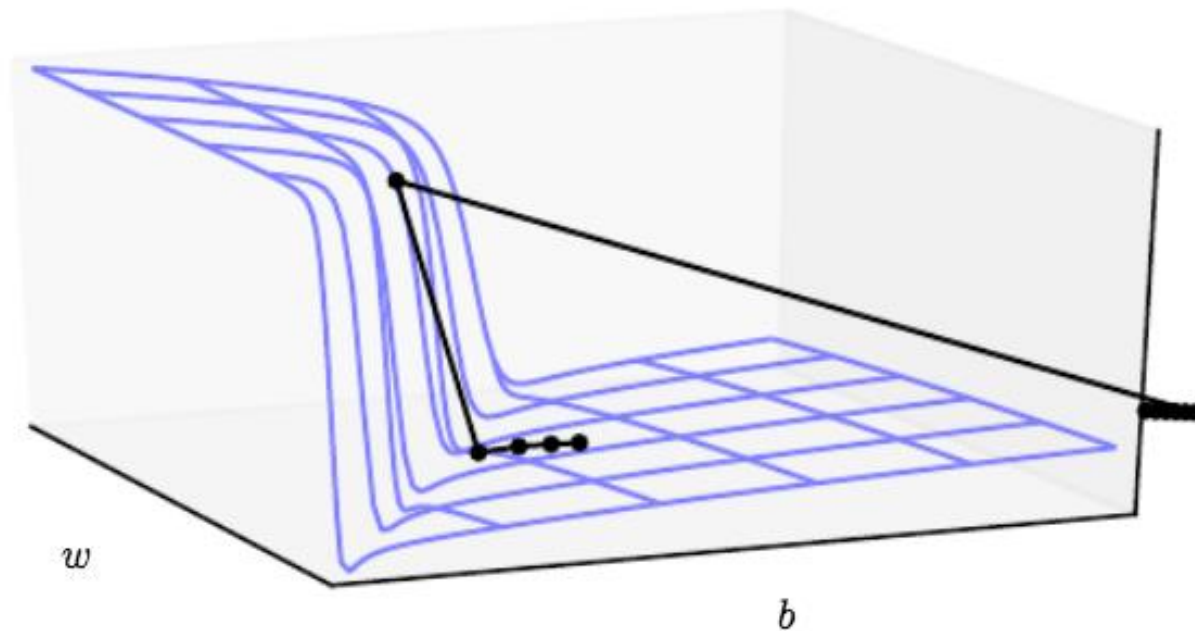


- **Solution #1: reduce the learning rate**

# Exploding Gradients

## Exploding gradients:

- One of the most often occurring learning problems.
- Due to large jumps parameter update becomes extremely unstable.



- **Solution #1: reduce the learning rate**
- **Solution #2: clip the gradients**

# Momentum

## Other ways to speed-up the training:

- Even if we address the vanishing gradient problem, the stochastic gradient descent (SGD) optimization is still quite slow.

# Momentum

## Other ways to speed-up the training:

- Even if we address the vanishing gradient problem, the stochastic gradient descent (SGD) optimization is still quite slow.
- We can accelerate the learning using the momentum method.

# Momentum

## Other ways to speed-up the training:

- Even if we address the vanishing gradient problem, the stochastic gradient descent (SGD) optimization is still quite slow.
- We can accelerate the learning using the momentum method.
- The momentum method introduces a speed variable, that keeps track of the direction and speed at which the parameters move through the parameter space.



# Momentum

## Standard gradient descent:

- Learning rule:

$$\theta = \theta - \epsilon \frac{\partial L}{\partial w}$$

# Momentum

## Standard gradient descent:

- Learning rule:

$$\theta = \theta - \epsilon \frac{\partial L}{\partial w}$$

## Gradient descent with momentum:

- Learning rule:

$$v = \alpha v - \epsilon \frac{\partial L}{\partial w}$$

# Momentum

## Standard gradient descent:

- Learning rule:

$$\theta = \theta - \epsilon \frac{\partial L}{\partial w}$$

## Gradient descent with momentum:

- Learning rule:

$$v = \alpha v - \epsilon \frac{\partial L}{\partial w}$$

$$\theta = \theta + v$$

# Momentum

## Standard gradient descent:

- Learning rule:

$$\theta = \theta - \epsilon \frac{\partial L}{\partial w}$$

## Gradient descent with momentum:

- Learning rule:

$$v = \alpha v - \epsilon \frac{\partial L}{\partial w}$$

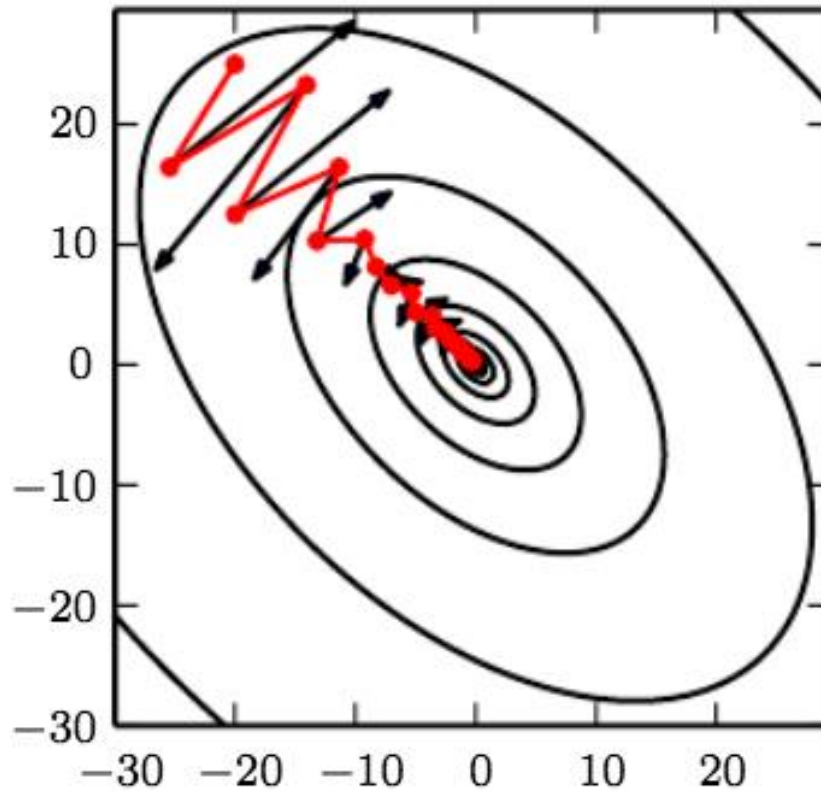
$$\theta = \theta + v$$

- **Makes it more difficult for the parameters to fluctuate a lot, which makes learning more stable**

# Momentum

→ SGD

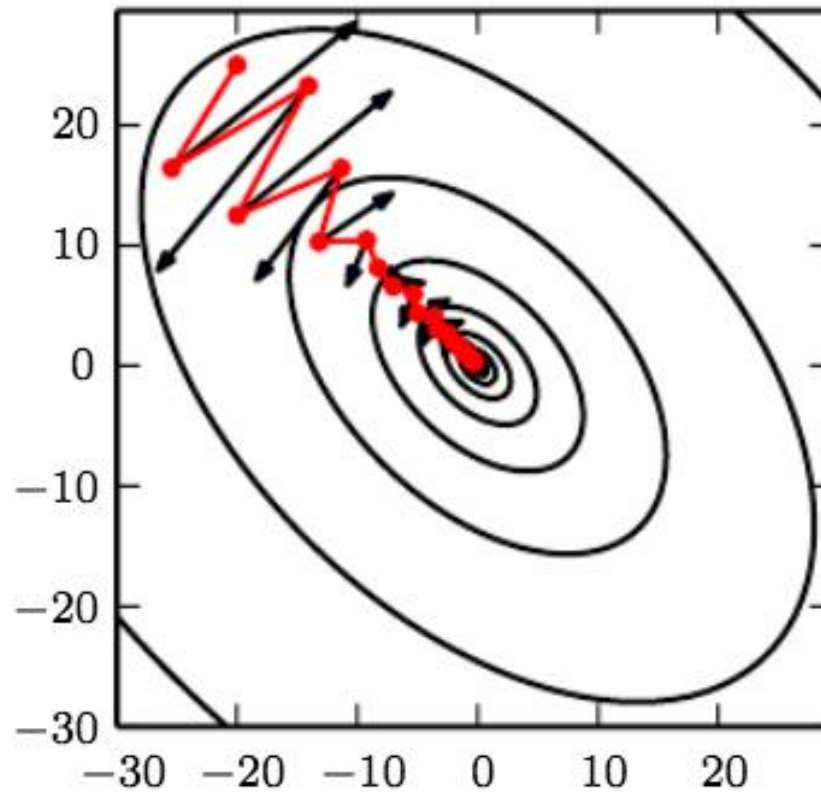
●—● SGD with Momentum



# Momentum

→ SGD

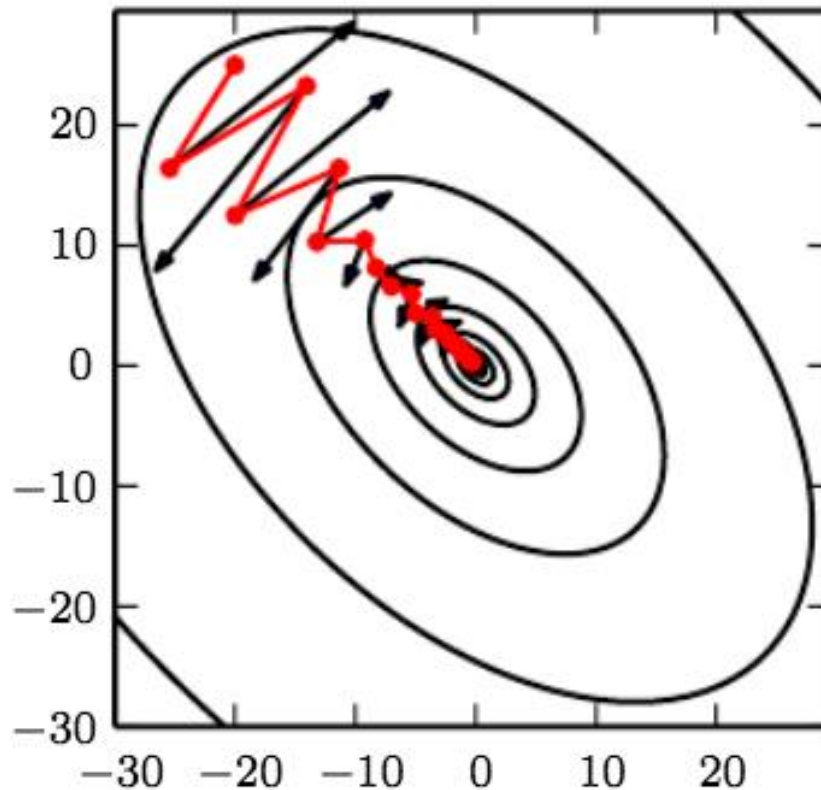
●—● SGD with Momentum



- Momentum helps to achieve more direct path towards local minimum

# Momentum

→ SGD  
●—● SGD with Momentum



- Momentum helps to achieve more direct path towards local minimum
- Therefore, learning becomes faster.

# Batch Normalization

## Batch Training Mode:

- SGD training is typically done in batch mode (e.g. by randomly selecting  $N$  samples from the training dataset, and averaging the gradient across them during backprop).



# Batch Normalization

## Batch Training Mode:

- SGD training is typically done in batch mode (e.g. by randomly selecting  $N$  samples from the training dataset, and averaging the gradient across them during backprop).

## Issues:

- Samples in different batches can be very different.
- Small changes to the network parameters amplify as the network becomes deeper
- This changes the internal node distribution in many layers.
- The layers need to continuously adapt to this new internal node distribution, which slows down the training.

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

**What's wrong with this approach?**

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

<b>Input:</b> Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$ ; Parameters to be learned: $\gamma, \beta$
<b>Output:</b> $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$
$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$

**What happens if we normalize the inputs to the sigmoid function?**

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

<b>Input:</b> Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$ ; Parameters to be learned: $\gamma, \beta$
<b>Output:</b> $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$
$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$

- **What happens if we normalize the inputs to the sigmoid function?**
- **We may lose representational power (e.g. ability to represent non-linear functions)**

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- **Scaling and shifting restores the original representational power.**

# Batch Normalization

## Batch Normalization:

- In every layer, normalize each feature in the mini-batch to have zero-mean and the variance of 1.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- **Scaling and shifting restores the original representational power.**
- **Two parameters gamma and beta learned during training.**

# Batch Normalization

## Backpropagation:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \quad // \text{ scale and shift}$$



# Batch Normalization

## Backpropagation:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

<b>Input:</b> Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$ ; Parameters to be learned: $\gamma, \beta$	
<b>Output:</b> $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$	
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$	// mini-batch mean
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$	// mini-batch variance
$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$	// normalize
$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$	// scale and shift

**Which gradients do we need to compute during a backward pass?**

# Batch Normalization

## Backpropagation:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \quad // \text{ scale and shift}$$

Need to send backwards.

# Batch Normalization

## Backpropagation:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  // mini-batch mean

$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$  // mini-batch variance

$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$  // normalize

$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$  // scale and shift

Intermediate  
gradients

# Batch Normalization

## Backpropagation:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Parameter gradients

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

**Intermediate gradients:**

$$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

## Intermediate gradients:

$$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot (x_i - \mu_{\mathcal{B}}) \cdot \frac{-1}{2} (\sigma_{\mathcal{B}}^2 + \epsilon)^{-3/2}$$

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

## Intermediate gradients:

$$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot (x_i - \mu_{\mathcal{B}}) \cdot \frac{-1}{2} (\sigma_{\mathcal{B}}^2 + \epsilon)^{-3/2}$$

$$\frac{\partial \ell}{\partial \mu_{\mathcal{B}}} = \left( \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_{\mathcal{B}})}{m}$$

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

## Gradient to send backwards:

$$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{2(x_i - \mu_{\mathcal{B}})}{m} + \frac{\partial \ell}{\partial \mu_{\mathcal{B}}} \cdot \frac{1}{m}$$



# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

## Parameter gradients:

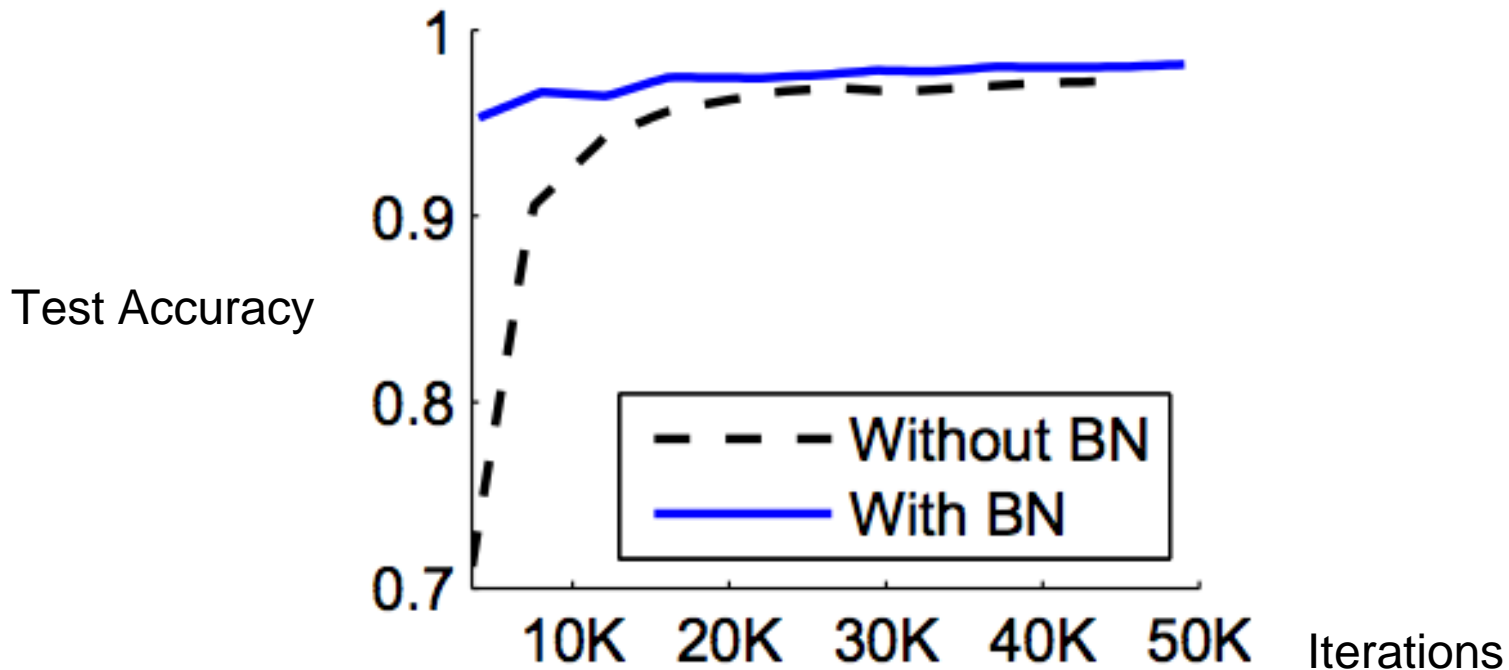
$$\frac{\partial \ell}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \hat{x}_i$$

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i}$$

# Batch Normalization

## Batch Normalization:

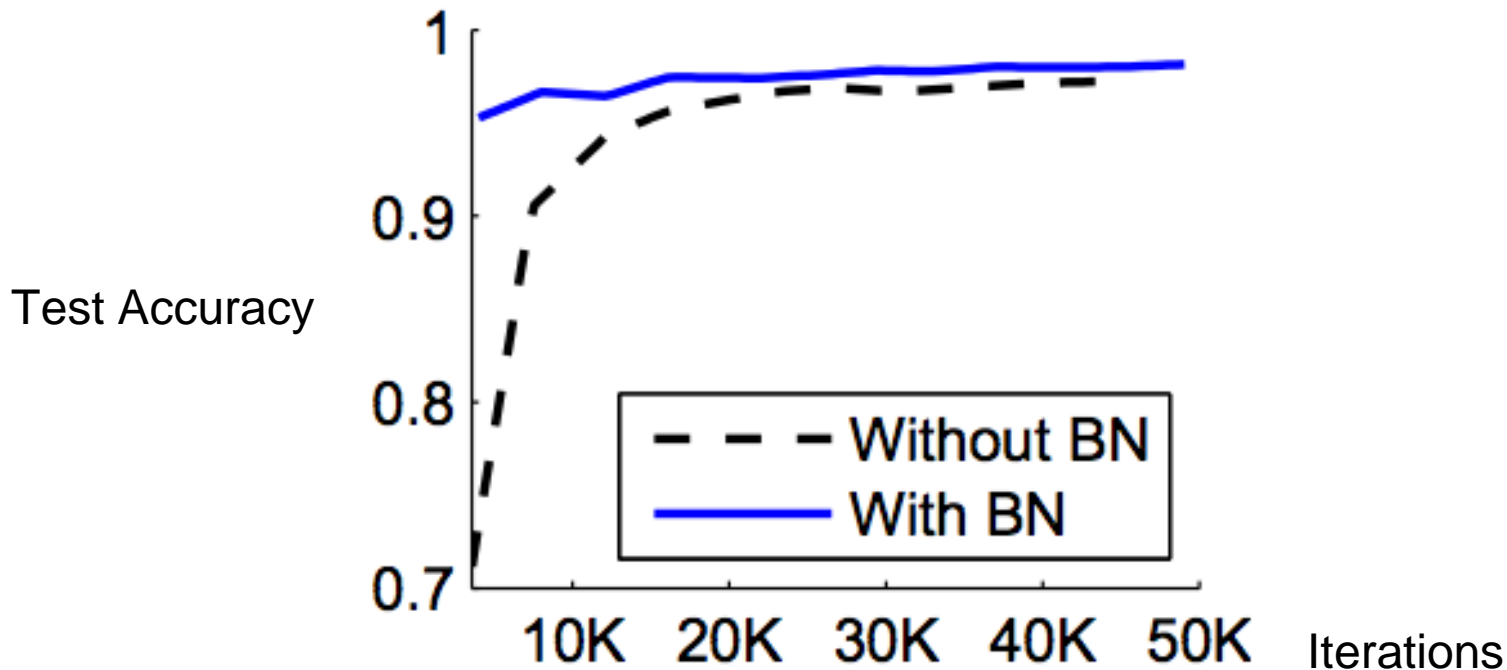
- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.



# Batch Normalization

## Batch Normalization:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.

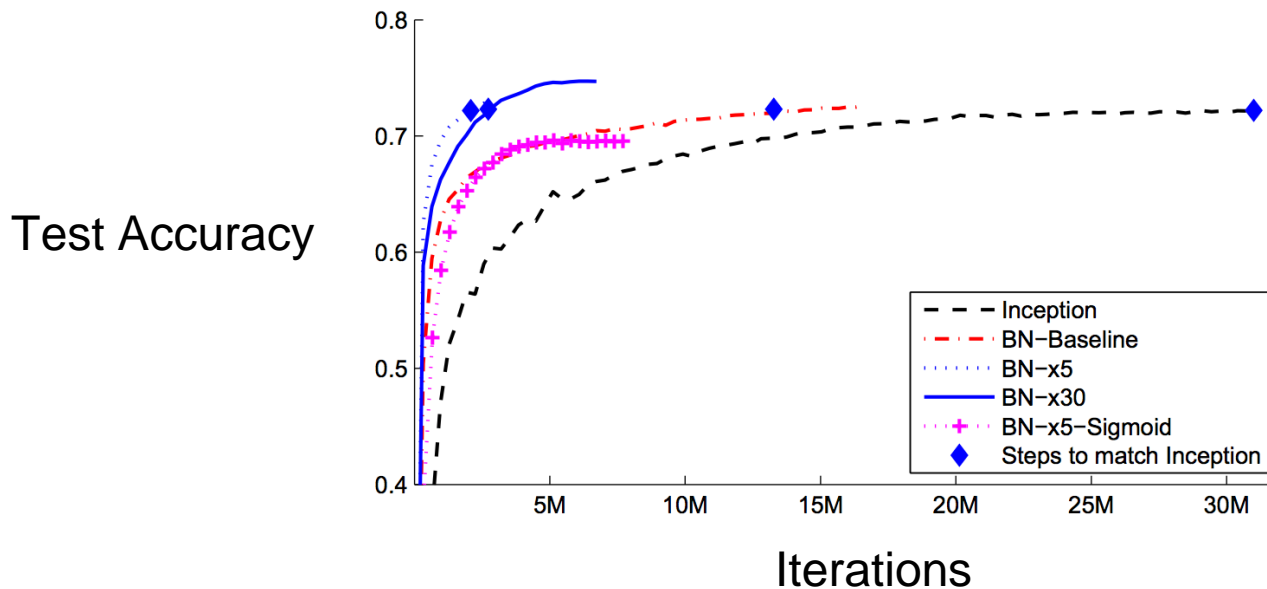


**Batch normalization makes training much faster!**

# Batch Normalization

## Batch Normalization:

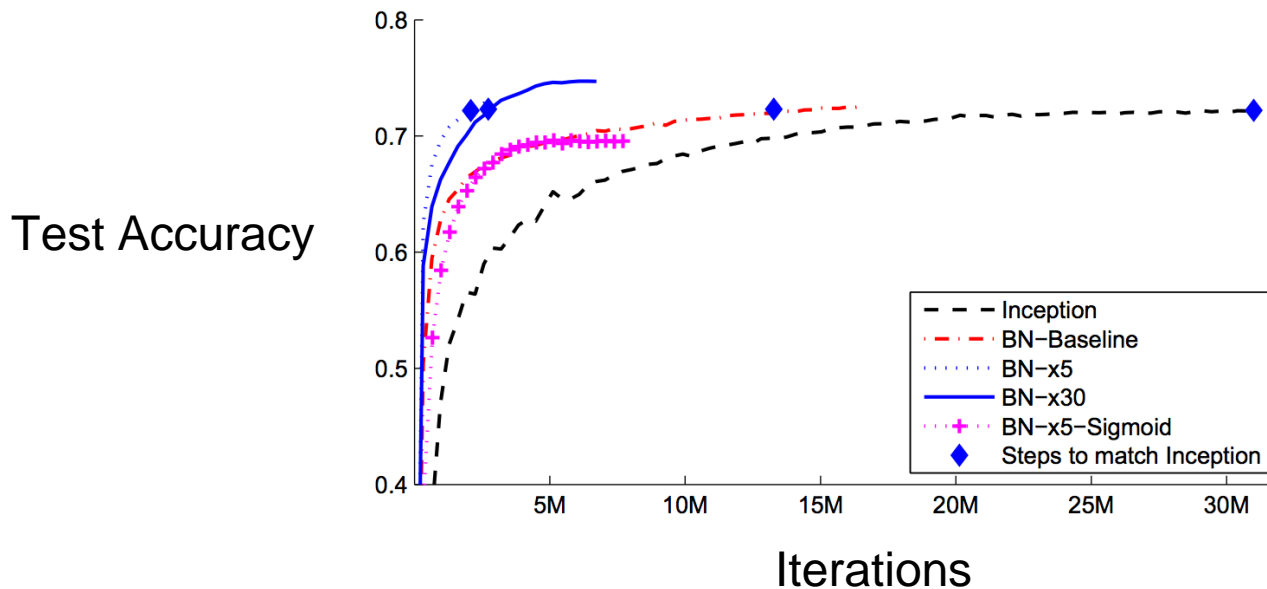
- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.



# Batch Normalization

## Batch Normalization:

- Unlike many other normalization schemes, batch normalization can be easily incorporated into backprop.



**Batch normalization makes training much faster!**

# Summary

## Key Take-Aways:

- Specifying the right loss (e.g. cross entropy) is important.
- Picking the right activation function (e.g. RELU) is also imperative.
- We can reduce vanishing gradient problem via deep supervision.
- Reducing the learning rate, and clipping gradients helps to prevent exploding gradient problem.
- Momentum methods allow faster and more stable learning.
- Batch normalization significantly speeds up the training.