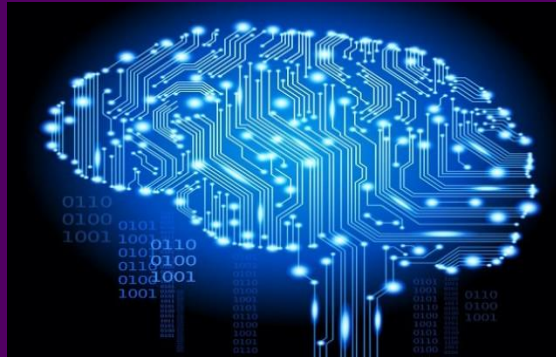


# Deep Learning

## MSiA 490-30



---

Theory and Applications

---



NORTHWESTERN  
UNIVERSITY



NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

## Feedback: News and Ideas



~5-10 min to share ideas/news



## Assignment schedule

- Assignment 1: Due April 26<sup>th</sup>
  - Basic NN concepts
  - Optimization
  - Discuss with 1 partner, submit individually
- Assignment 2 (Group): Due May 18<sup>th</sup>
  - Computer Vision
  - Preparation for project
- Mini Quiz 1: May 8<sup>th</sup>: Basic NN, optimization and generalization multiple choice
- Mini Quiz 2: June 9<sup>th</sup> (2<sup>nd</sup> half): Generative, recurrent and applications multiple choice

April 2017

Su	Mo	Tu	We	Th	Fr	Sa
26	27	28	29	30	31	1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29
30	1	2	3	4	5	6

May 2017

Su	Mo	Tu	We	Th	Fr	Sa
30	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	1	2	3
4	5	6	7	8	9	10

June 2017

Su	Mo	Tu	We	Th	Fr	Sa
28	29	30	31	1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	1
2	3	4	5	6	7	8



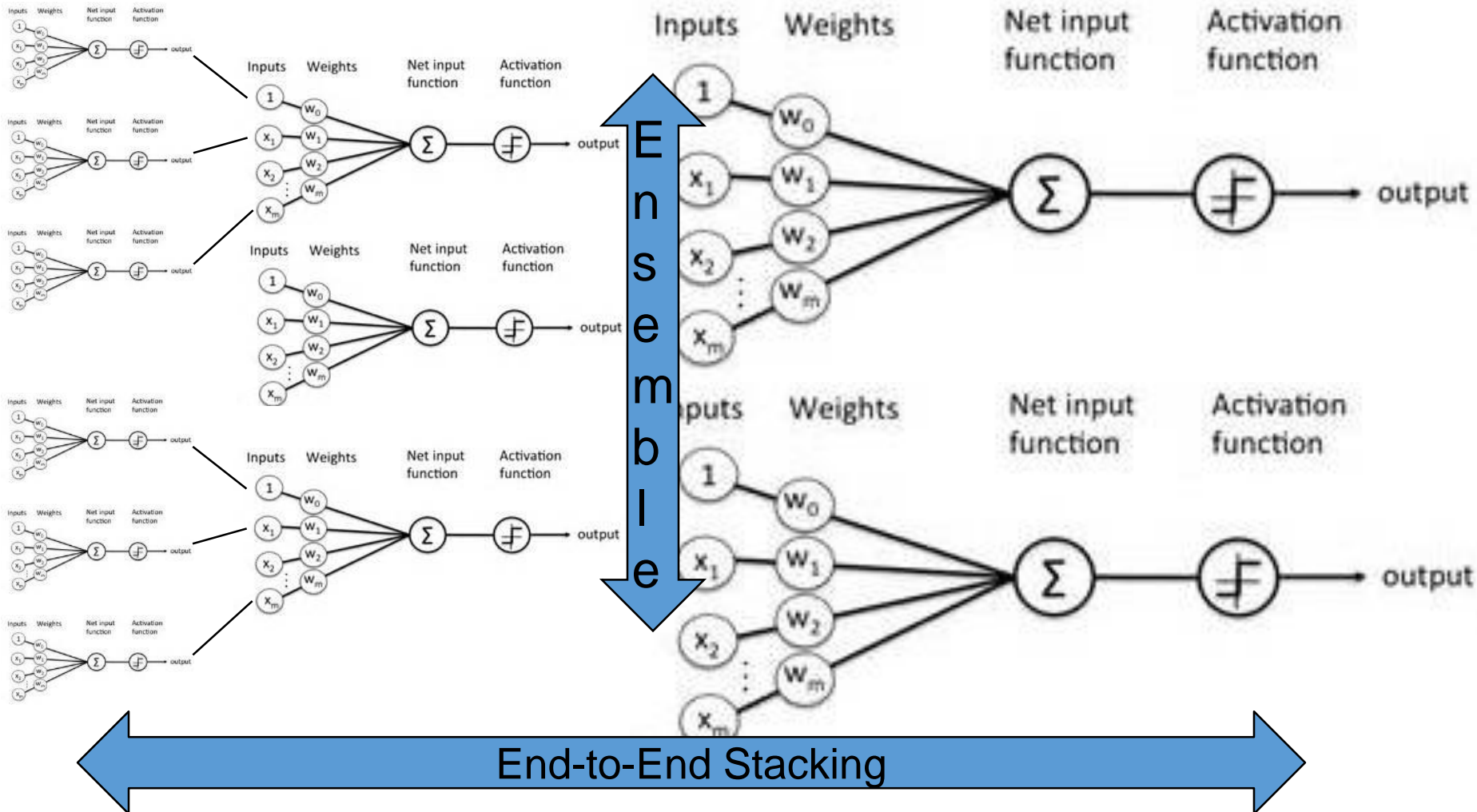
## Last time

- Introduced multilayer nets
  - Apply optimization rules in layers
- This time
  - Why Optimization may slow down/fail
  - Some simple workarounds
  - Advanced optimization methods



NORTHWESTERN  
UNIVERSITY

## “Layered” Logistic Regression



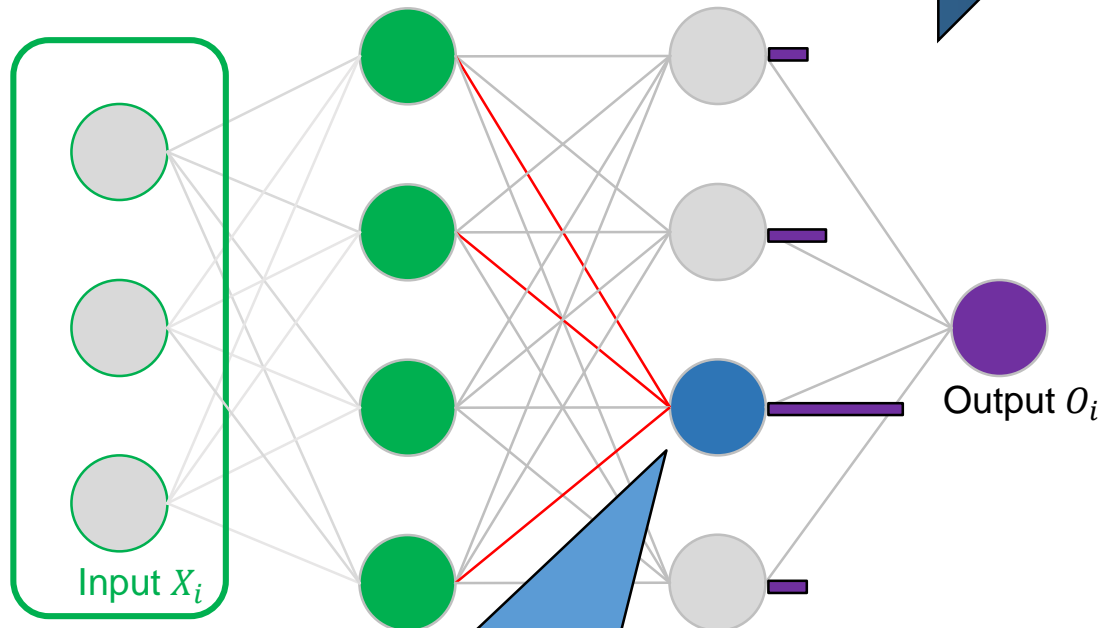


NORTHWESTERN  
UNIVERSITY

## Forward propagation

Forward propagation

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24



Each activation is simply a stacked logistic regression.

Bar length = level of activation



## Backpropagation

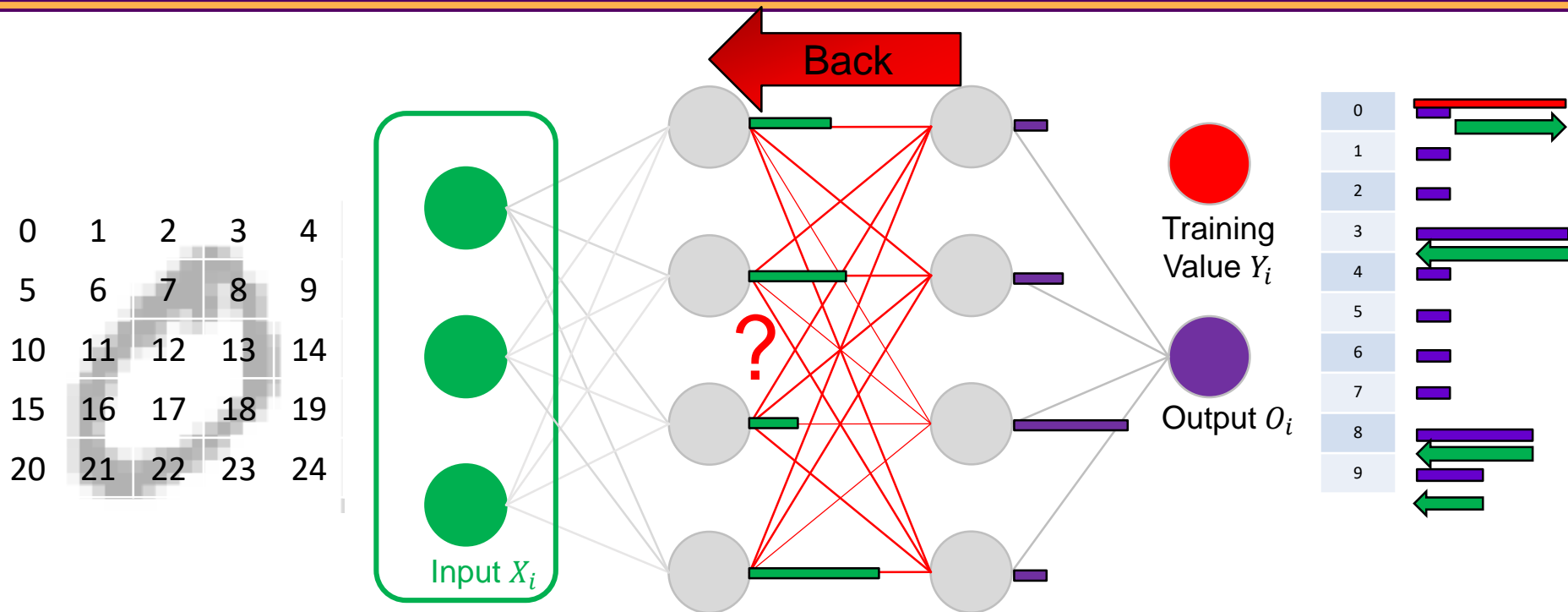
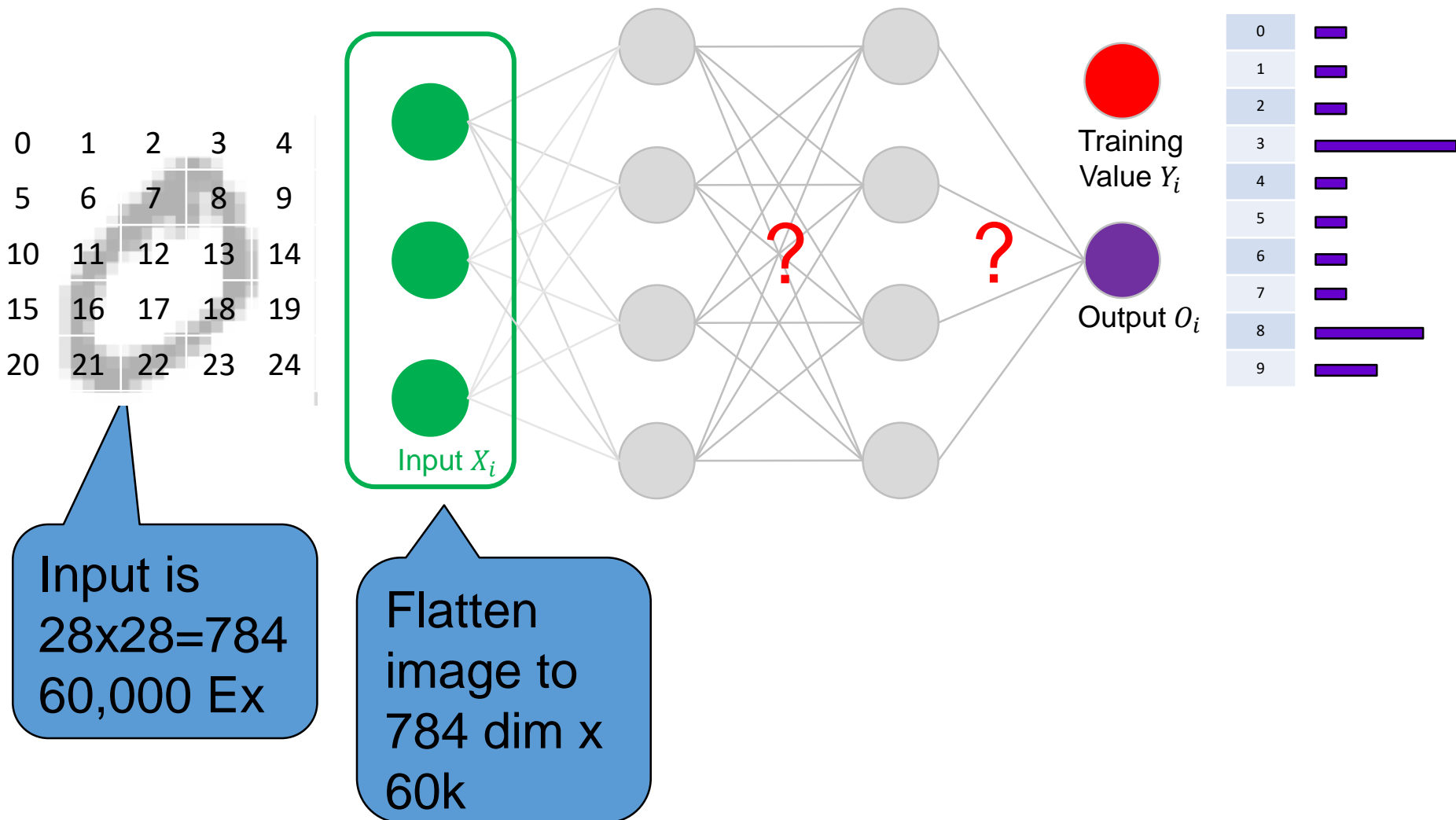


Figure out how to adjust  
input(**green**) weights(**red**) to  
match target activations(**purple**)



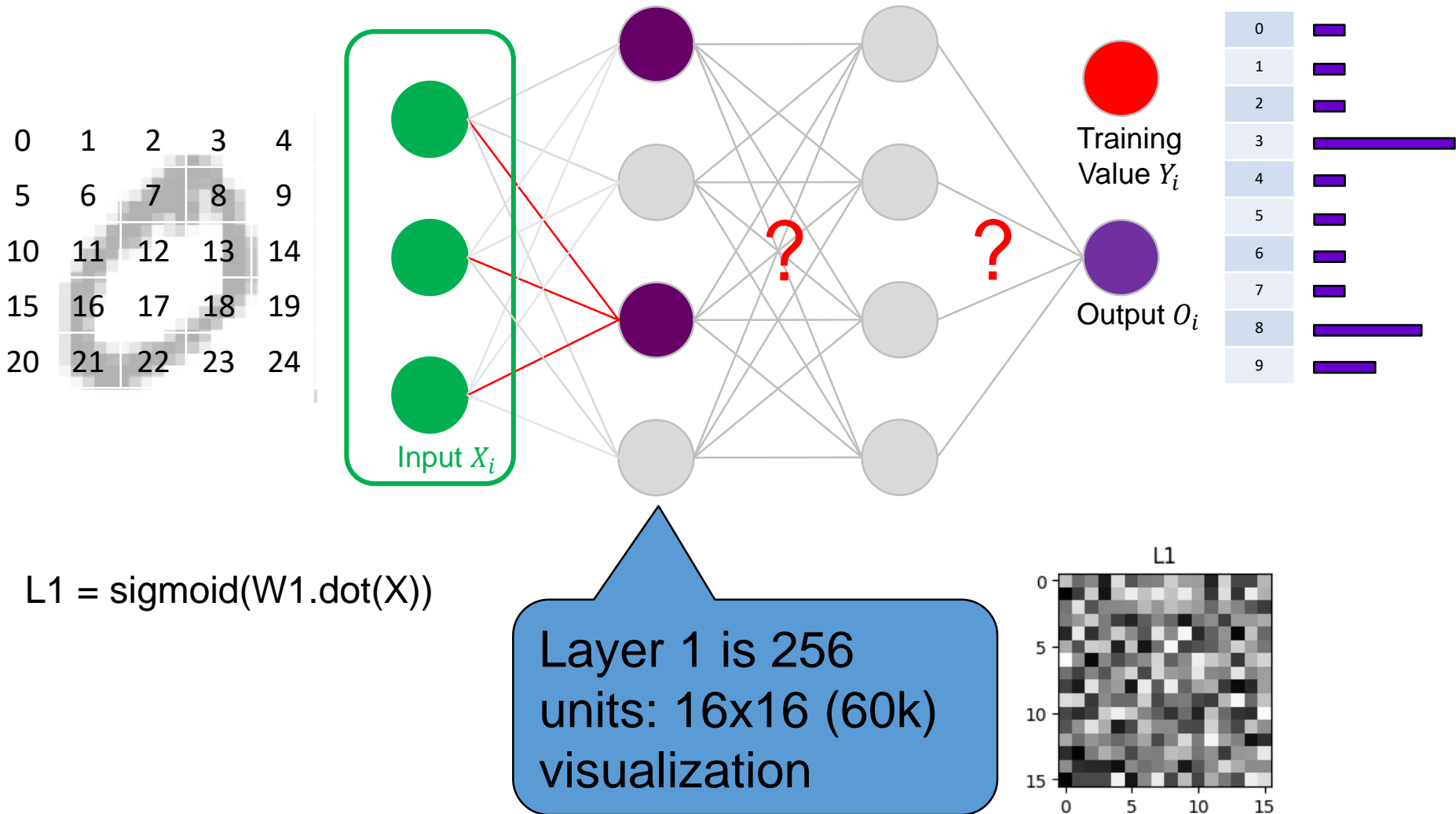
# Feedforward networks





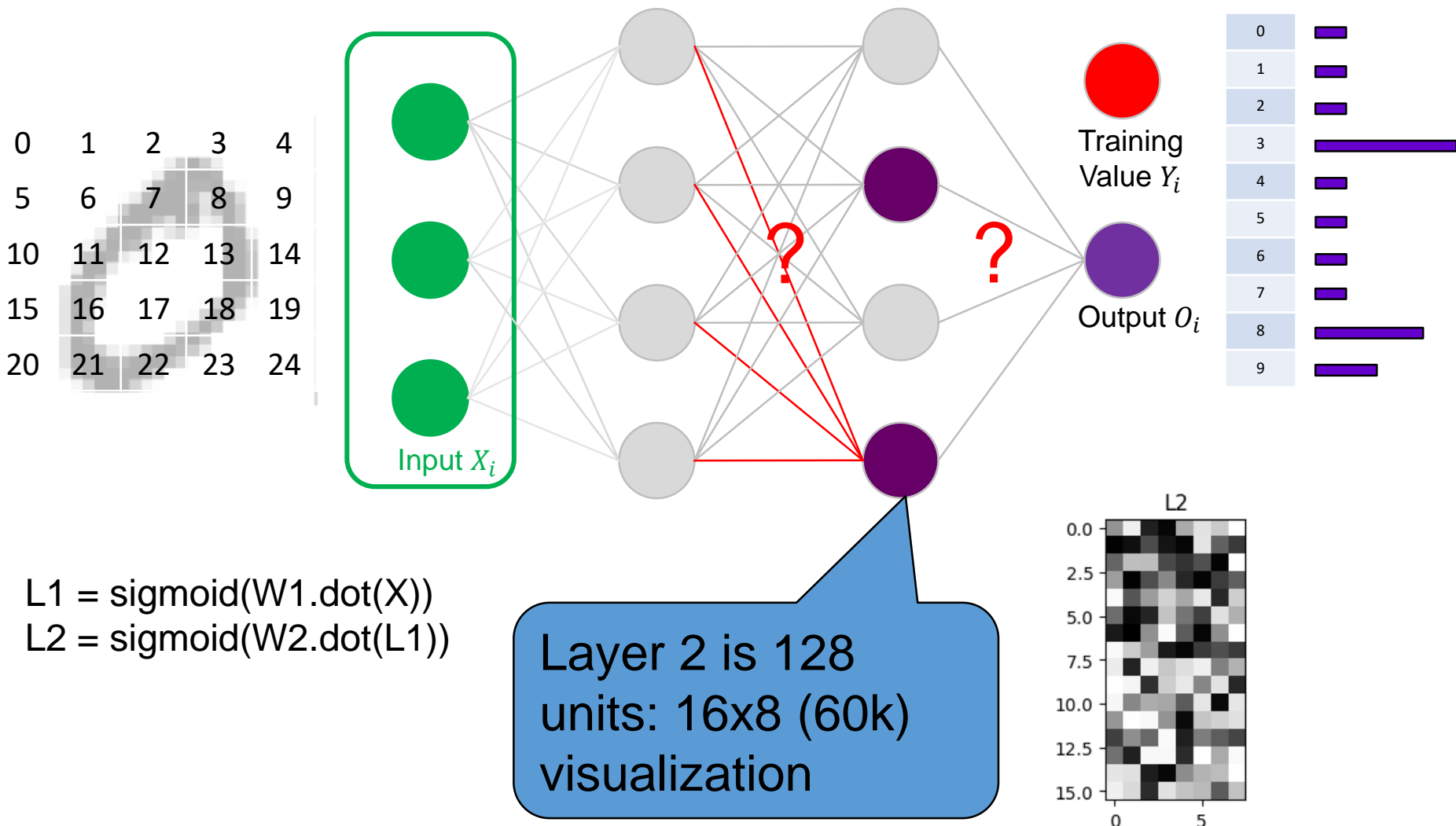


## Visualizing NN





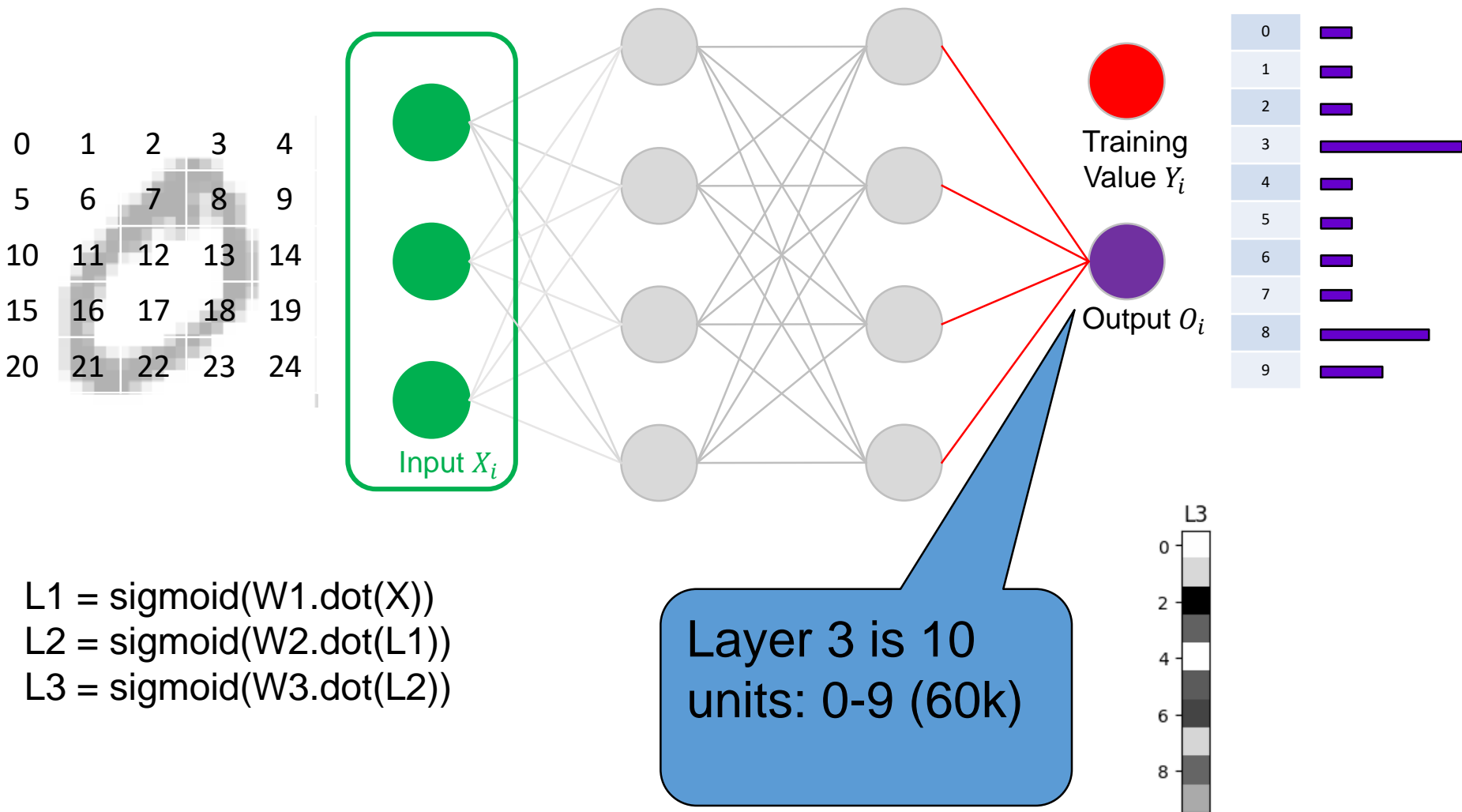
## Visualizing NN





NORTHWESTERN  
UNIVERSITY

## Visualizing NN

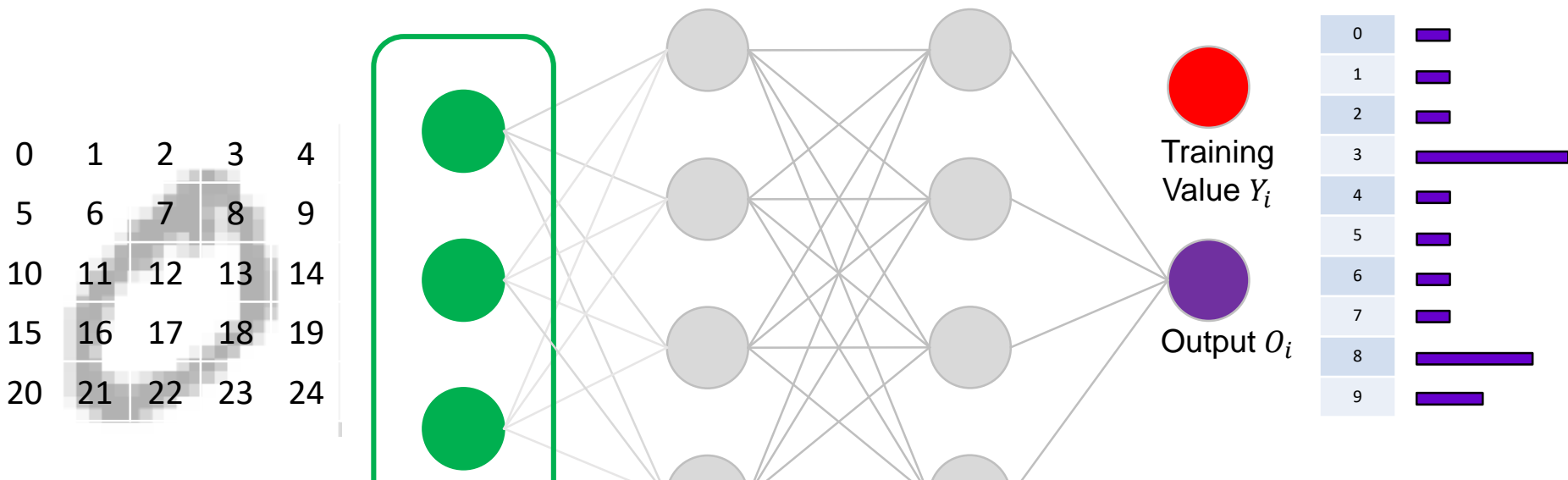




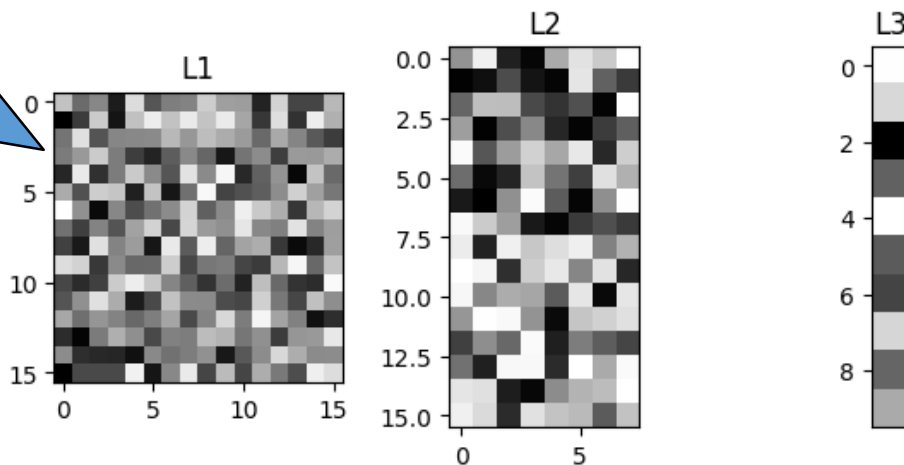
NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

## Visualizing NN



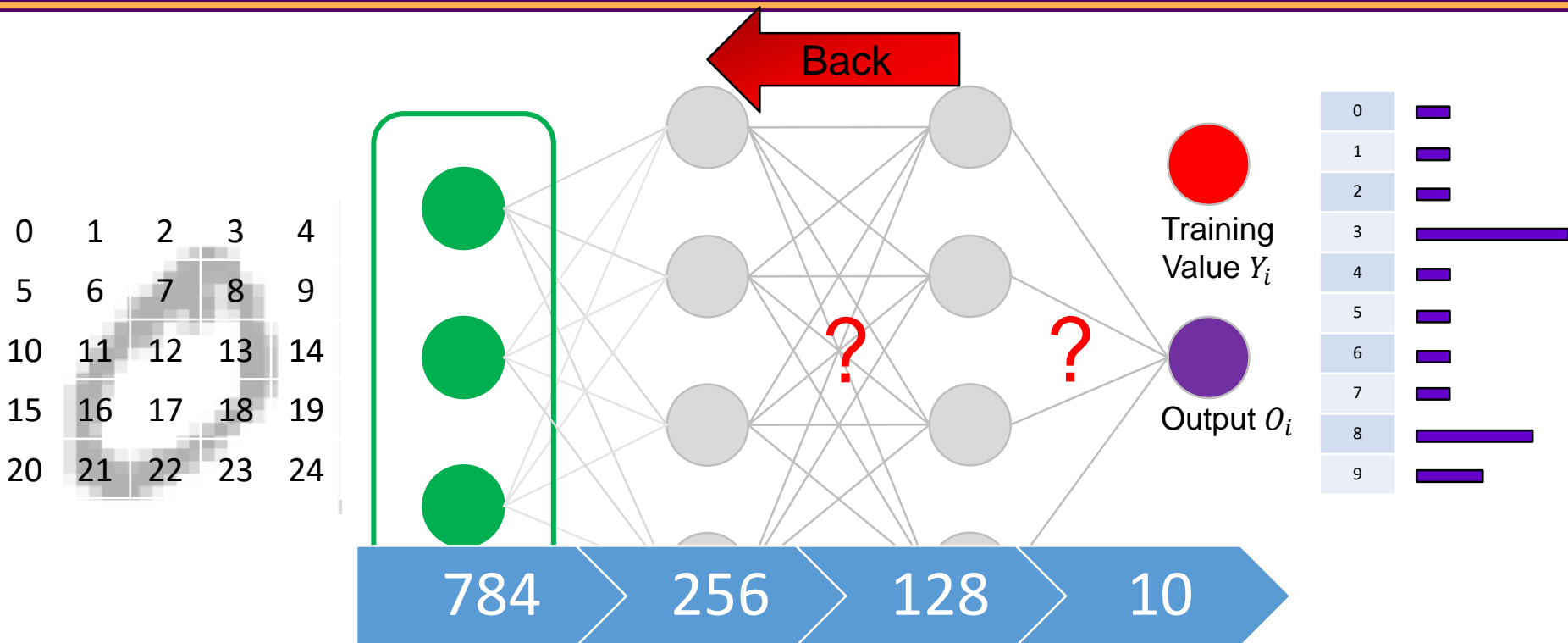
Mean  
activations  
over 60k  
images



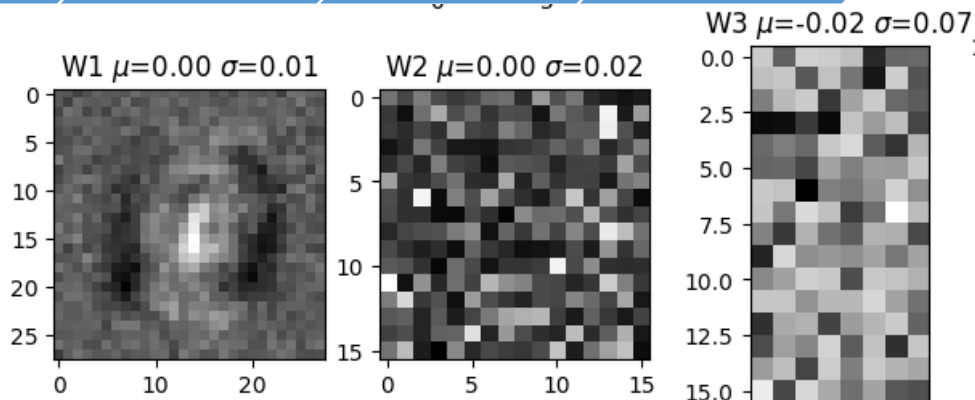


NORTHWESTERN  
UNIVERSITY

## Visualizing NN



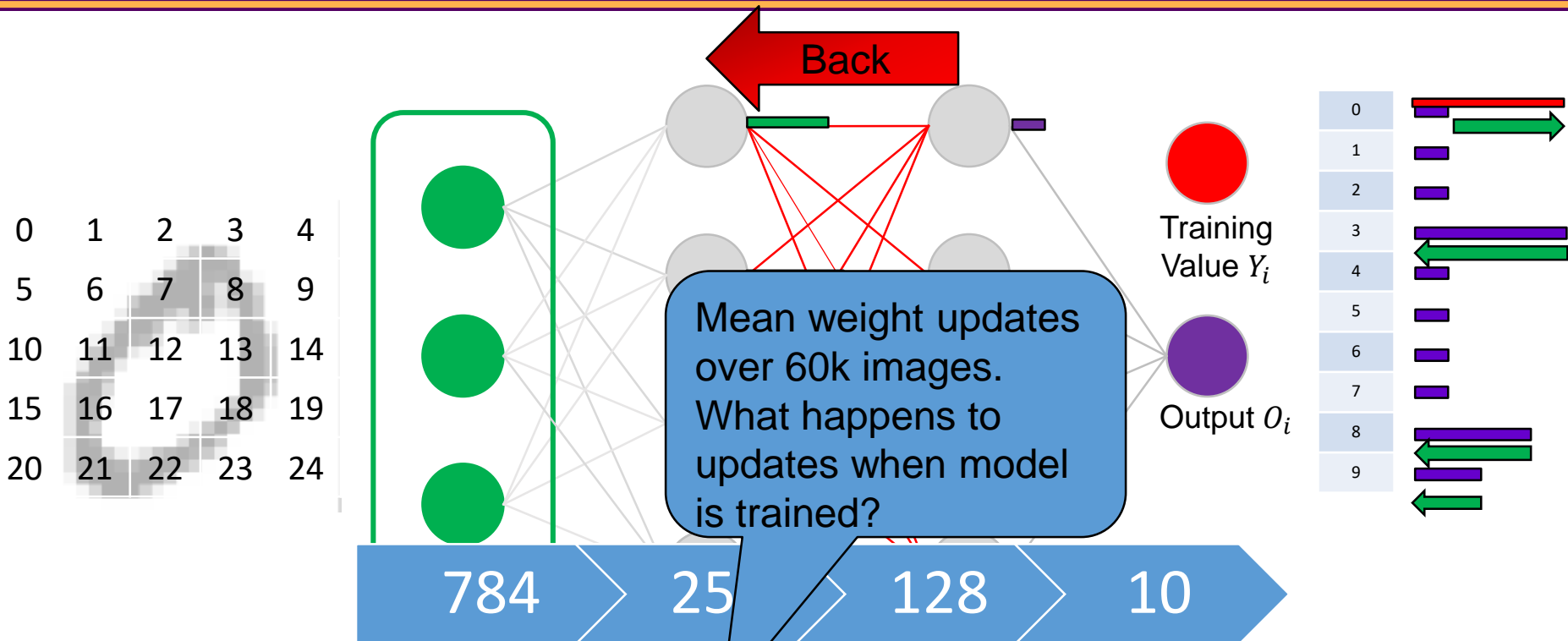
Mean weights over 60k images. Remember:  $w_i x_i$  means input and weights share same shape





NORTHWESTERN  
UNIVERSITY

## Visualizing NN



# Backward pass

$$dW3 = (L3 - T) * L3 * (1 - L3)$$

$$dW2 = W3.T.dot(dW3)*(L2*(1-L2))$$

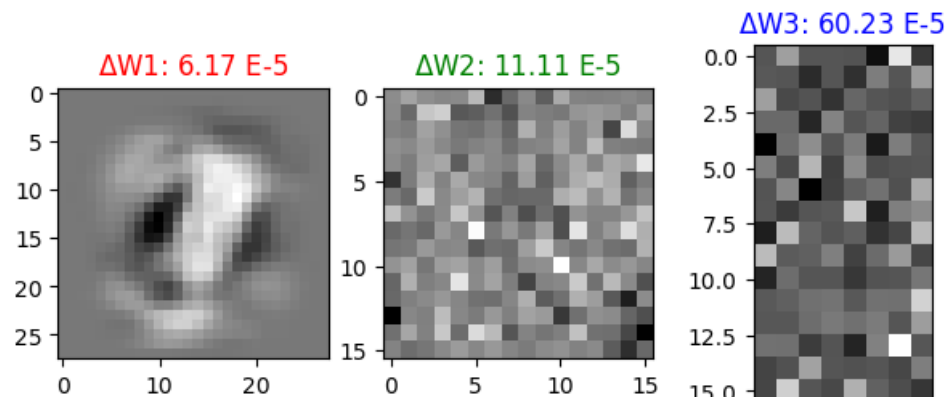
$$dW1 = W2.T.dot(dW2)*(L1*(1-L1))$$

# Update

$$W3 -= lr * np.dot(dW3, L2.T)$$

$$W2 -= lr * np.dot(dW2, L1.T)$$

$$W1 -= lr * np.dot(dW1, X.T)$$

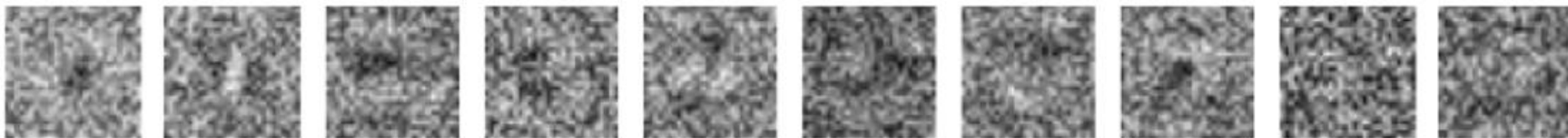




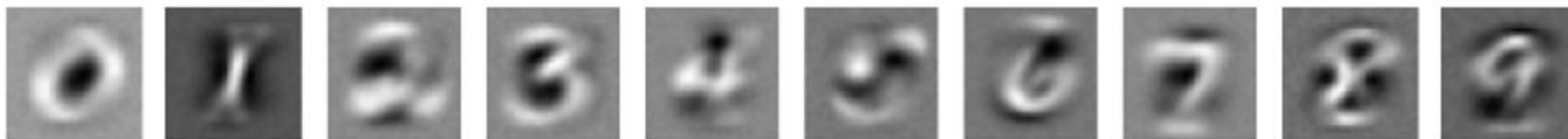
## Effects of Regularization

- L1 regularization encourages stronger feature learning

No regularization

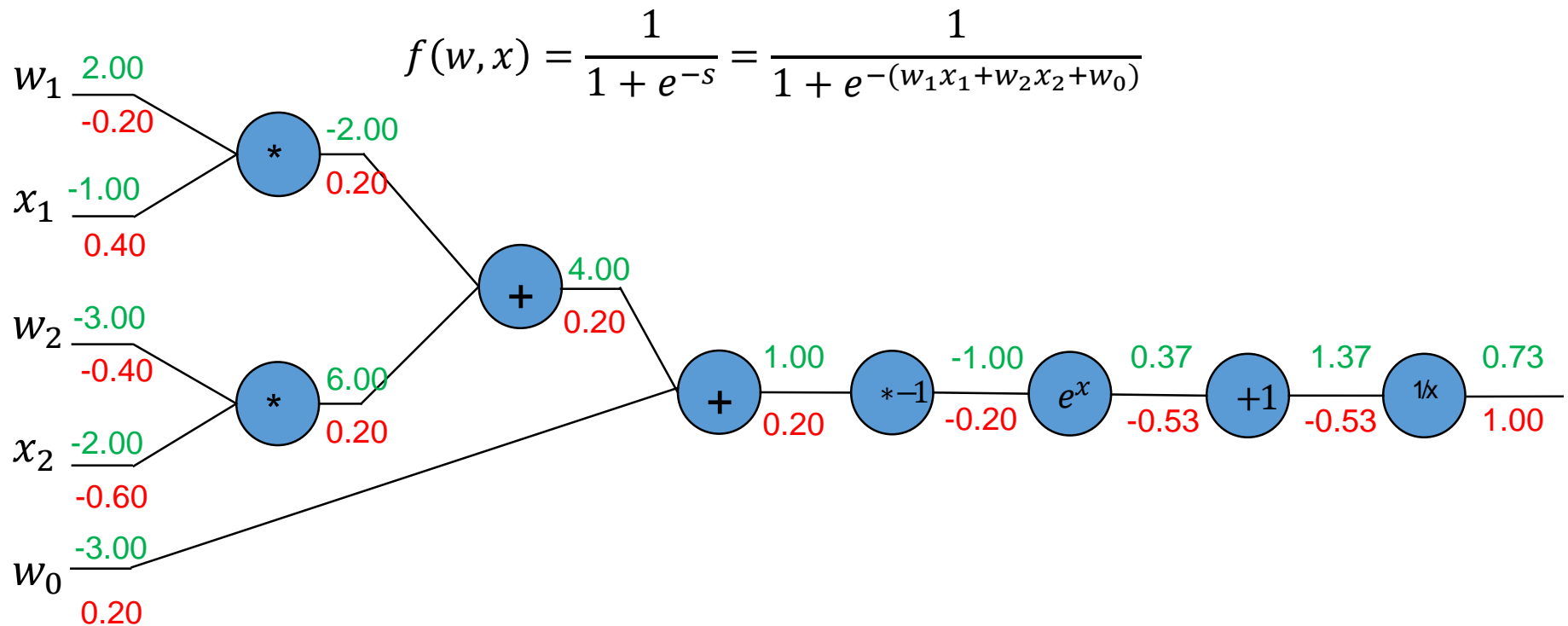


L1 regularization





## Backprop in “depth”



$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

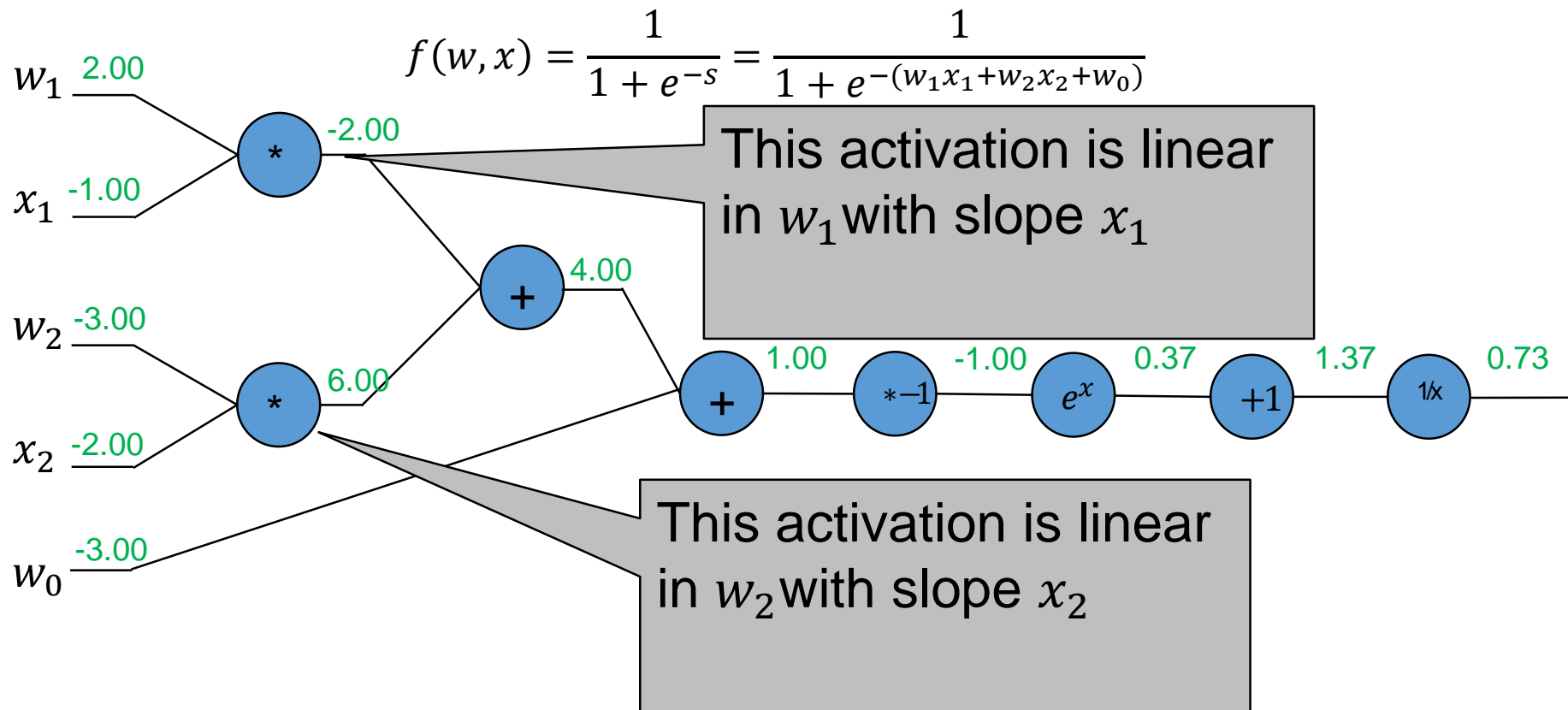
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$





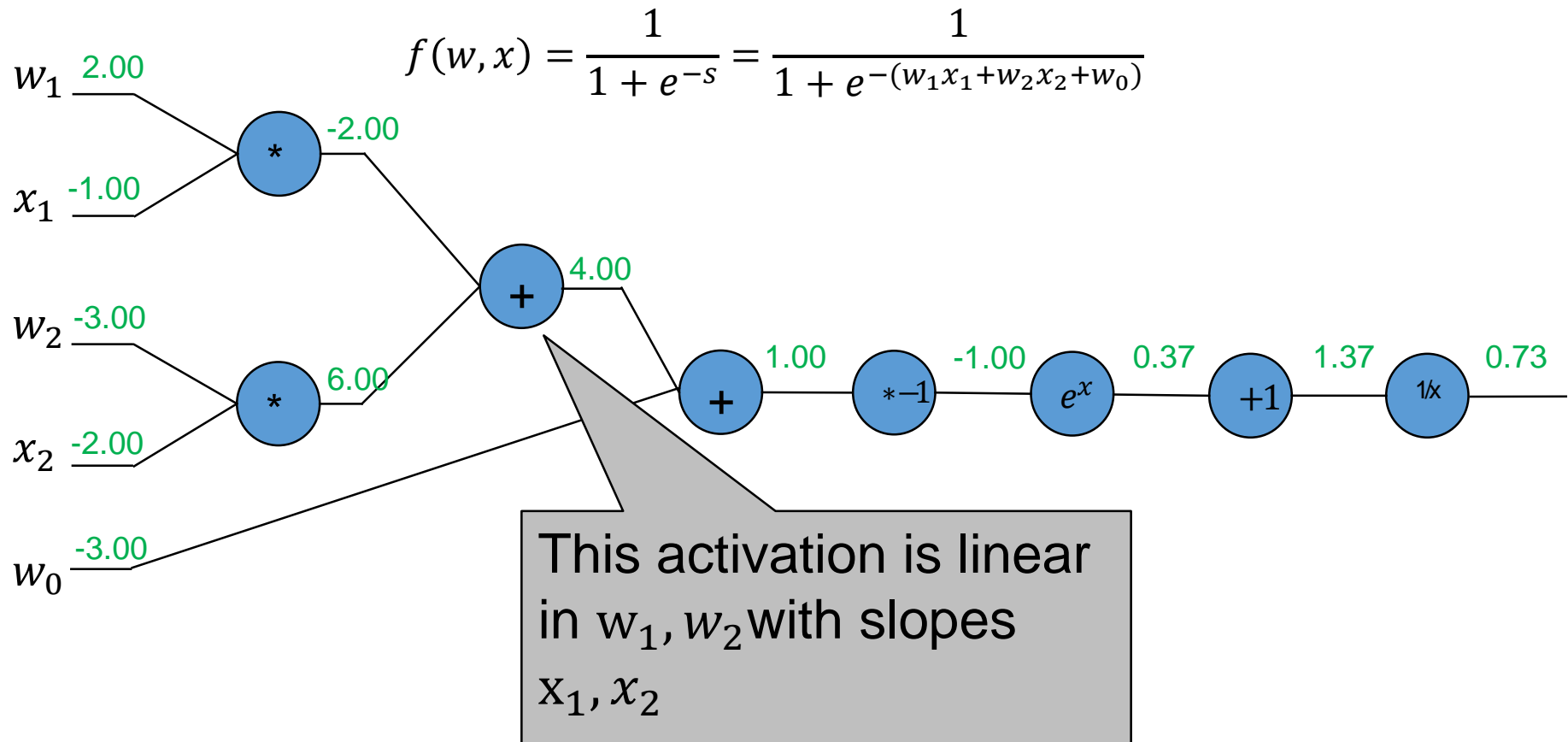
## Backprop in “depth”: Forward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



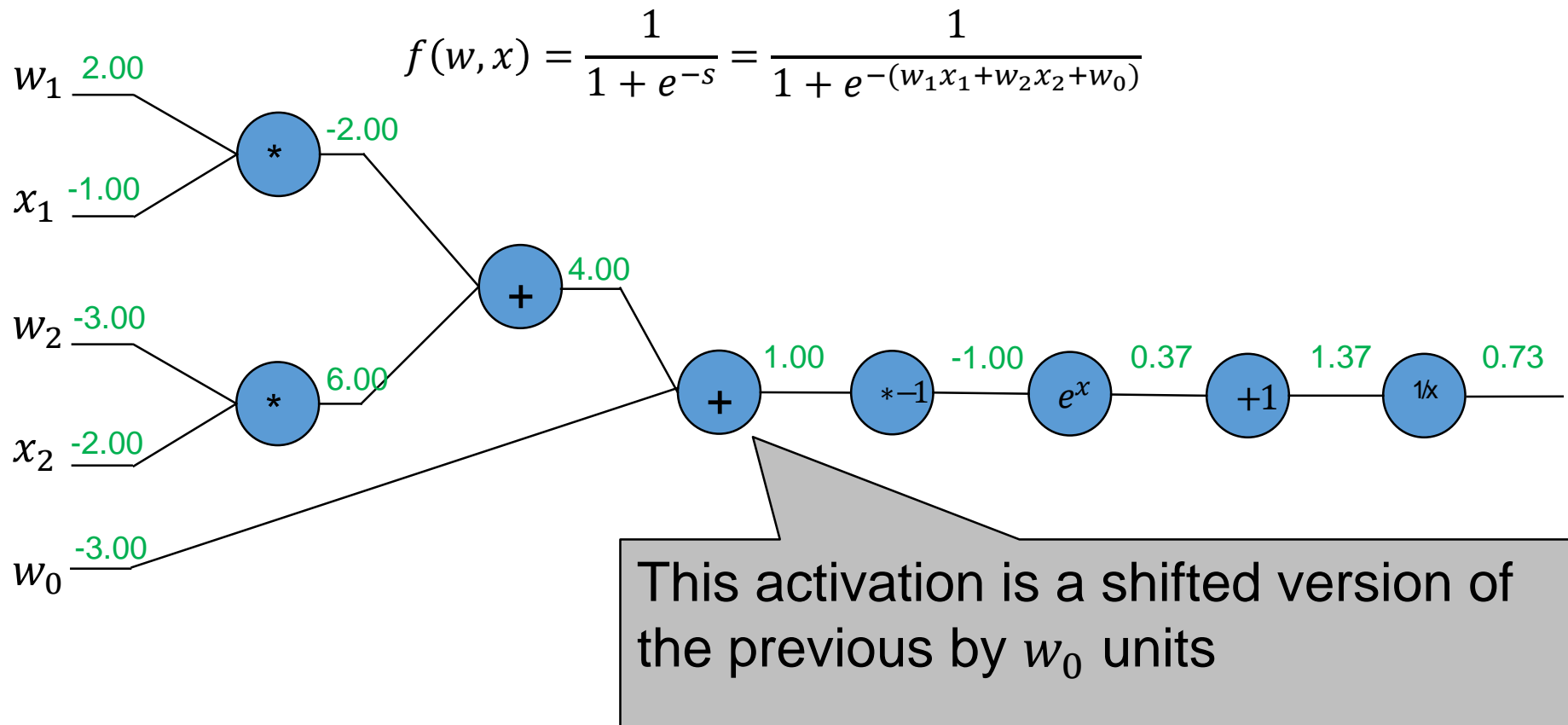
## Backprop in “depth”: Forward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



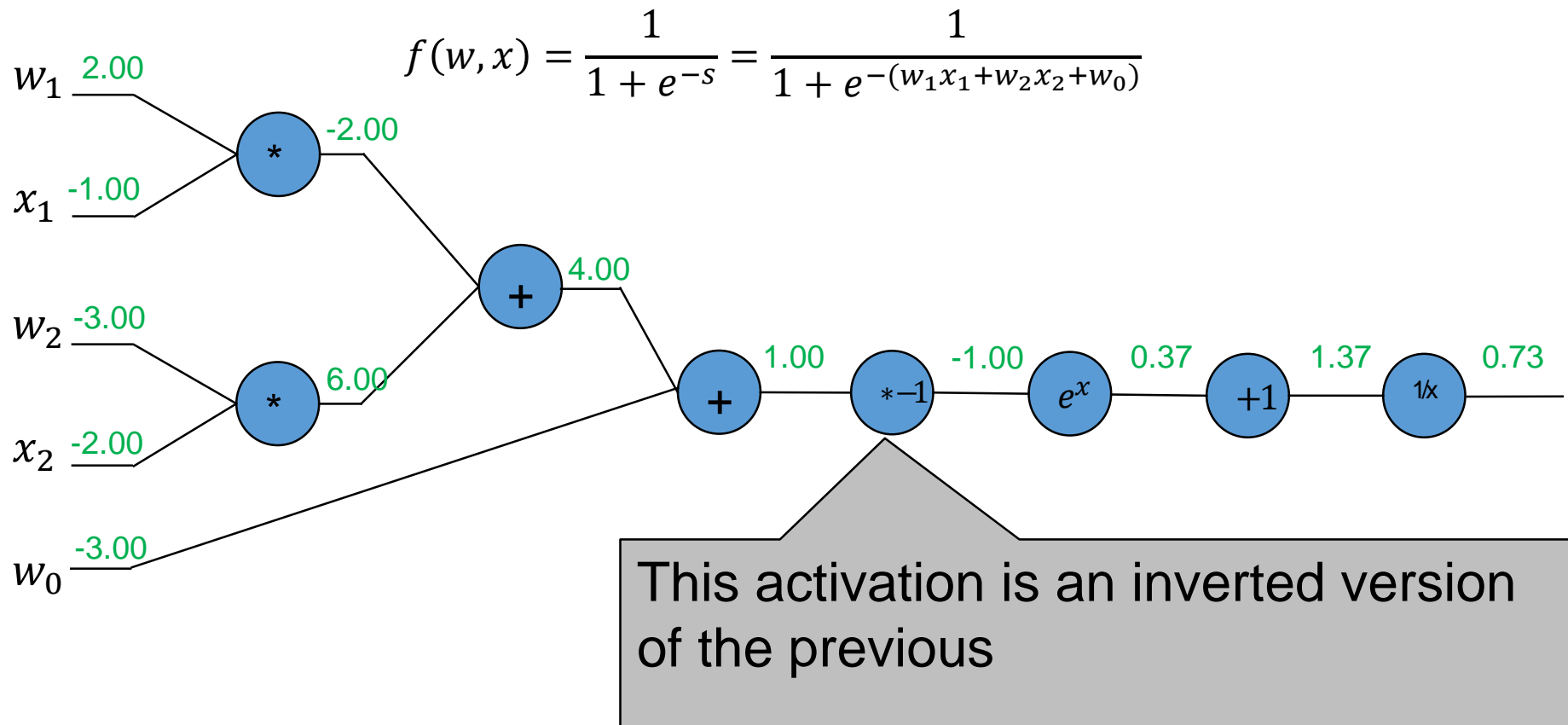
## Backprop in “depth”: Forward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



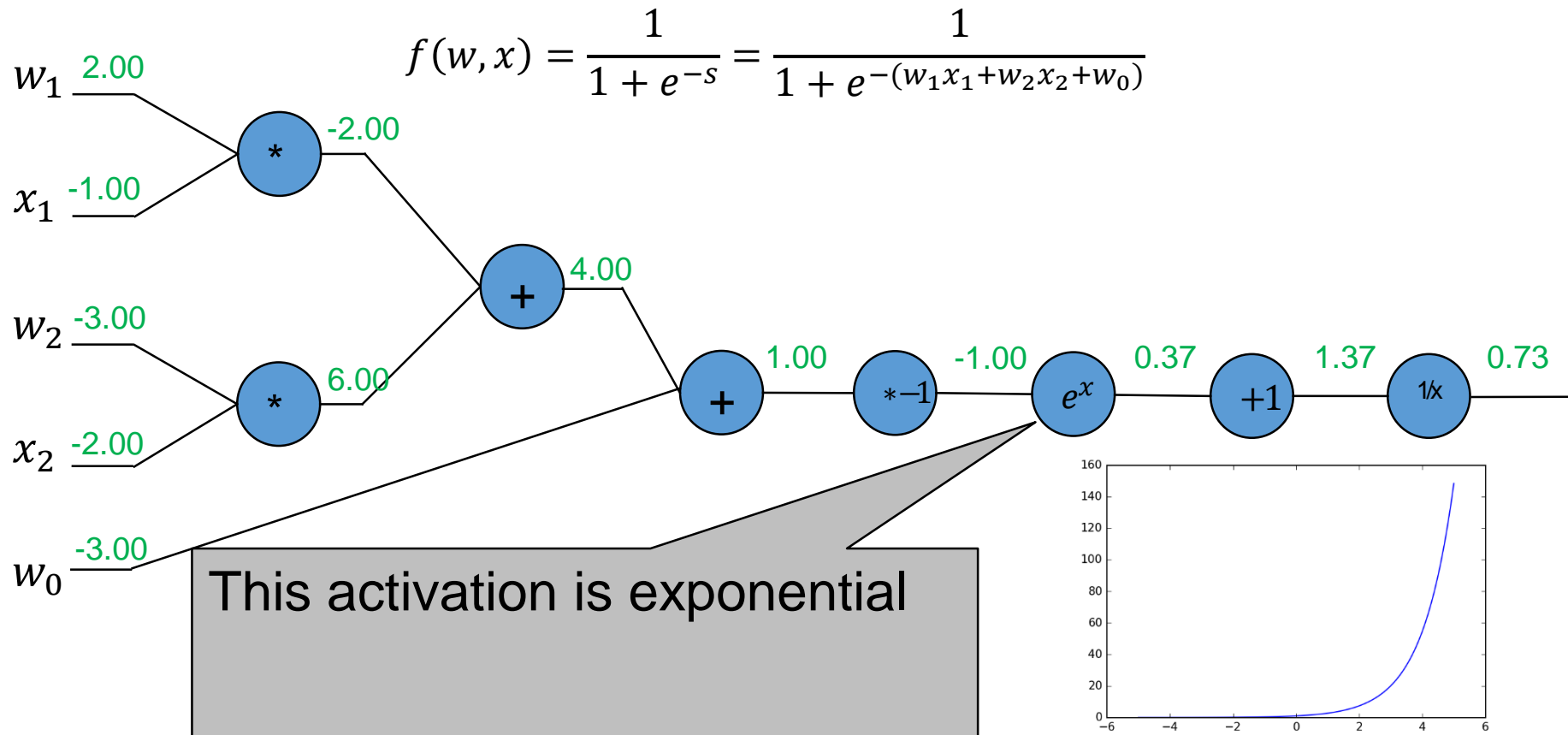
## Backprop in “depth”: Forward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Forward pass



$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

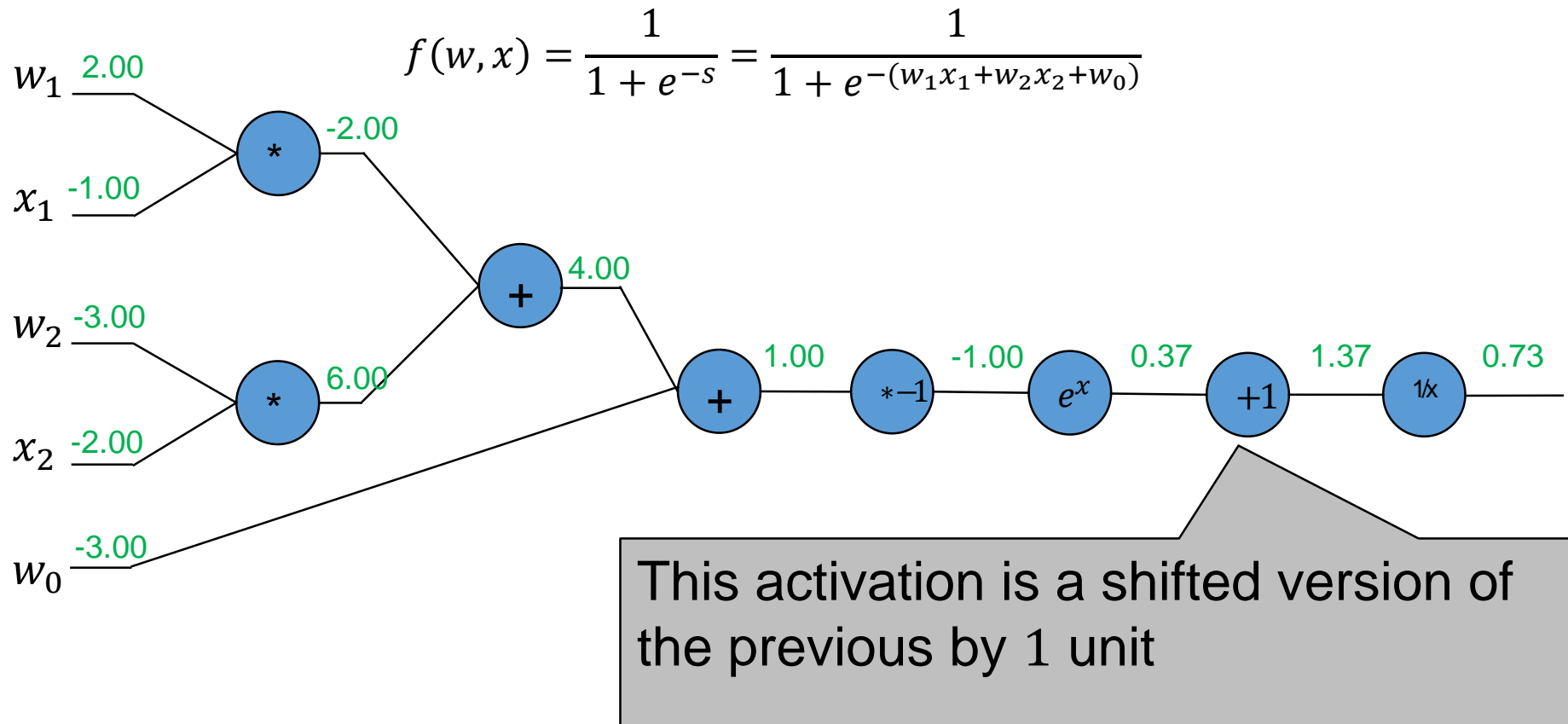
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$



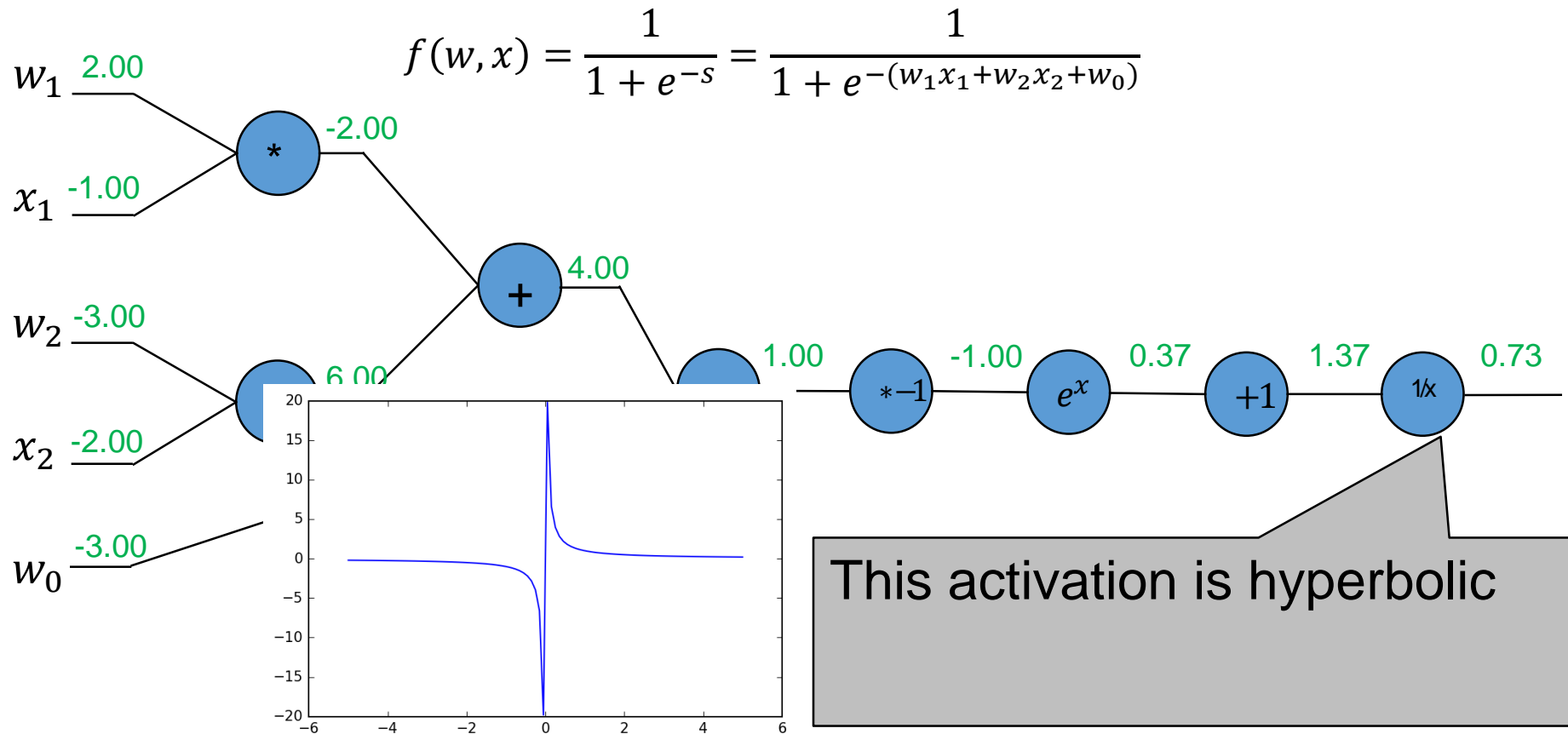
## Backprop in “depth”: Forward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Forward pass



$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

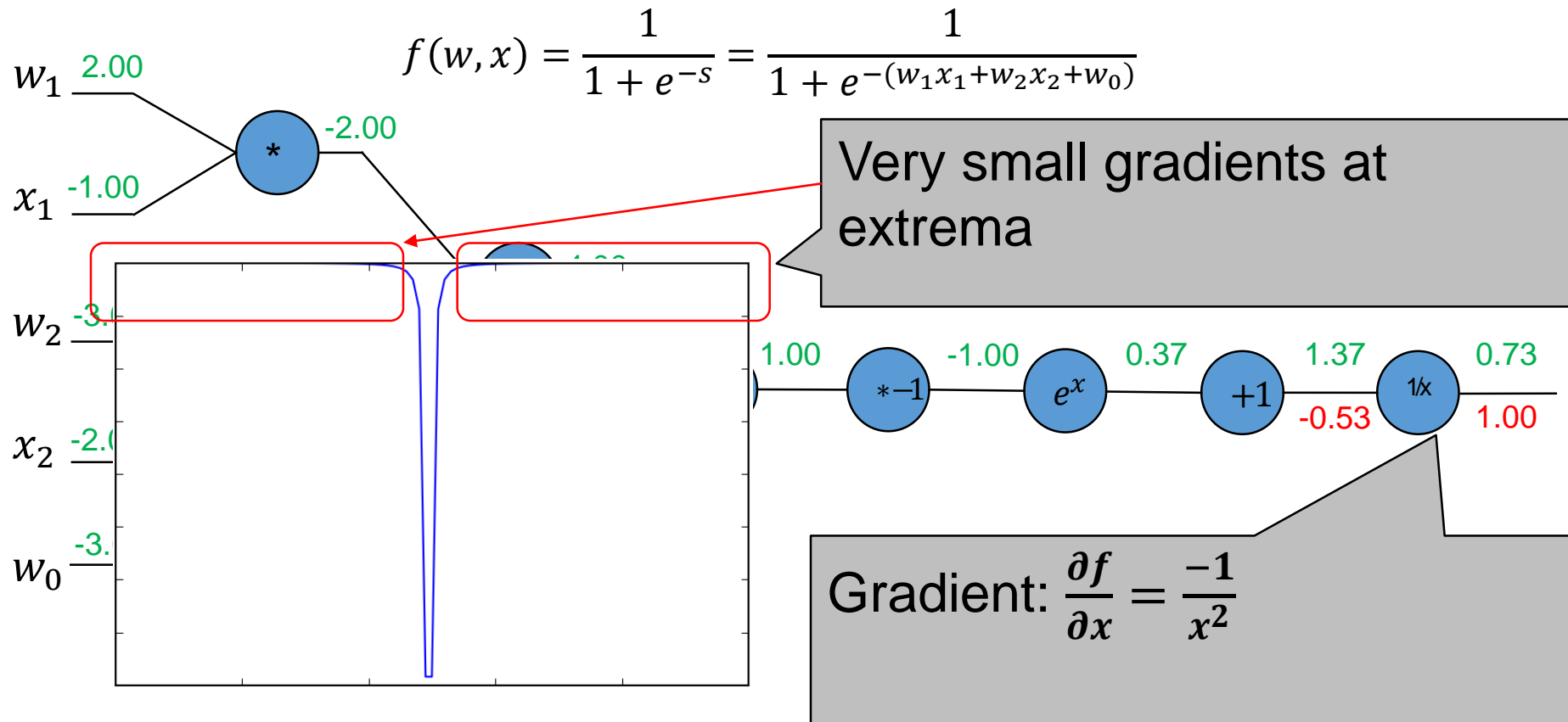
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$



## Backprop in “depth”: Backward pass

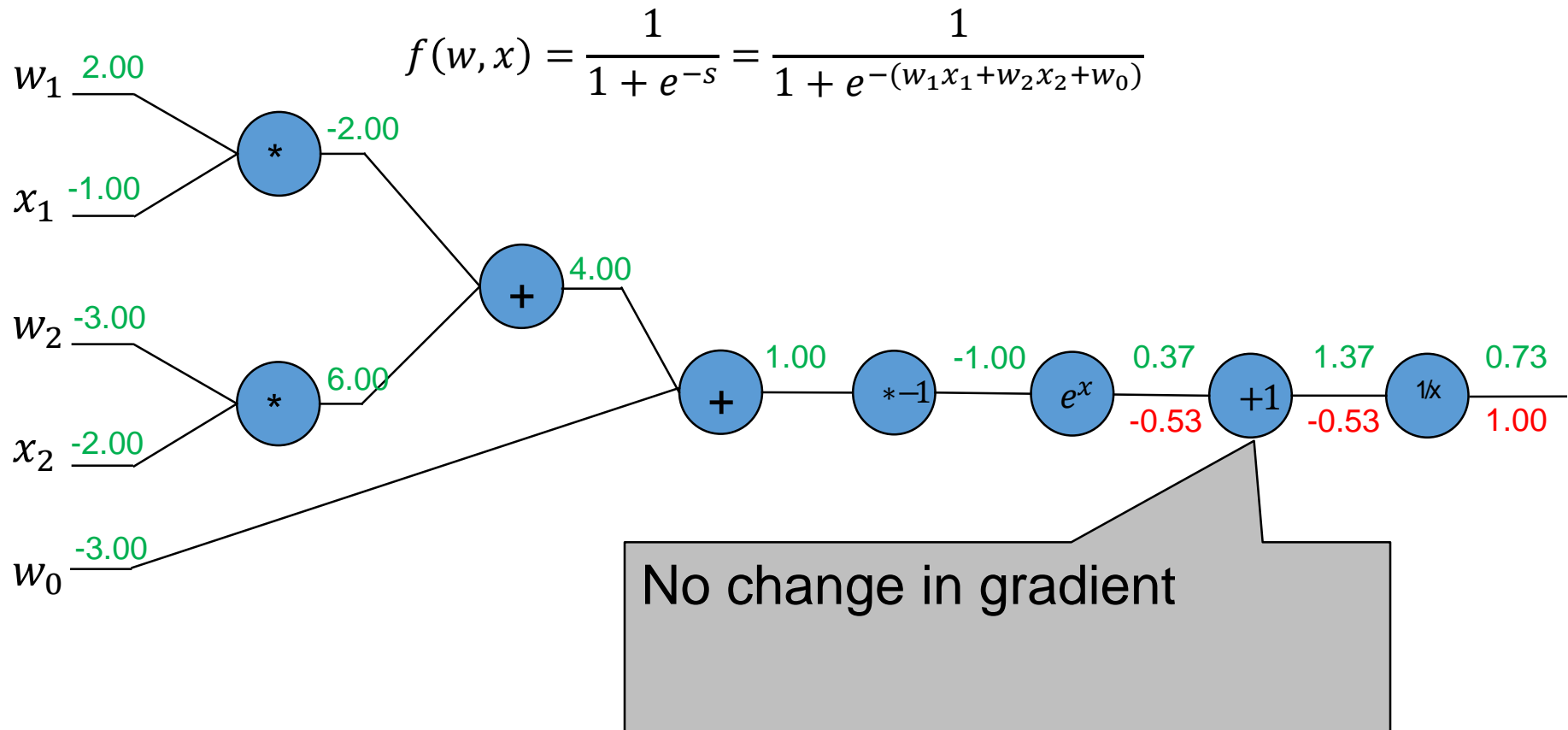


$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$





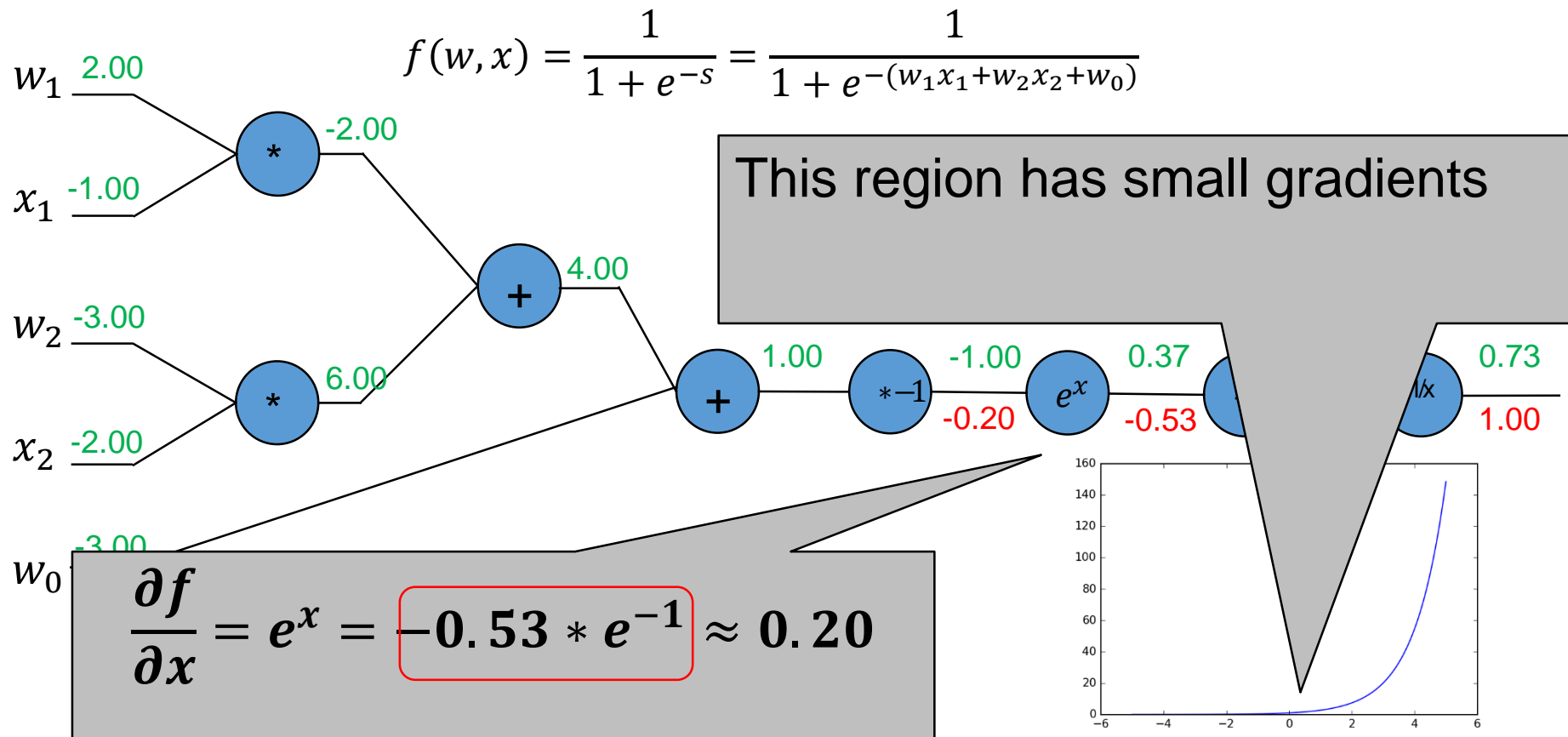
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Backward pass



$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

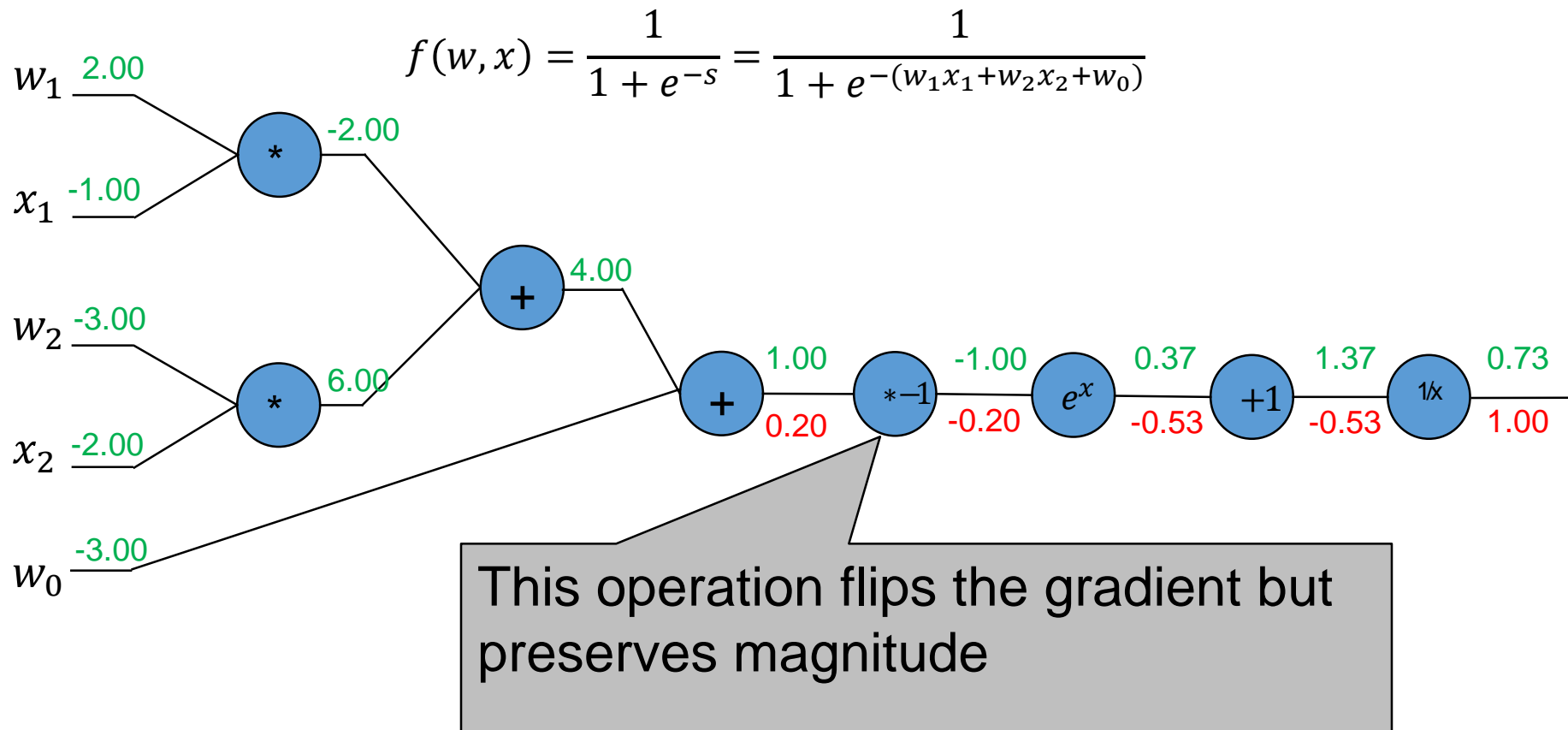
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$



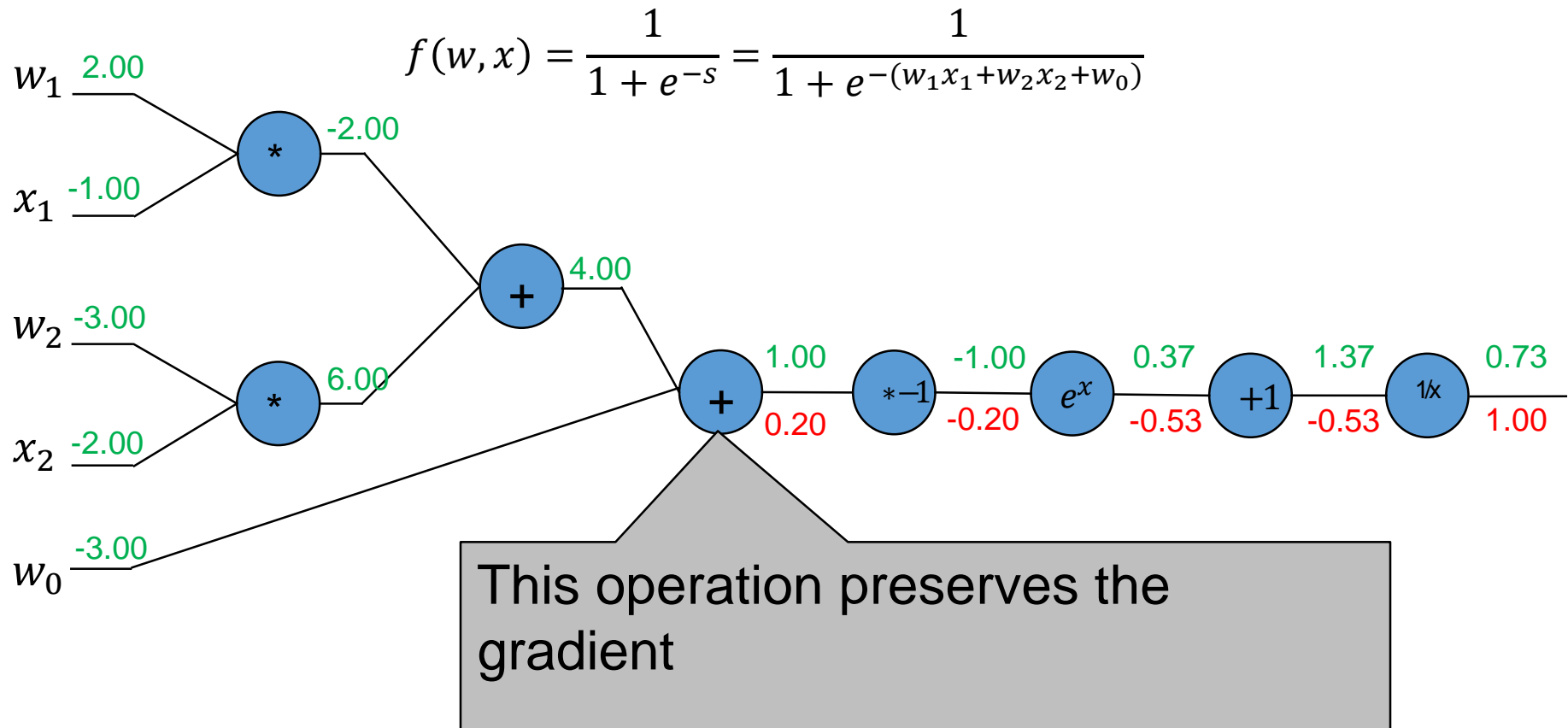
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



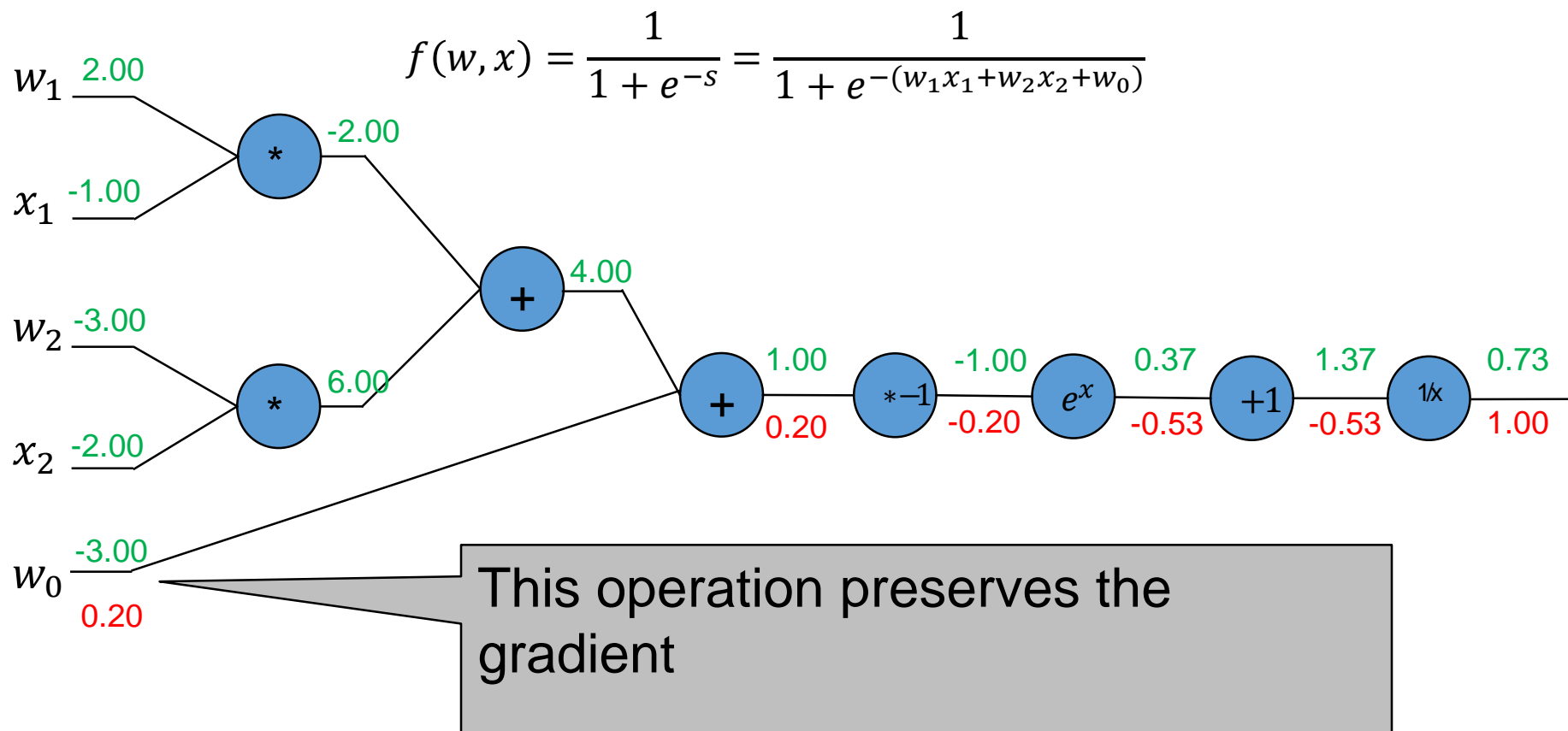
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



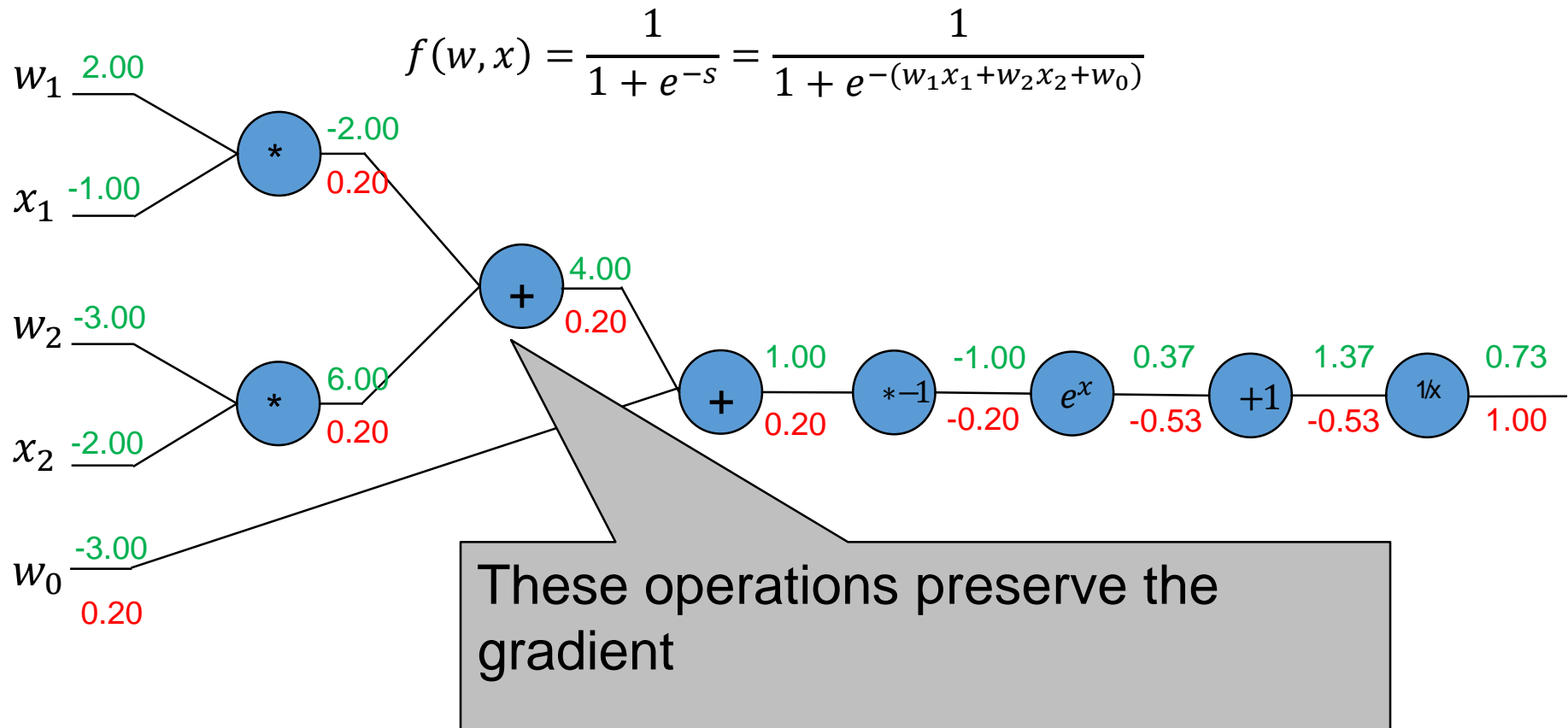
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



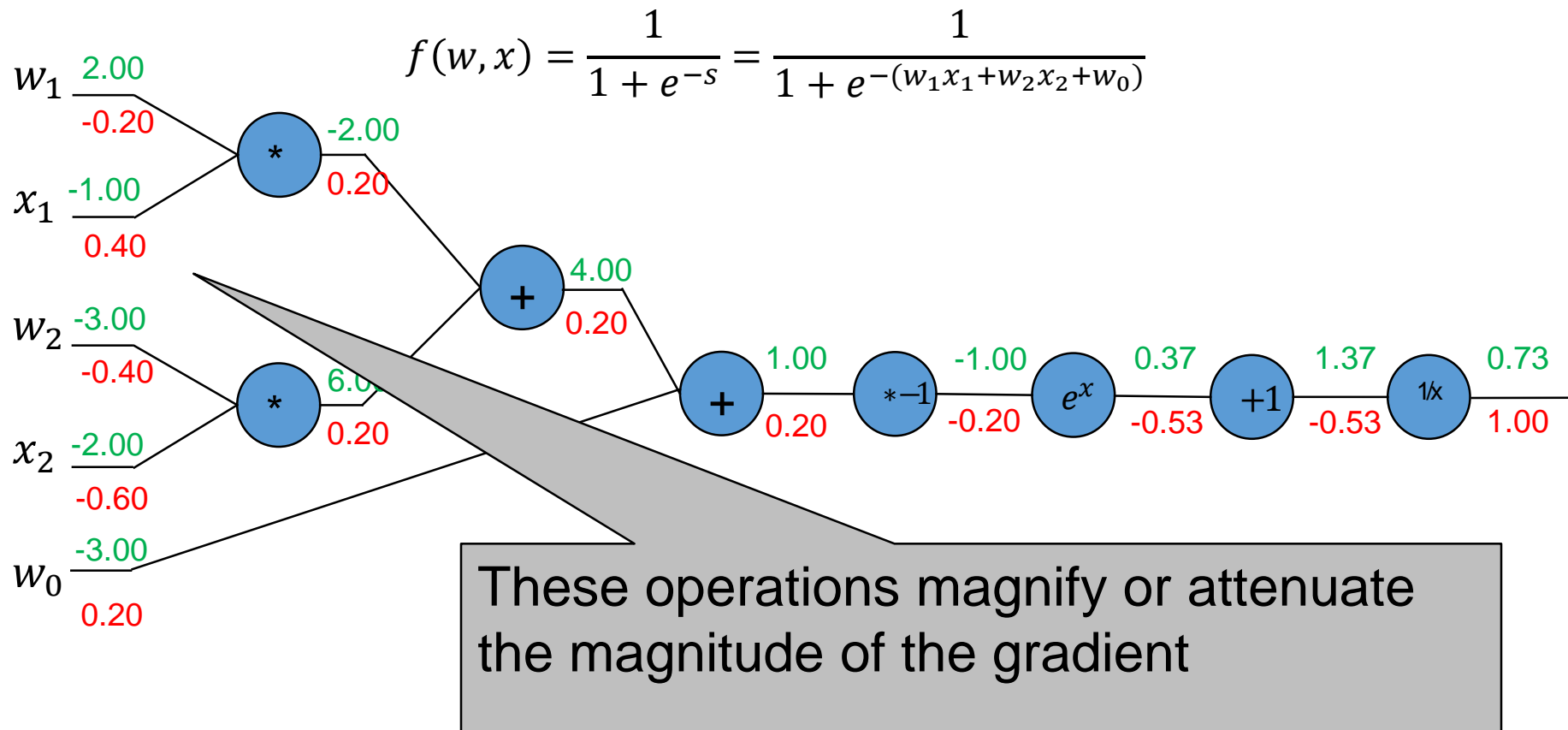
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



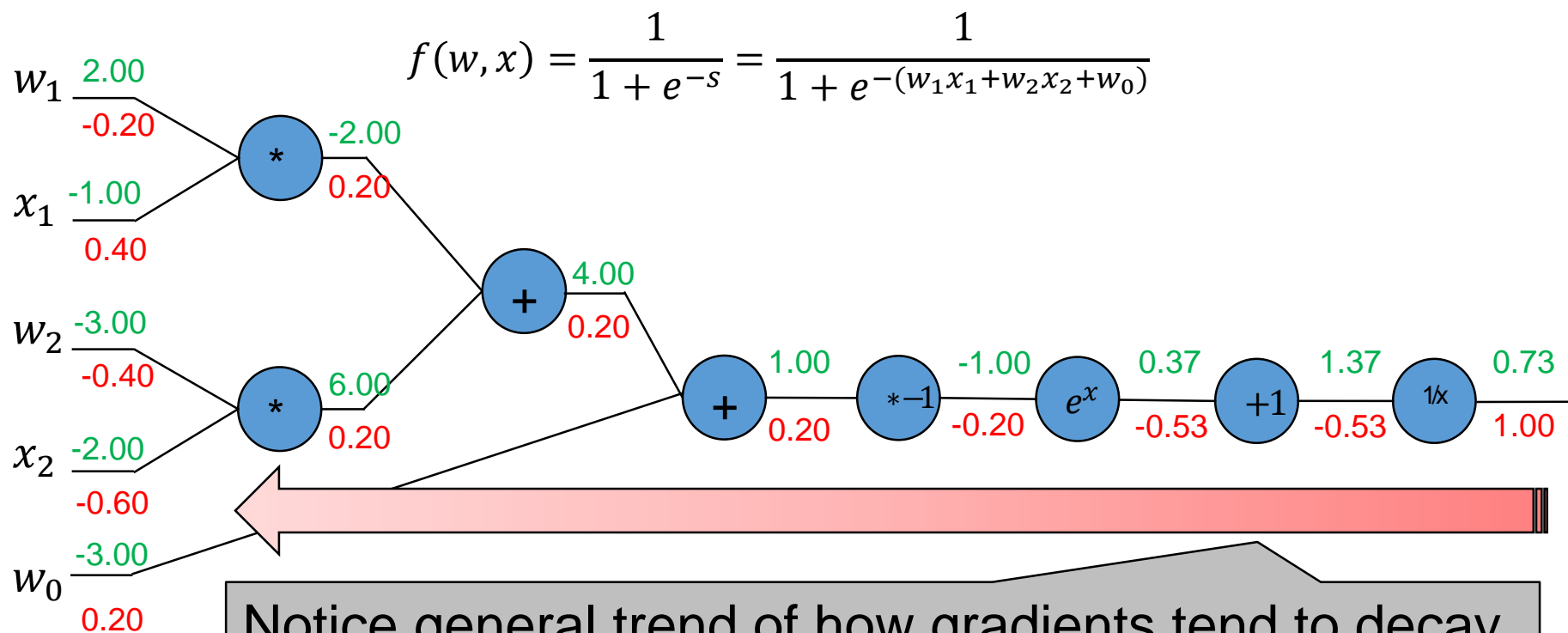
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Backward pass



Notice general trend of how gradients tend to decay in magnitude through layers of backprop... why?

$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

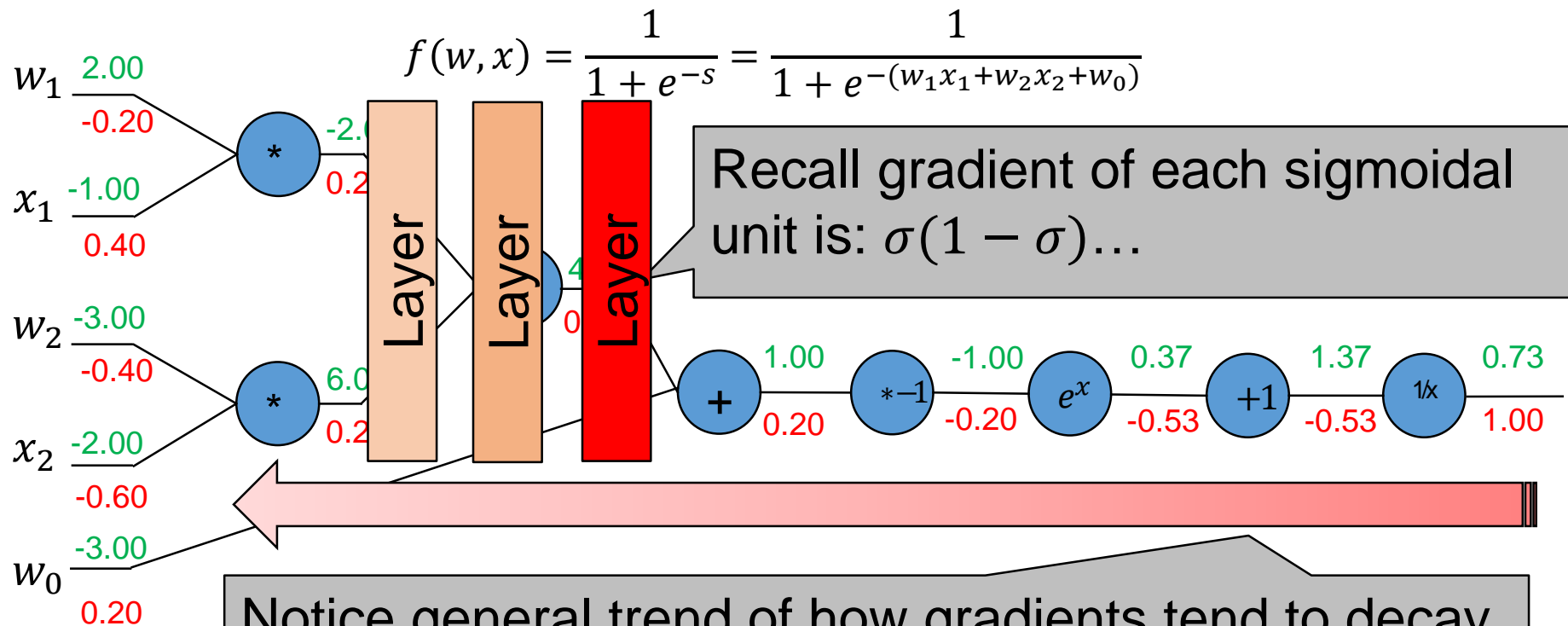
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$





## Backprop in “depth”: Backward pass



$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

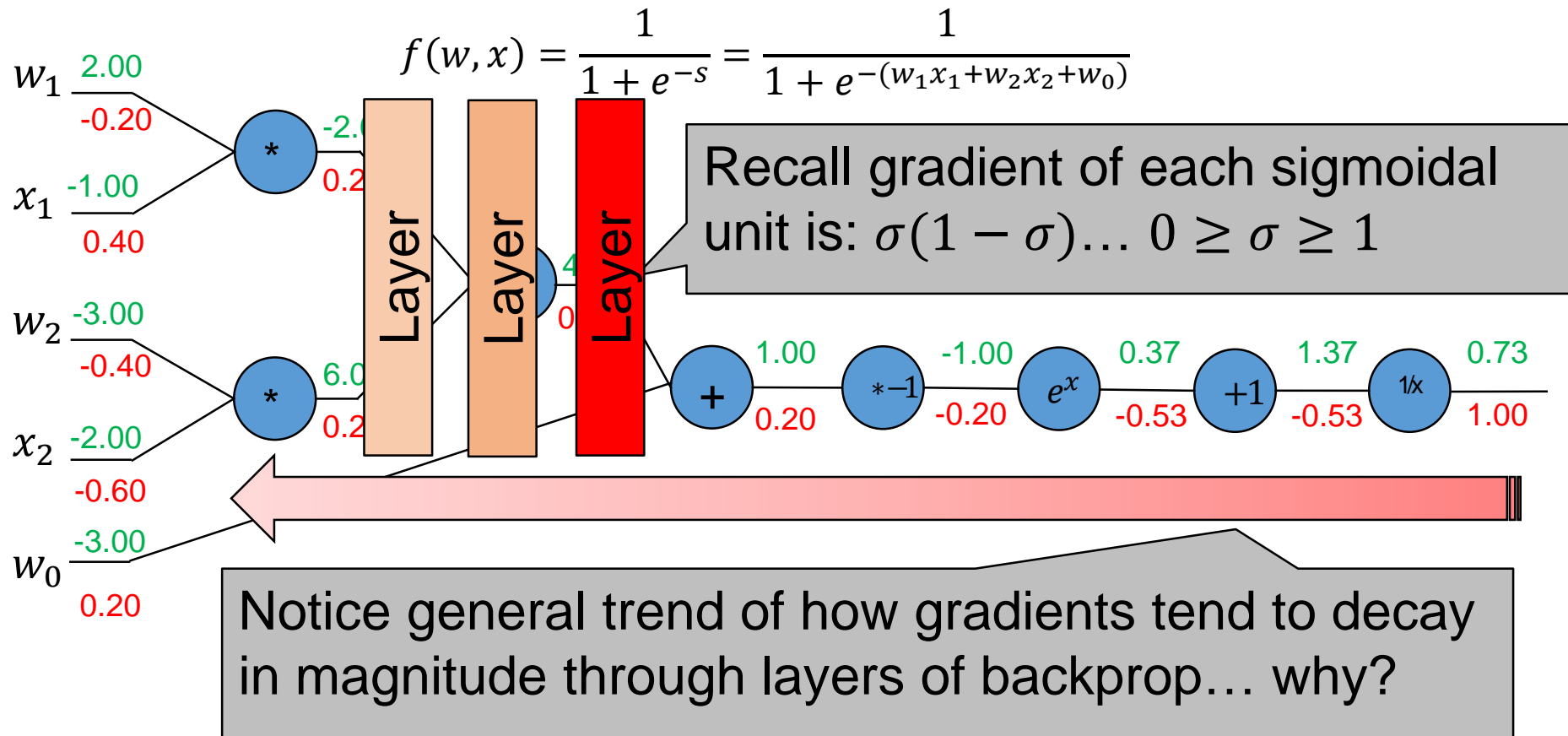
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$



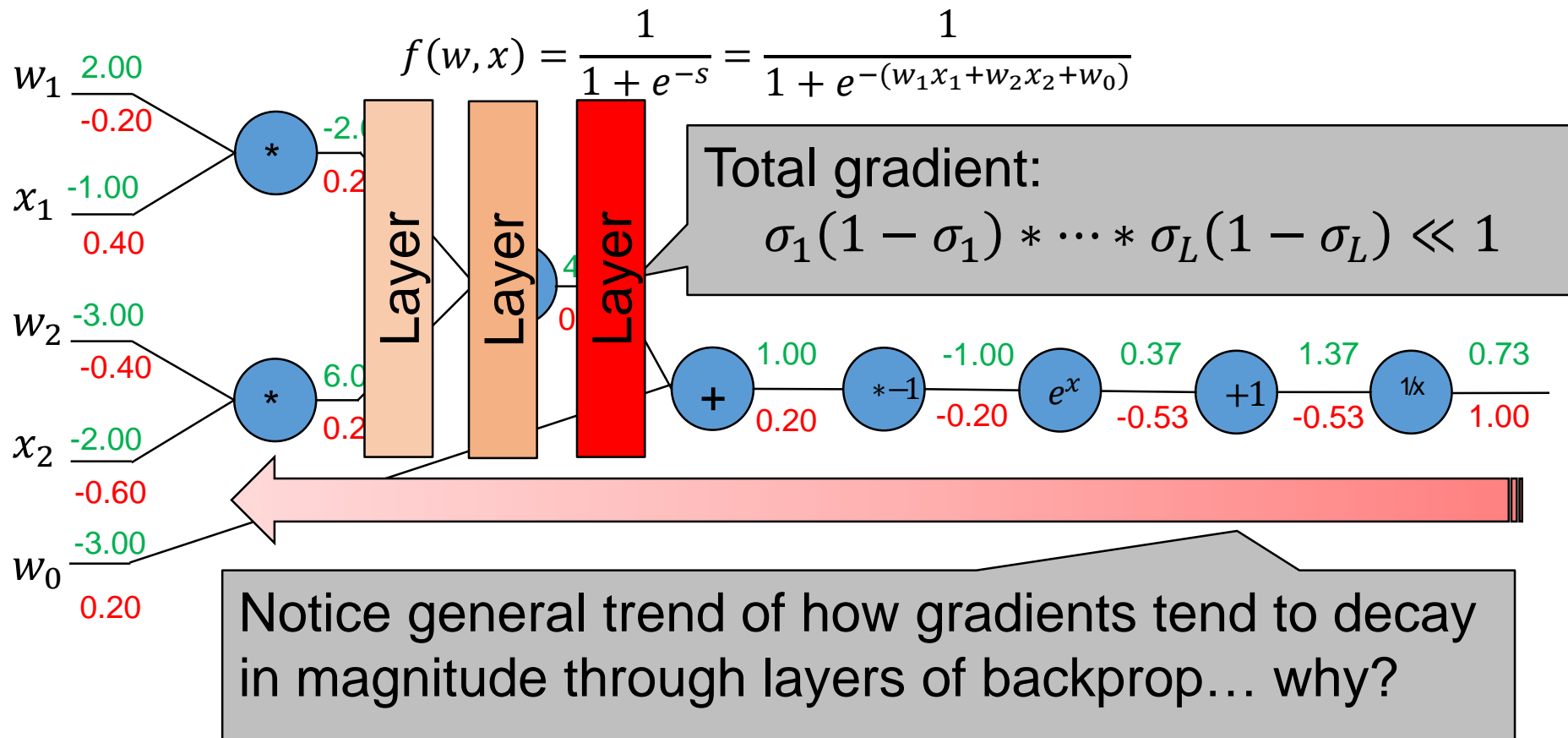
## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Backward pass



$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$	$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = \frac{-1}{x^2}$
$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$	$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$



## Backprop in “depth”: Backward pass

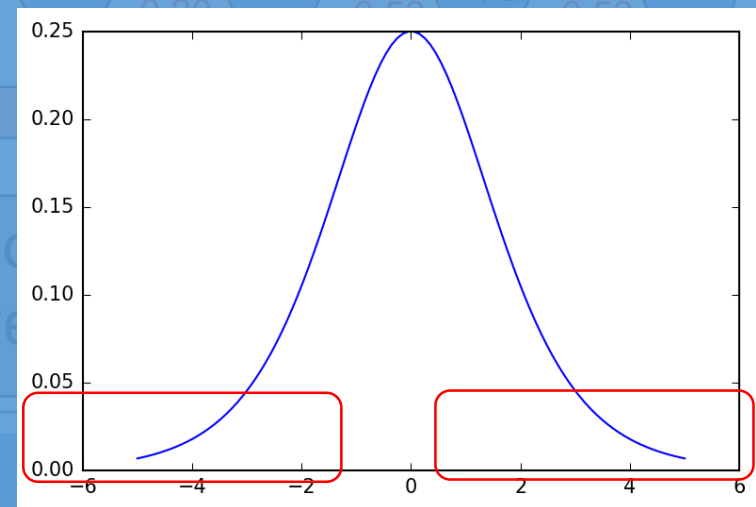
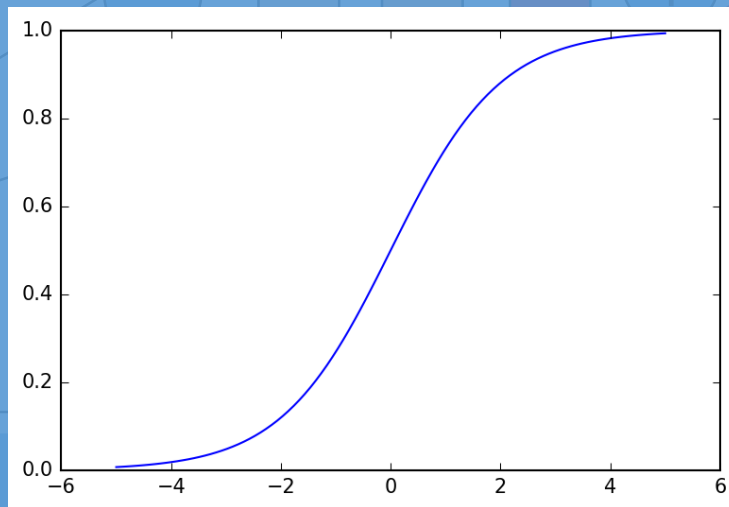
$$f(w, x) = \frac{1}{1 + e^{-s}} = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + w_0)}}$$

The overall effect is that gradients decay with each layer of sigmoidal activations.

→ Training slows down a lot at extrema

$$\sigma(s)$$

$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$



$$f_a(x) = ax, \quad \frac{\partial f}{\partial x} = a$$

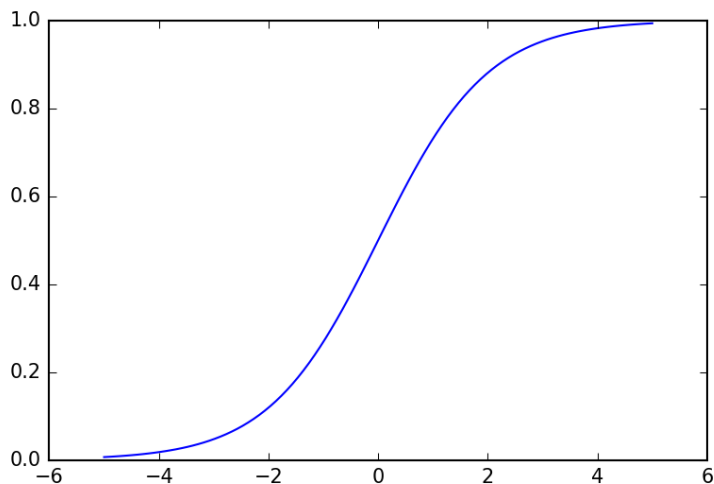
$$f_c(x) = c + x, \quad \frac{\partial f}{\partial x} = 1$$



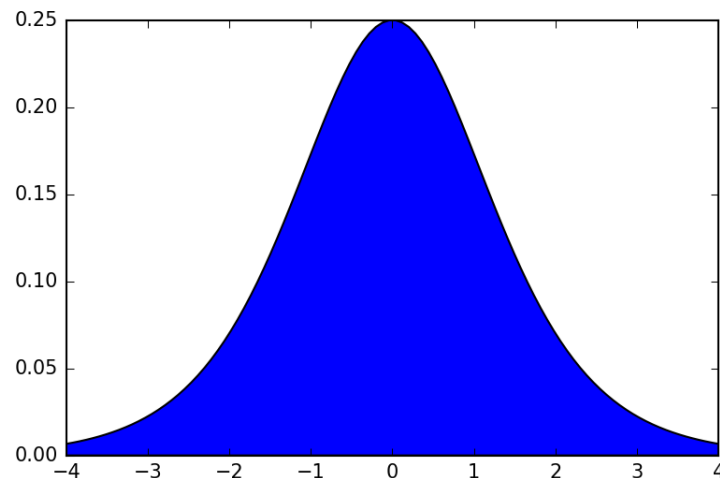
## Active regions of Sigmoids

- Sigmoids are relatively effective in active region
  - > 95% of gradient density in  $\sigma'(s)$  is in  $[-4,4]$
- Idea: keep sigmoid activations in this active region

$$\sigma(s)$$



$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$





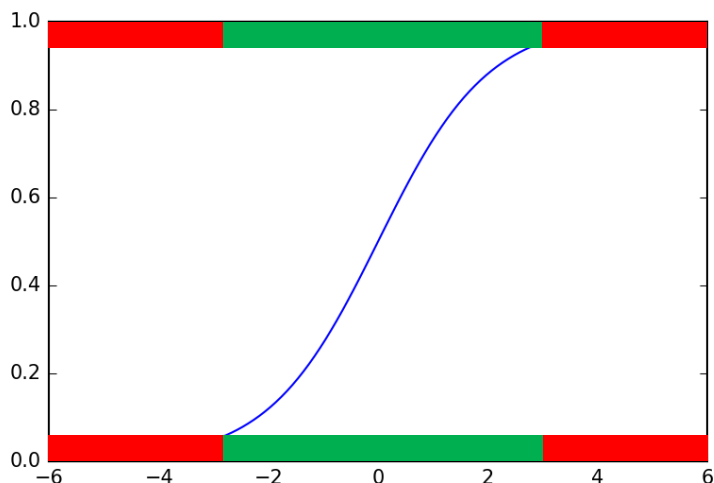
## Active regions of Sigmoids

→ Idea: keep sigmoid activations in this active region

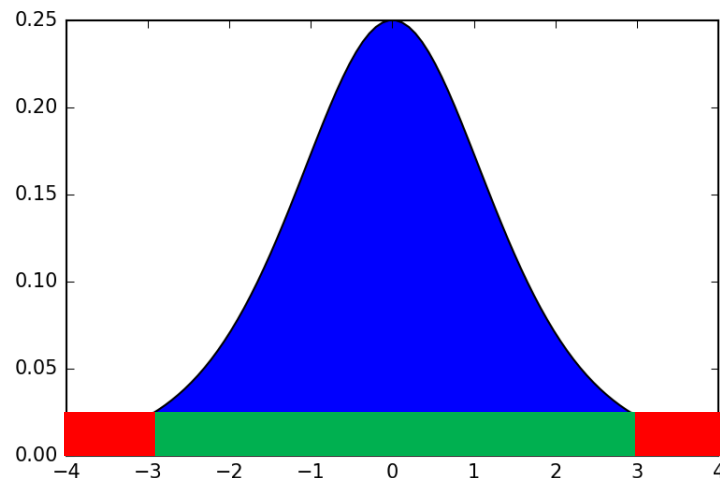
$$f(w, x) = \frac{1}{1 + e^{-s}} = \frac{1}{1 + e^{-(w_1x_1 + w_2x_2 + w_0)}}$$

This means that:  $-3 \leq w_1x_1 + w_2x_2 + w_0 \leq 3$

$\sigma(s)$

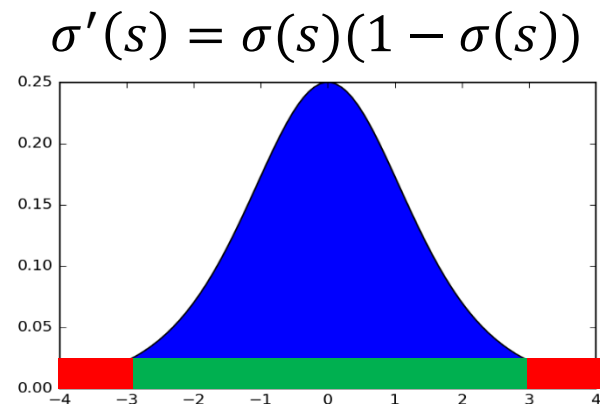
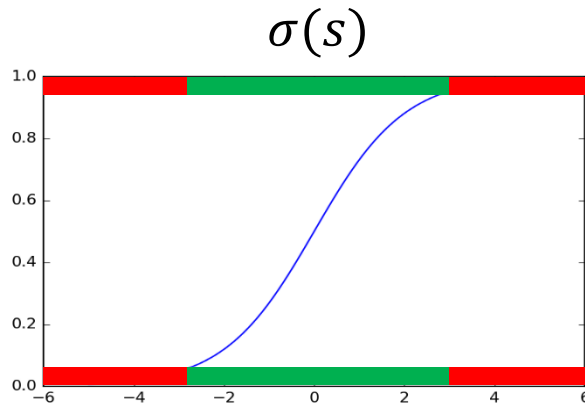
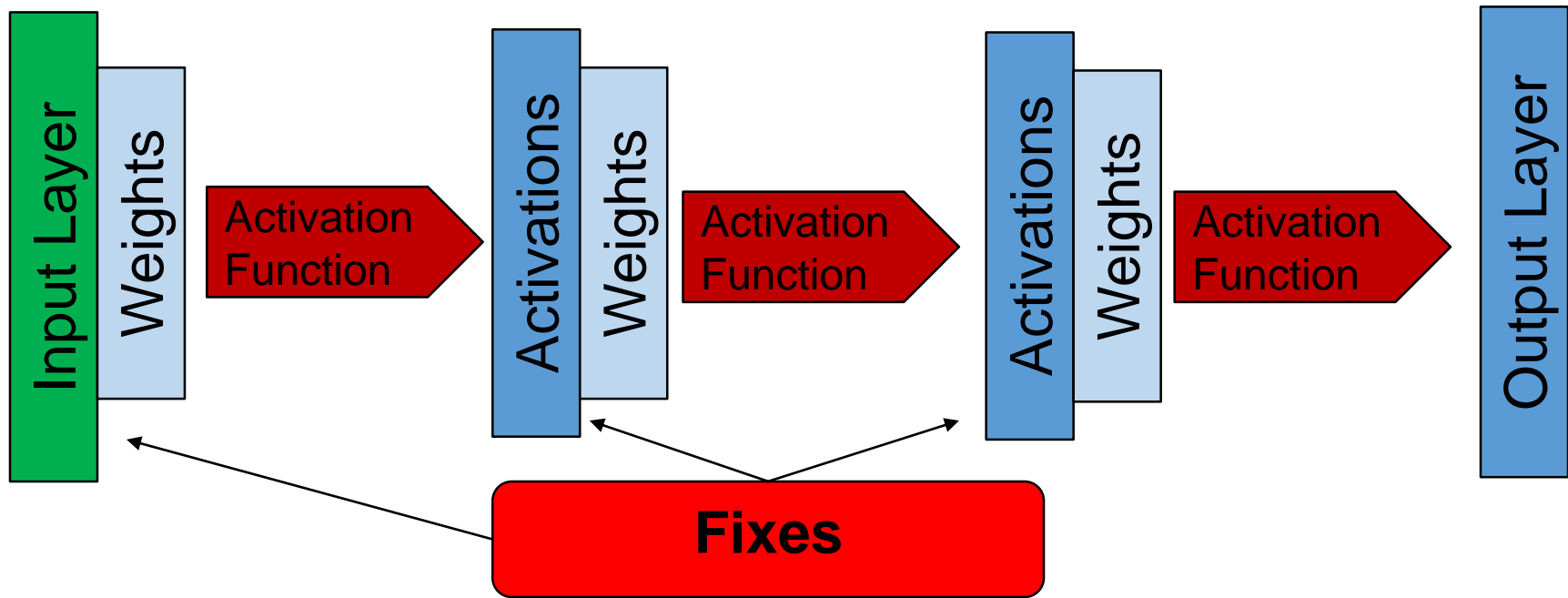


$\sigma'(s) = \sigma(s)(1 - \sigma(s))$





## Keeping Sigmoids in active region





## Simple fix: Clip gradients

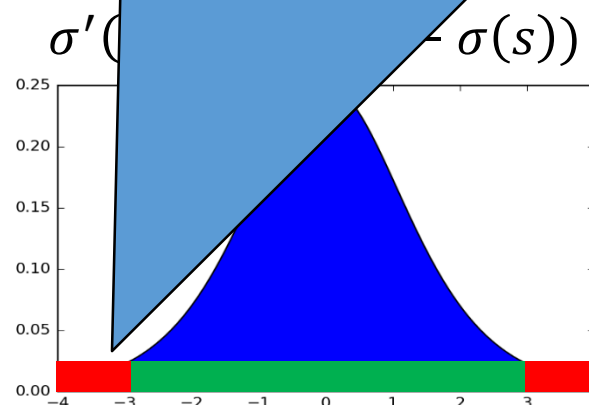
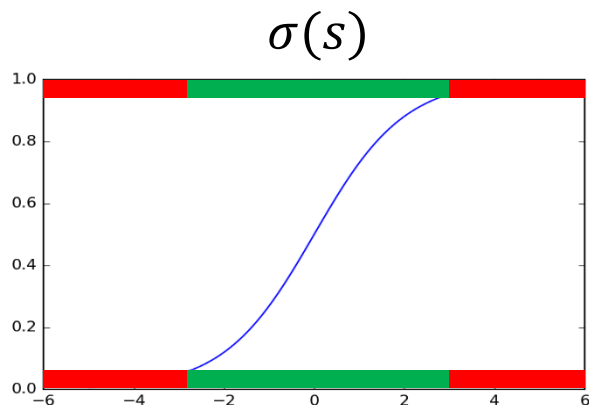
Very simple fix:

If  $\sigma'(s) < 0.01 \rightarrow \sigma'(s) = 0.01$

We can set this limit arbitrarily

→ This ensures that there will always be a non-saturated gradient

**Class:** Are there any cons to this approach?



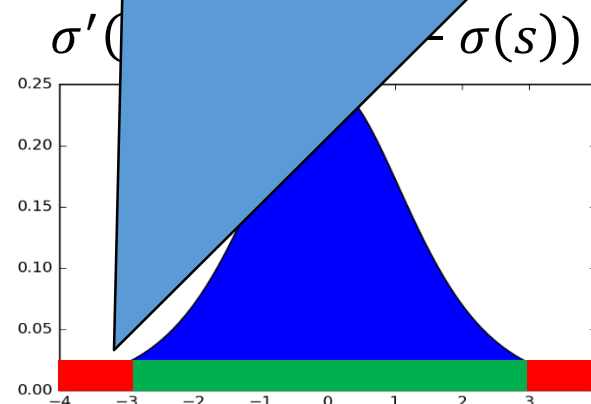
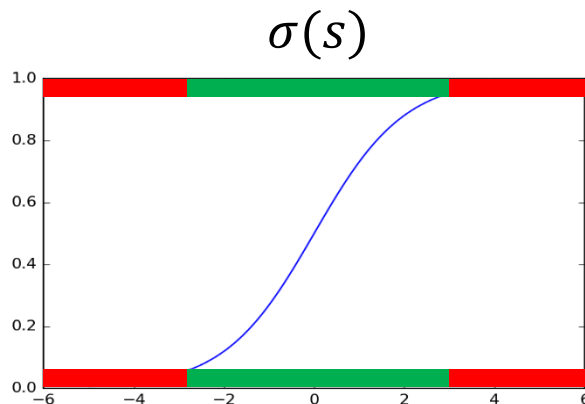




## Simple fix: Clip gradients

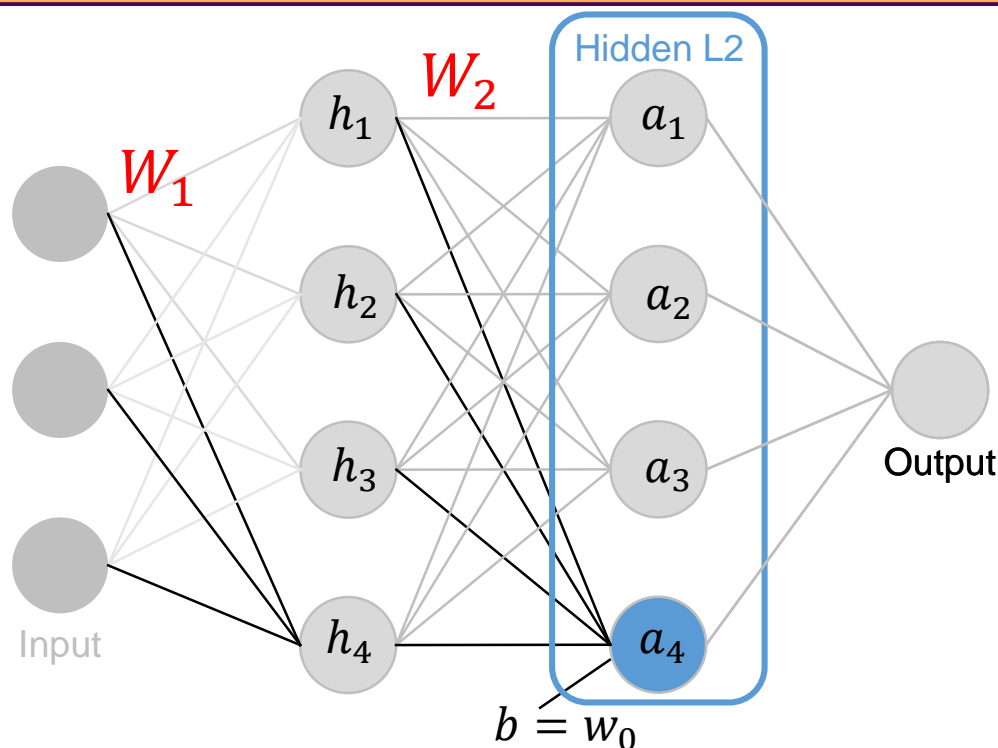
**Class:** Are there any cons to this approach?

- If the threshold is set too low, we recover the original sigmoid with the vanishing gradients
- If the threshold is set too high, we get only the linear region and lose the benefit of a nonlinearity,
  - E.g. The matrices can be multiplied out into a single linear layer, thus losing effective depth





## Weight constraints



$$\sigma \left( \begin{matrix} W_1 & h_1 \\ \dots & h_2 \\ \dots & h_3 \\ W_n & h_4 \end{matrix} + w_0 \right) = \begin{matrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{matrix}$$

Recall:

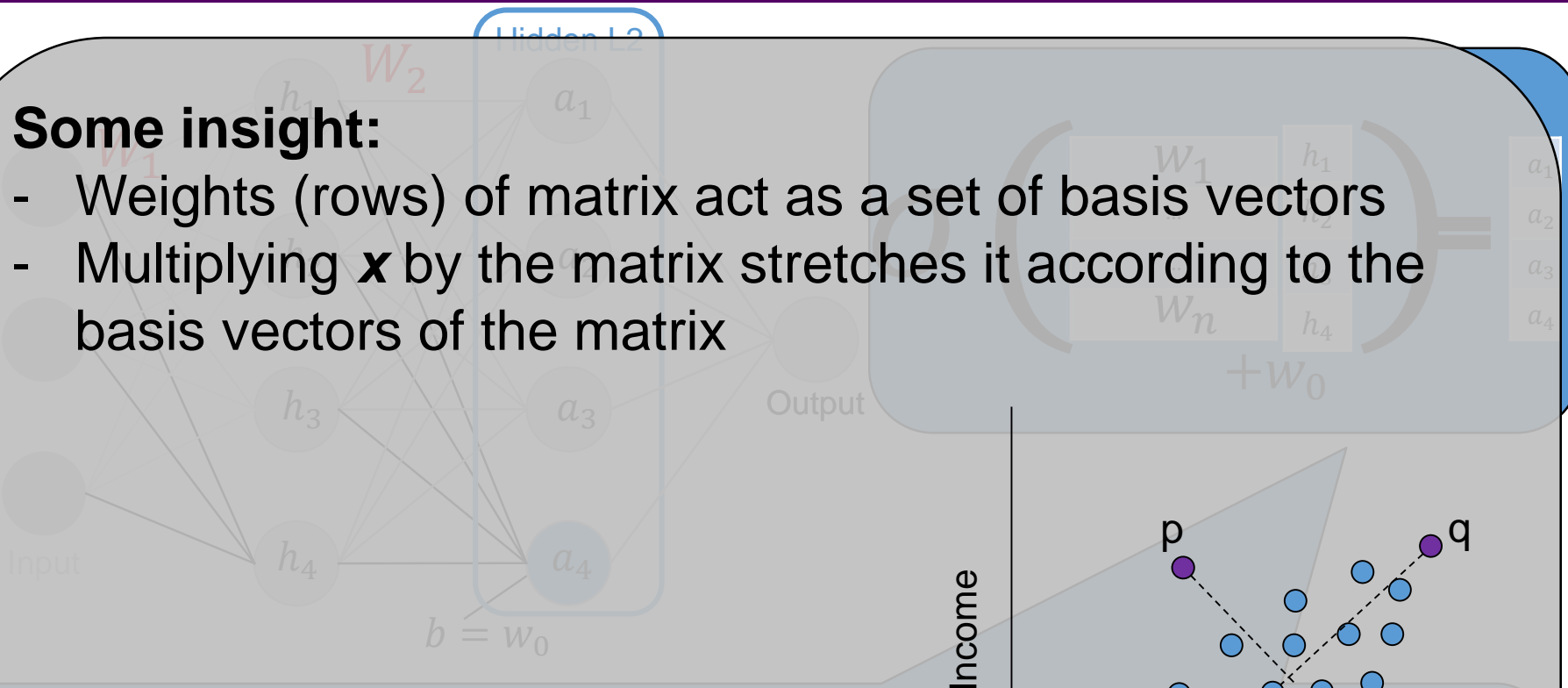
- Weights are row vectors
- Each activation is a dot product of weight and input  $x$
- We want  $-3 \leq w \cdot x \leq 3$  so that sigmoid is in active region



## Bases of a weight matrix

### Some insight:

- Weights (rows) of matrix act as a set of basis vectors
- Multiplying  $\mathbf{x}$  by the matrix stretches it according to the basis vectors of the matrix



### Recall:

- Weights are row vectors
- Each activation is a dot product of weight and input  $\mathbf{x}$
- We want  $-3 \leq \mathbf{w} \cdot \mathbf{x} \leq 3$  so that sigmoid is in active region



## Eigenvectors of $W$

### Some insight:

- Weights (rows) of matrix act as a set of basis vectors
- Multiplying  $x$  by the matrix stretches it according to the basis vectors of the matrix

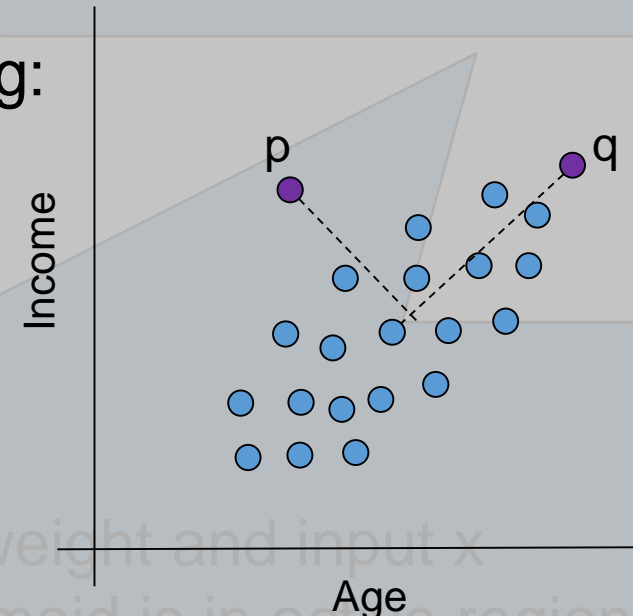
- We can find these bases by solving:

$$Ax = \lambda x$$

$$b = w_0$$

### Recall:

- Weights are row vectors
- Each activation is a dot product of weight and input  $x$
- We want  $-3 \leq w \cdot x \leq 3$  so that sigmoid is in active region





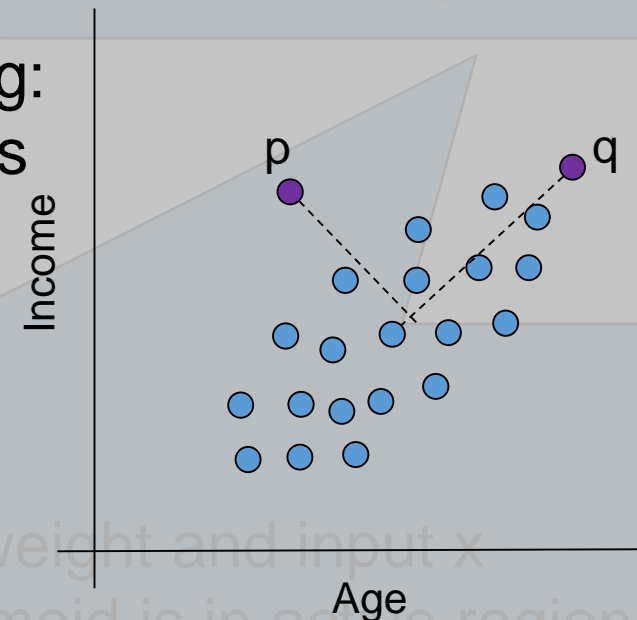
## Eigenvectors of $W$

### Some insight:

- Weights (rows) of matrix act as a set of basis vectors
- Multiplying  $x$  by the matrix stretches it according to the basis vectors of the matrix

- We can find these bases by solving:  
 $Ax = \lambda x \rightarrow$  Eigenvalues, eigenvectors

$v_1$	$\lambda_1$			
...		$\lambda_2$		
...			...	
$v_n$				$\lambda_n$





## Eigenvectors of $W$

### Some insight:

- Weights (rows) of matrix act as a set of basis vectors
- Multiplying  $x$  by the matrix stretches it according to the basis vectors of the matrix

- We can find these bases by solving:  
 $Ax = \lambda x \rightarrow$  Eigenvalues, eigenvectors

$v_1$	$\lambda_1$			
...		$\lambda_2$		
...			...	
$v_n$				$\lambda_n$

Each  $\lambda_i$  acts as a stretch factor. If  $\lambda_i < 1$  or  $\lambda_i > 1$ , then resulting vector shrinks or grows.



## Eigenvectors of $W$

### Some insight:

- Weights (rows) of matrix act as a set of basis vectors
- Multiplying  $x$  by the matrix stretches it according to the basis vectors of the matrix

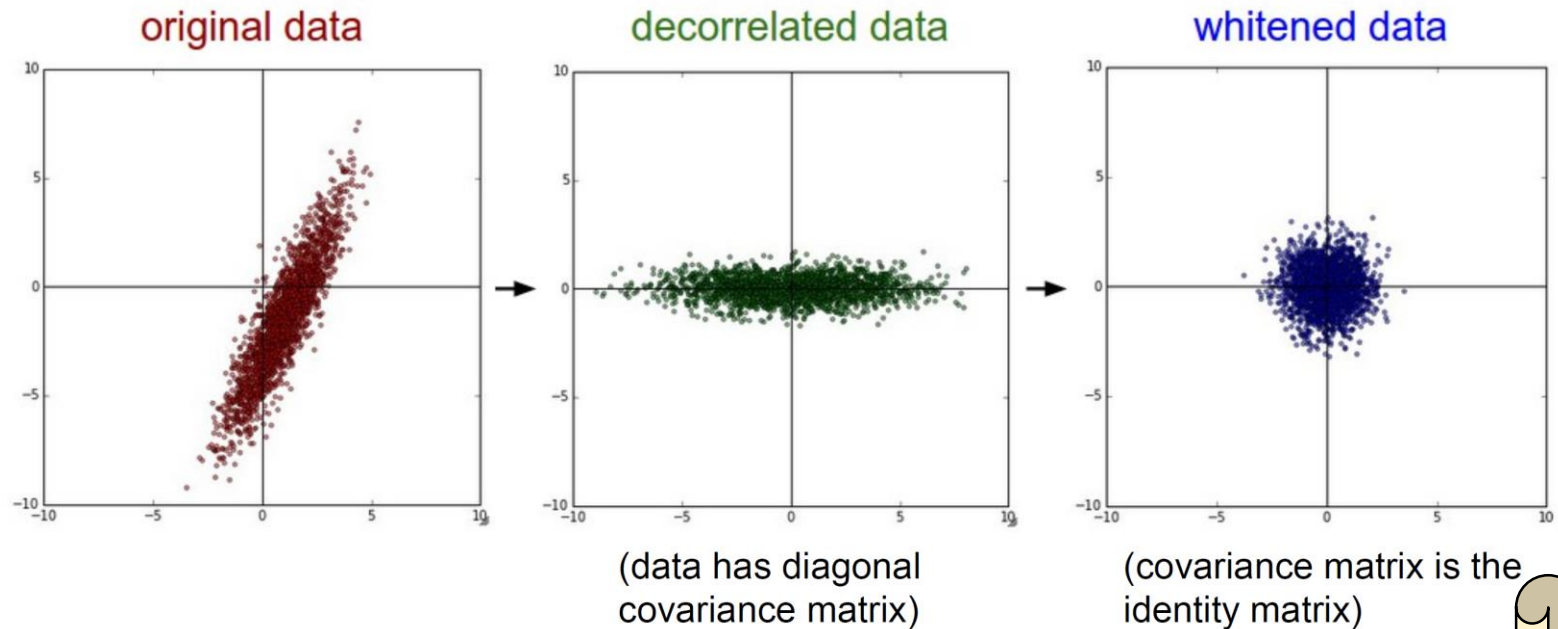
- We can find these bases by solving:  
 $Ax = \lambda x \rightarrow$  Eigenvalues, eigenvectors

$v_1$	$\lambda_1$			
...		$\lambda_2$		
...			...	
$v_n$				$\lambda_n$

Each  $\lambda_i$  acts as a stretch factor. If  $\lambda_i < 1$  or  $\lambda_i > 1$ , then resulting vector shrinks or grows. We want  $\lambda_i \approx 1$ !



## Keeping sigmoids in active region



### Practical tip:

- Subtract data mean(bias=0), Divide by standard deviation
- Whiten data (PCA/ZCA) if possible

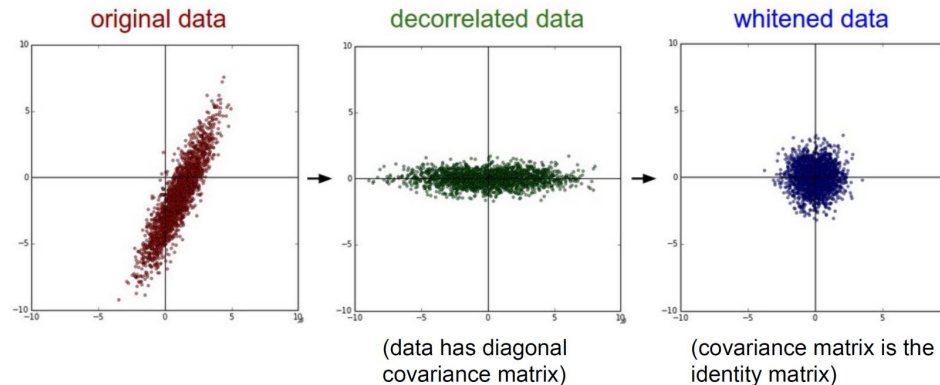
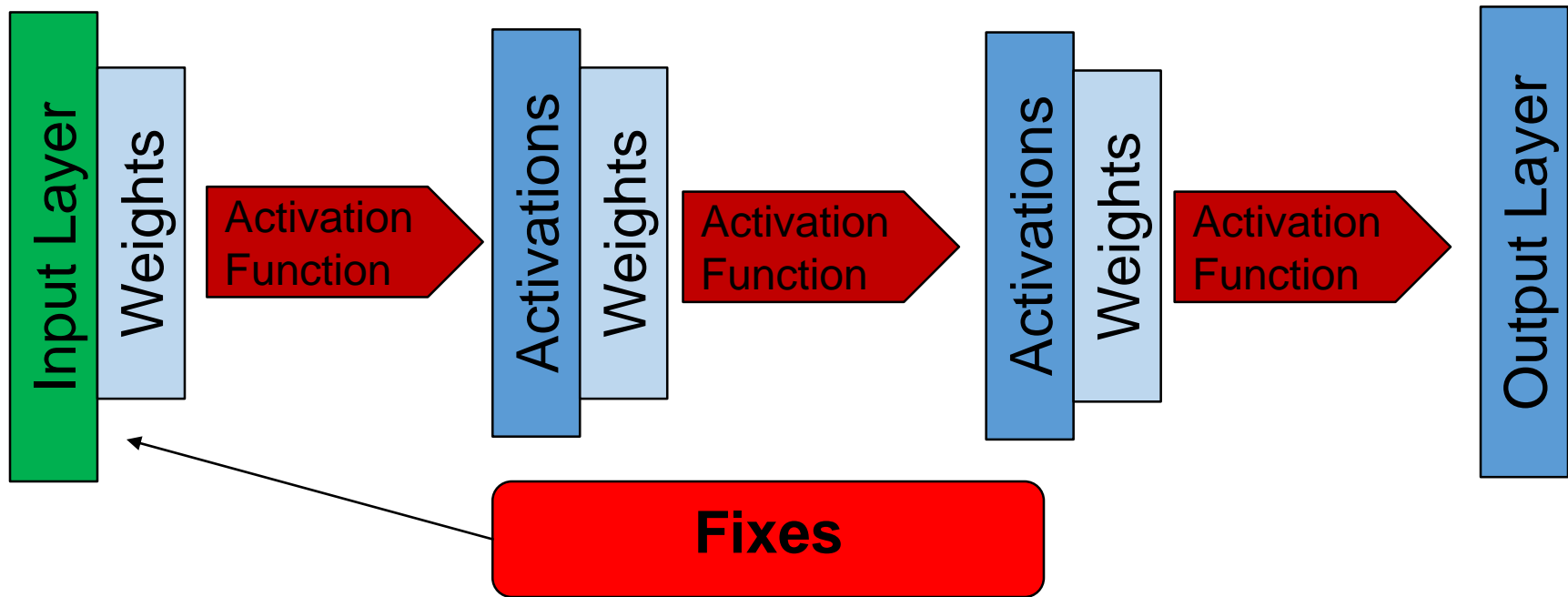
See: [http://ufldl.stanford.edu/wiki/index.php/Exercise:PCA\\_and\\_Whitening](http://ufldl.stanford.edu/wiki/index.php/Exercise:PCA_and_Whitening)

- Keeps activations at each layer in an active region



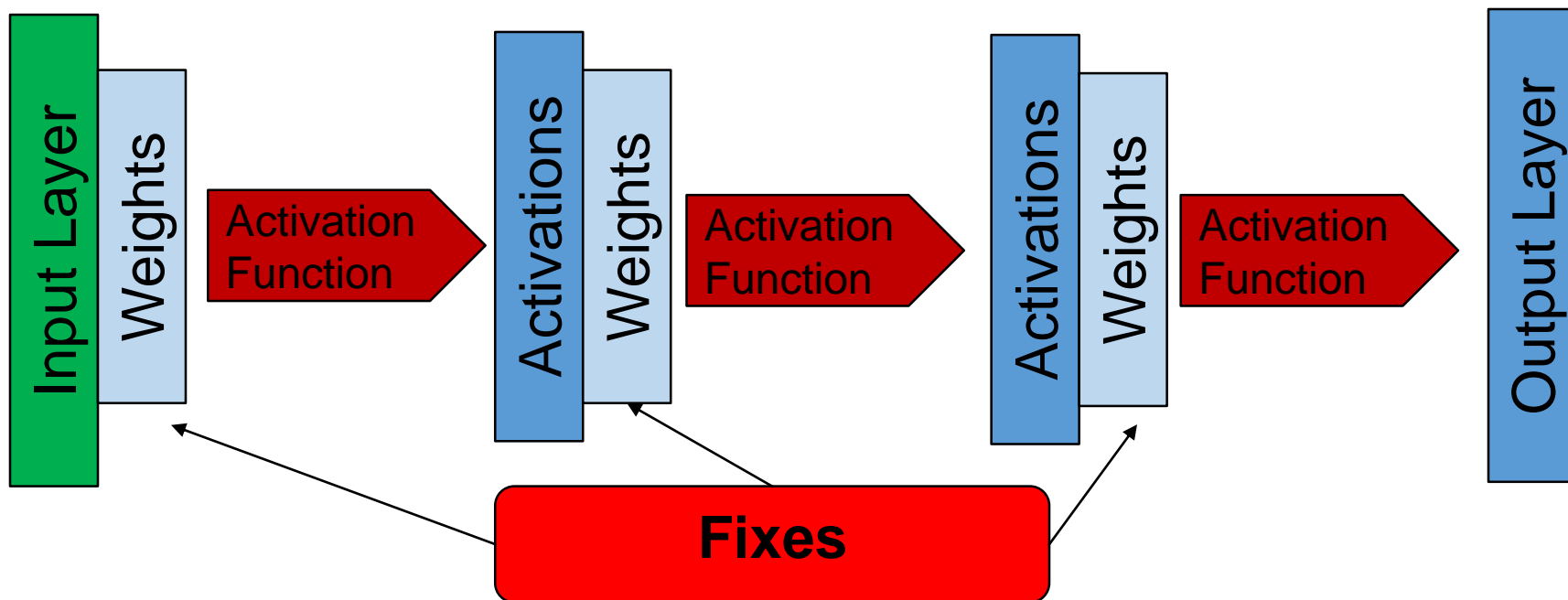


## Keeping sigmoids in active region





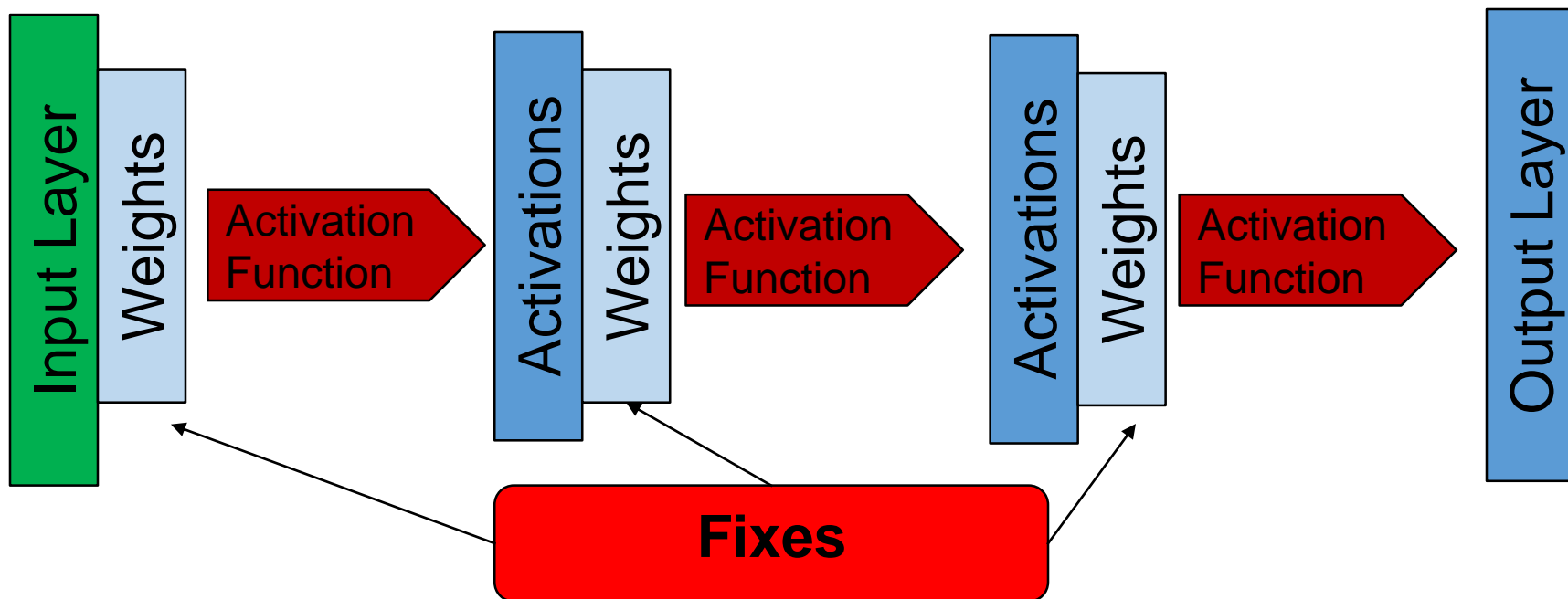
## Keeping sigmoids in active region



- What about weights?
- Weights are usually initialized randomly
- Poor initialization can lead to small/large eigenvalues and shrinkage/explosion of gradients
- **Practical fix:** Initialize weights with  $\mu = 0, \sigma \ll 1$



## Keeping sigmoids in active region

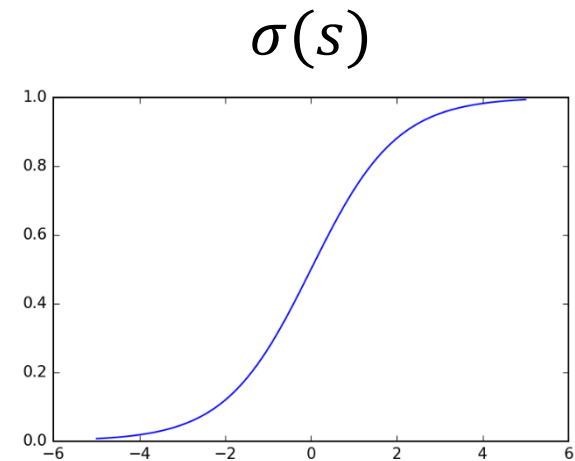


- $f(w, x) = \sum_i w_i x_i$ . Scale of output is determined by fan-in. Suppose  $w_i x_i$  is Gaussian distributed with  $\mu_i, \sigma_i$ .  $\sigma^2 = \sum_{i=1}^n \sigma_i^2 = n\sigma^2$ . We overshoot by  $\sqrt{n}$ .
- **Theory:** Normalize inputs by fan-in  $\sigma_i/\sqrt{n}$  to maintain unit variance in output.



## Nature of the logistic activation

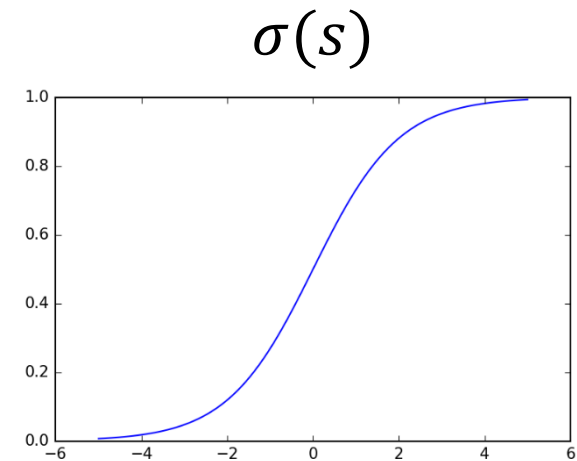
- $0 \leq \sigma(\sum_i w_i x_i) \leq 1$





## Nature of the logistic activation

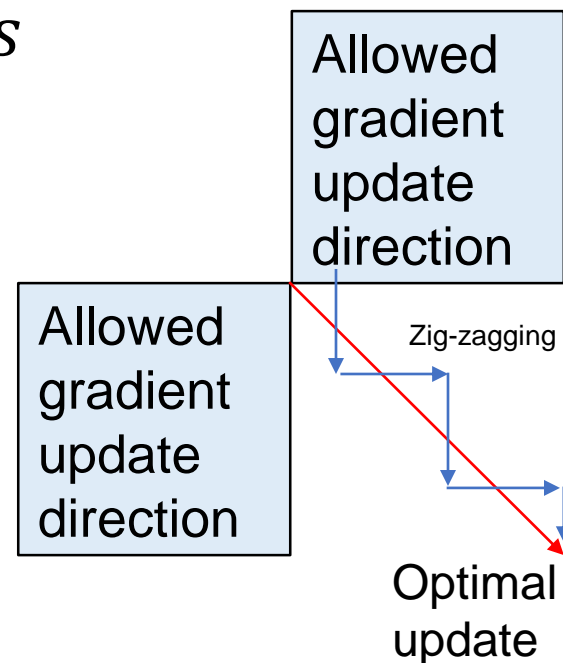
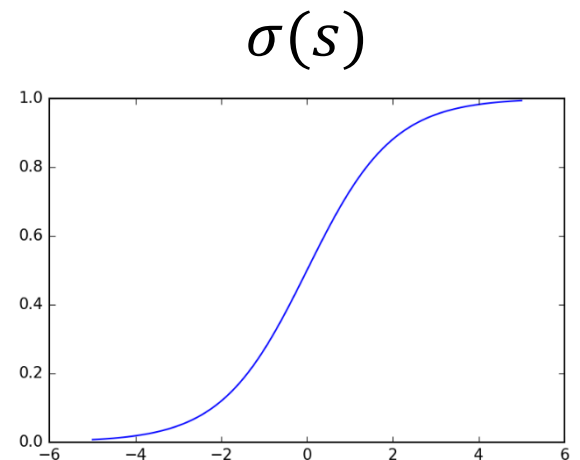
- $0 \leq \sigma(\sum_i w_i x_i) \leq 1$
- Inputs at next layer are always positive
- How does that constrain the gradients?  $\sigma'(s) = \sigma(s)(1 - \sigma(s))ds$





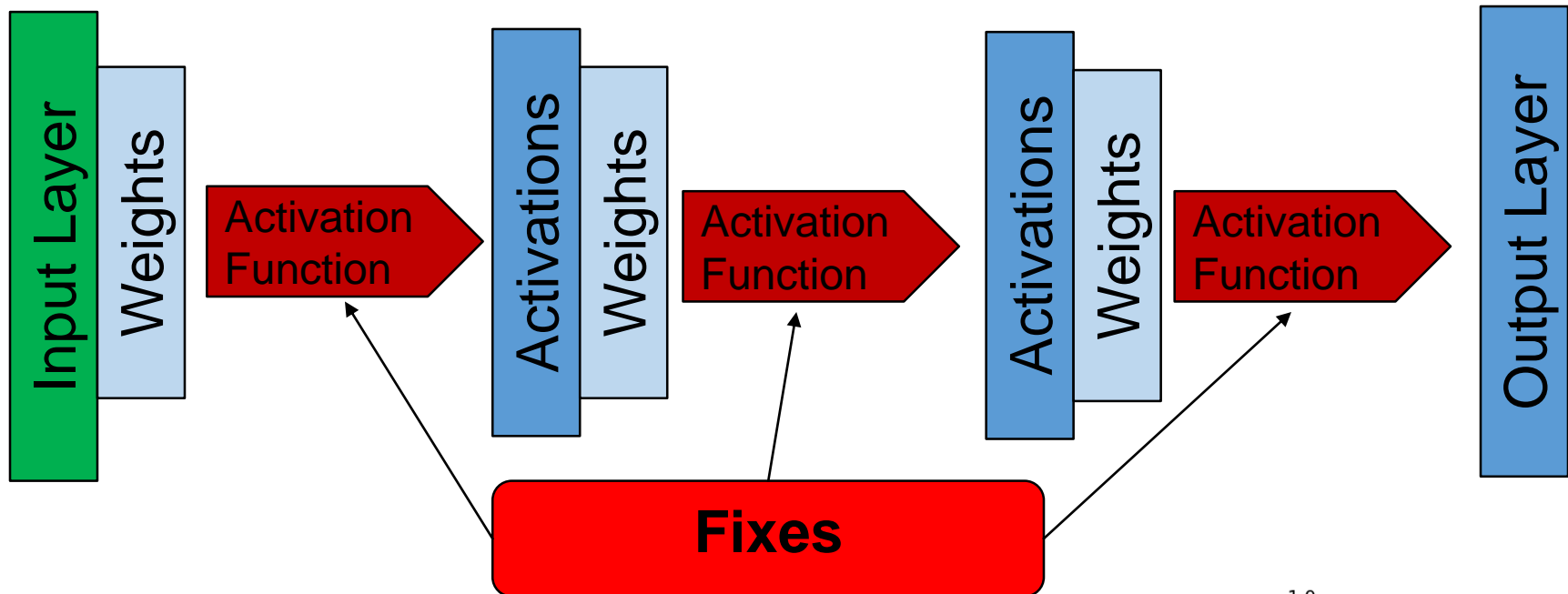
## Nature of the logistic activation

- $0 \leq \sigma(\sum_i w_i x_i) \leq 1$
- Inputs at next layer are always positive
- How does that constrain the gradients?  $\sigma'(s) = \sigma(s)(1 - \sigma(s))ds$
- Depending on the sign of  $ds$ , gradients are always all positive or all negative





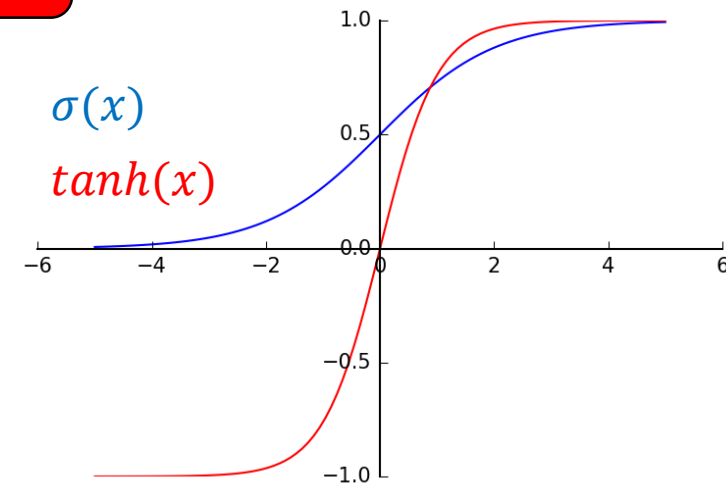
## Tanh(x)



**Advice:** Choose a mean-centered  
Activation function

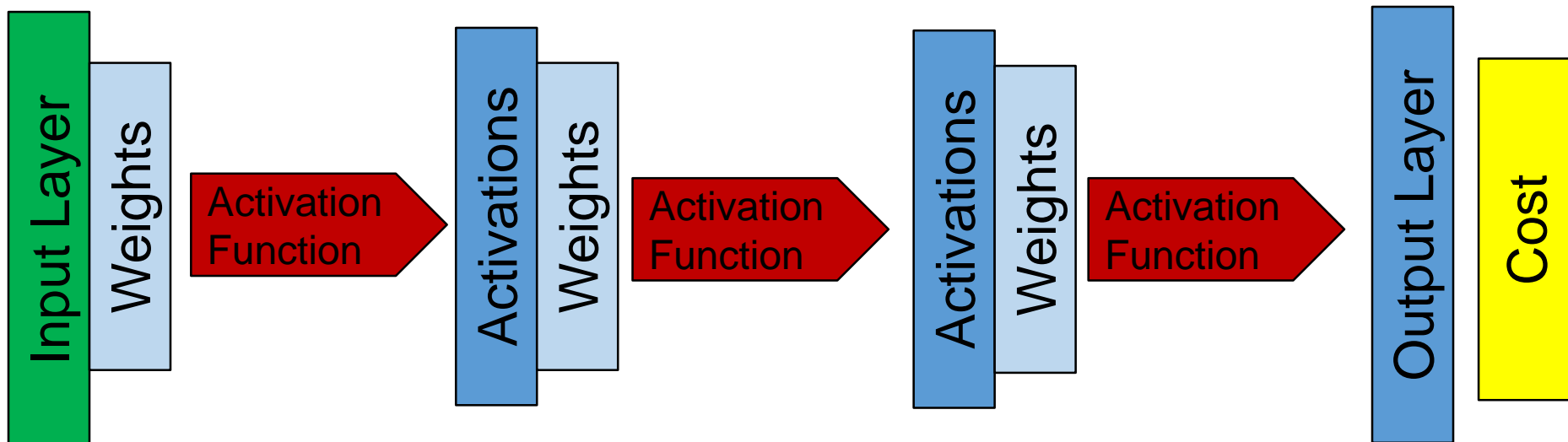
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(2x) - 1$$

- Addresses data centering issue (if input is zero mean)
- Addresses zigzag issue (activation can be negative)
- Does not address vanishing gradients





## Cost function

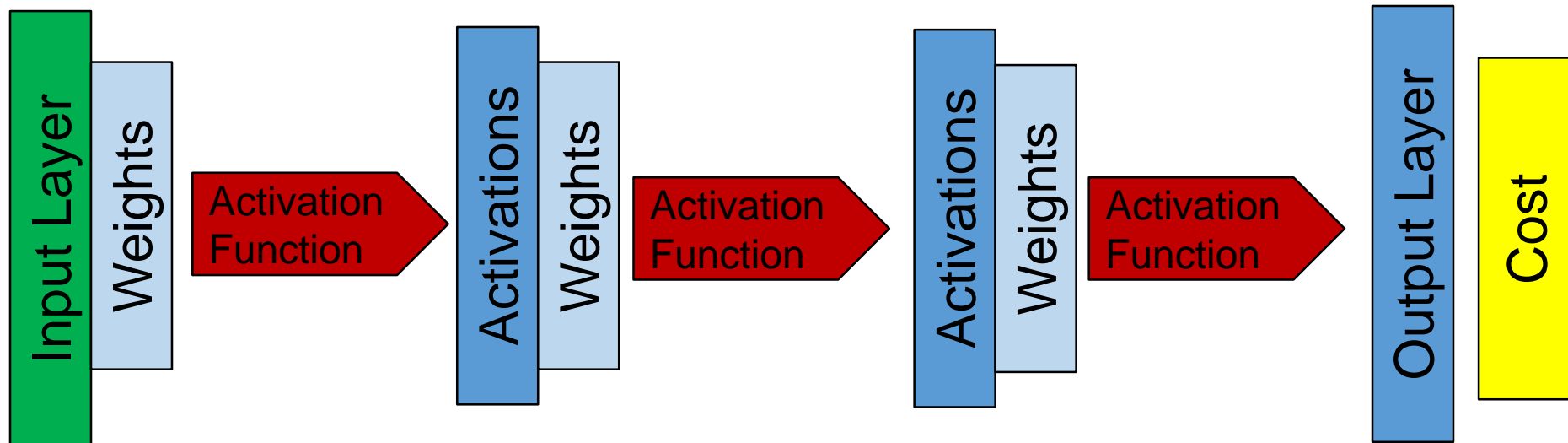


- What about the cost function?
- We've been using  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - t_i)^2$
- Derivative of MSE:  $2(y_i - t_i) dy$

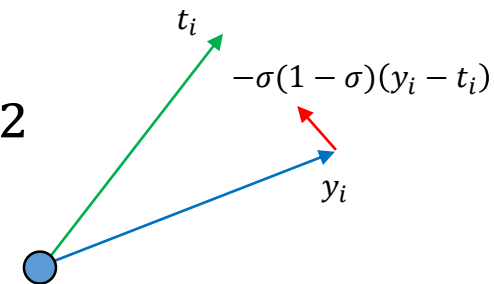




## Cost function



- What about the cost function?
- We've been using  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - t_i)^2$
- Derivative of MSE:  $2(y_i - t_i) dy$

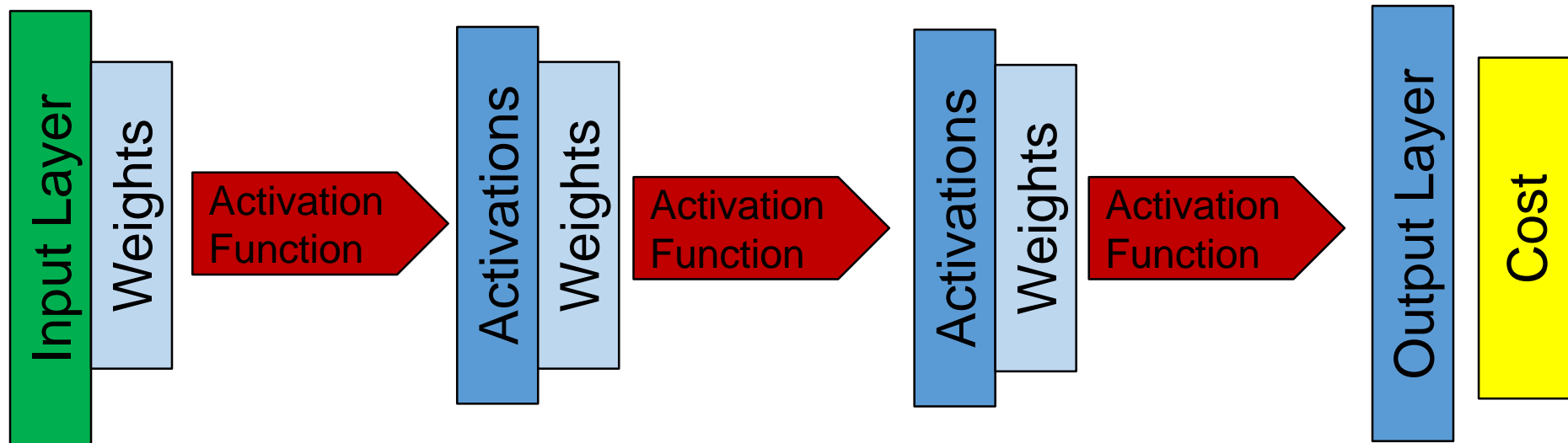


Measure of  
disagreement between  
prediction  $y$  and target  $t$

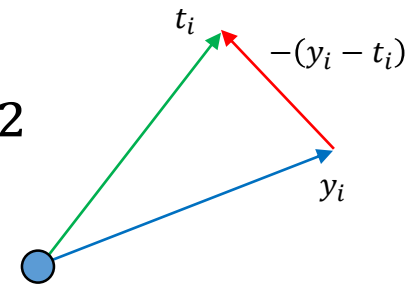
$dy = \sigma(1 - \sigma)$   
Derivative of activation function.  
Remember that this is **SLOW**.



## Cost function



- What about the cost function?
- We've been using  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - t_i)^2$
- Derivative of MSE:  $2(y_i - t_i)$

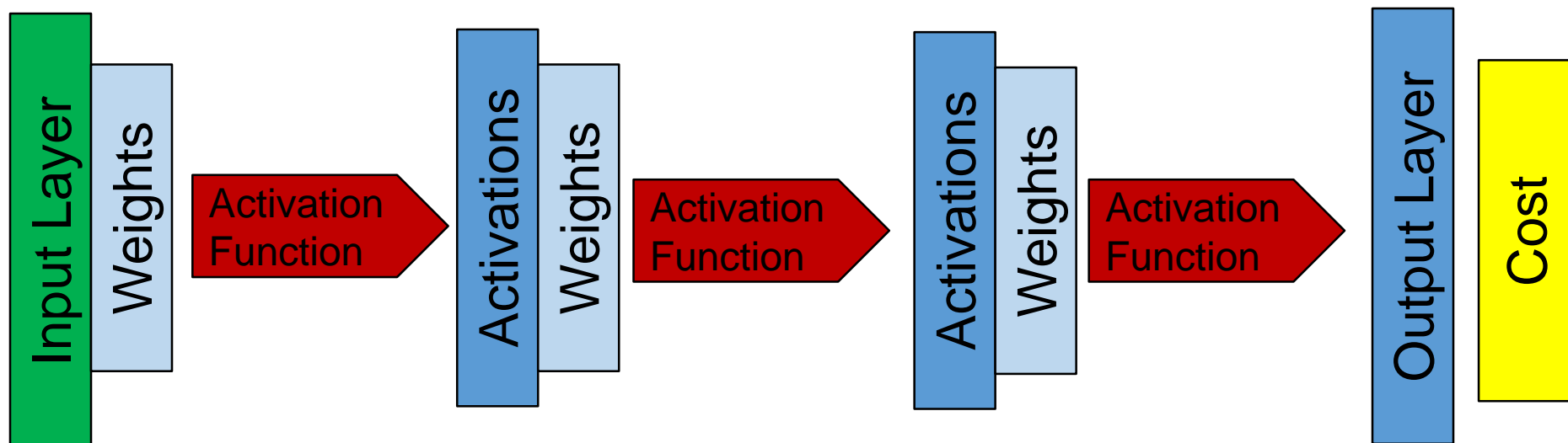


Measure of disagreement between prediction  $y$  and target  $t$ .  
Let's KEEP this.

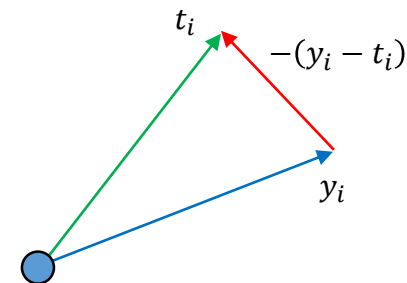
And find a way to eliminate or at least speed up this...



## Keeping sigmoids in active region

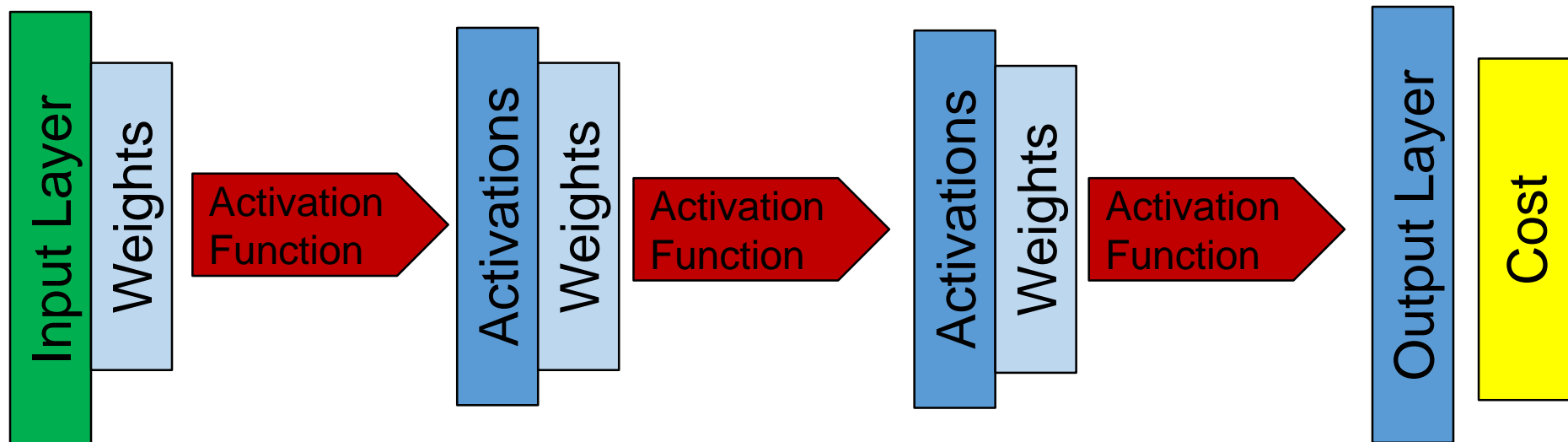


- Desired gradient:  $\partial E = y - t$





## Keeping sigmoids in active region



- Desired gradient:  $\partial E = y - t$
- Cross entropy cost:  $E = -t \ln y - (1 - t) \ln(1 - y)$

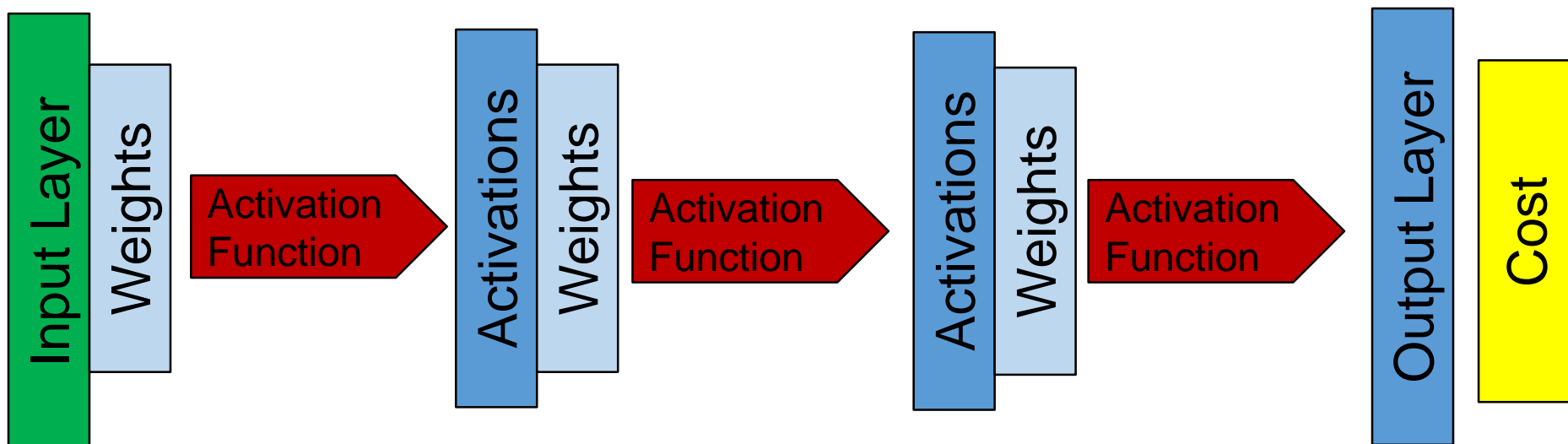


## Cross entropy in depth

- $E = -\frac{1}{n} \sum_x t \ln \sigma(z) + (1 - t) \ln(1 - \sigma(z))$
- $\frac{\partial E}{\partial w_j} = -\frac{1}{n} \sum_x \left( \frac{t}{\sigma(z)} - \frac{1-t}{1-\sigma(z)} \right) \sigma'(z) x_j$
- $= -\frac{1}{n} \sum_x \left( \frac{t(1-\sigma(z)) - (1-t)\sigma(z)}{\sigma(z)(1-\sigma(z))} \right) \sigma(z)(1 - \sigma(z)) x_j$
- $= -\frac{1}{n} \sum_x \left( \frac{t - t\sigma(z) - \sigma(z) + t\sigma(z)}{\sigma(z)(1-\sigma(z))} \right) \sigma(z)(1 - \sigma(z)) x_j$
- $= \frac{1}{n} \sum_x (\sigma(z) - t) x_j$
- **→ Slow  $\sigma'(z)$  cancels out!**
- See Nielsen Ch 3 for details: <http://neuralnetworksanddeeplearning.com/chap3.html>
- Also read Nielsen Ch 3 to learn where cross entropy comes from



## Keeping sigmoids in active region



- Cross entropy cost:  $E = -t \ln y - (1 - t) \ln(1 - y)$

- Cross entropy gradient:  $\partial E = y - t$

MSE gradient:  $\partial E = (y - t) dy$

$dy = \sigma(1 - \sigma)$   
Derivative of activation  
function

- Gradient is proportional to difference

- No learning slowdown due to sigmoid gradients



## But how did we get cross entropy?

- Using  $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = a(1 - a)$
- We'd like the  $a(1 - a)$  to cancel out of the cost, so
- $\frac{\partial E}{\partial a} = \frac{a - y}{a(1 - a)}$

$a - y$  is just the discrepancy between the prediction  $y$  and actual value  $t$ .  $a(1 - a)$  cancels out with the derivative.



## But how did we get cross entropy?

- Using  $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = a(1 - a)$
- We'd like the  $a(1 - a)$  to cancel out of the cost to get the gradient to be  $a - t$ , so
- $$\frac{\partial E}{\partial a} = \frac{a-t}{a(1-a)} = \frac{a+at-at-t}{a(1-a)} = \frac{a(1-t)-t(1-a)}{a(1-a)} = \frac{1-t}{1-a} - \frac{t}{a}$$
- $$\int \frac{\partial E}{\partial a} da = \int \frac{a-t}{a(1-a)} da = \int \frac{1-t}{1-a} da - \int \frac{t}{a} da$$
$$= -t \ln a - (1 - t) \ln(1 - a)$$





## Softmax

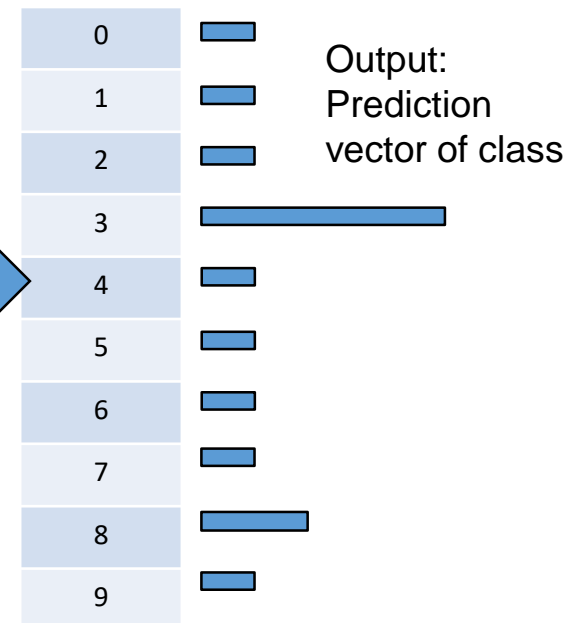
- Recall sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- Softmax is a generalization of the sigmoid:

$$y_i = \frac{e^{z_i}}{\sum_{j \in \text{group}} e^{z_j}}$$

- Softmax generalizes 2-class logistic regression/sigmoids to n classes
- Suitable for mutually exclusive classes



Input: 28x28 pixel  
grayscale image





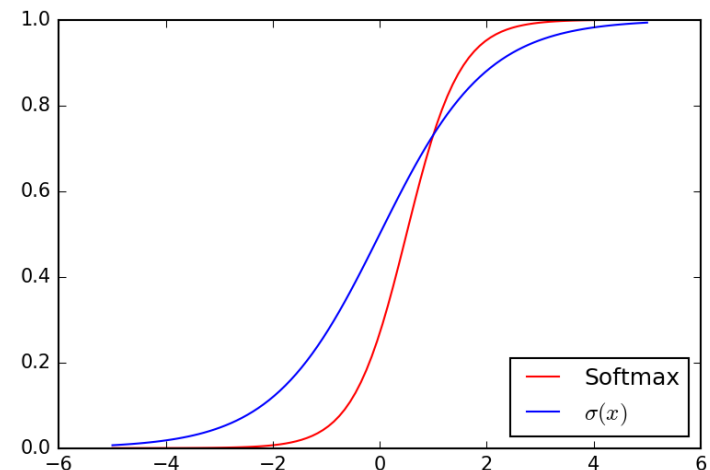
## Softmax

- Recall sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- Softmax is a generalization of the sigmoid:

$$y_i = \frac{e^{z_i}}{\sum_{j \in \text{group}} e^{z_j}}$$

- Consider for two classes with probabilities  $(z, 1-z)$ :

- $y_i = \frac{e^z}{e^z + e^{1-z}}$



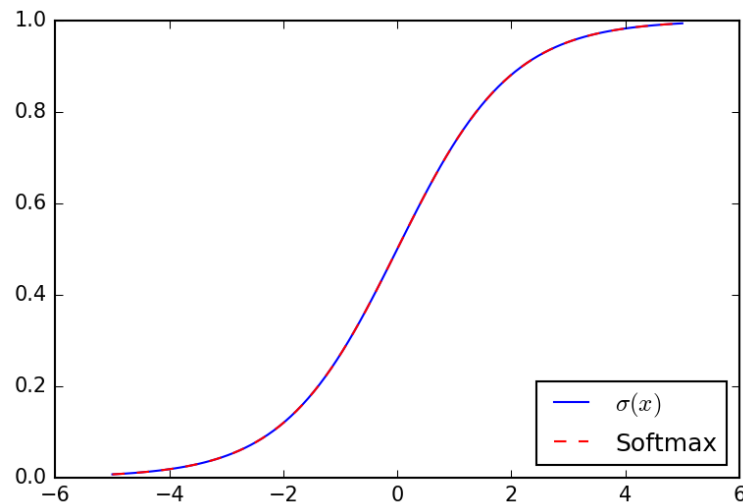


## Softmax

- Recall sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- Softmax is a generalization of the sigmoid:

$$y_i = \frac{e^{z_i}}{\sum_{j \in \text{group}} e^{z_j}}$$

- Consider for two classes with probabilities  $(z, 1-z)$ :
- $y_i = \frac{e^z}{e^z + e^{1-z}}$
- With  $\hat{z} = (z + 1)/2$ , we recover the Sigmoid function.





## Softmax

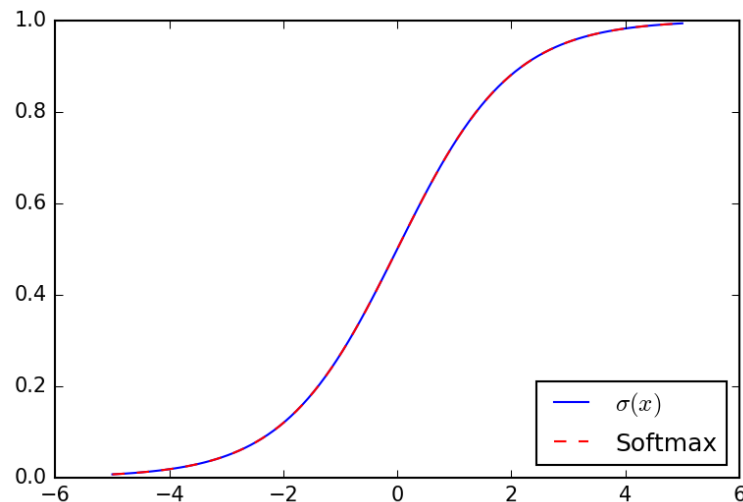
- Recall sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- Softmax is a generalization of the sigmoid:

$$y_i = \frac{e^{z_i}}{\sum_{j \in \text{group}} e^{z_j}}$$

- Consider for two classes with probabilities  $(z, 1-z)$ :

- $y_i = \frac{e^z}{e^z + e^{1-z}}$
- With  $\hat{z} = (z + 1)/2$ , we recover the Sigmoid function.
- As we found with the sigmoid,

$$y'_i = y_i(1 - y_i)$$





## Softmax

Cross entropy is the “best” cost function to use for the softmax:

$$C_j = - \sum_j t_j \ln y_j$$

For the familiar two-class case:

$$C = -t_j \ln y_j - (1 - t) \ln(1 - y)$$

$$y_i = \frac{e^z}{e^z + e^{1-z}} \text{ Class 1}$$

$$\text{Class 2} = 1 - \text{Class 1}$$

Sigmoid function.

- As we found with the sigmoid,

$$y'_i = y_i(1 - y_i)$$

- Consider for two classes with probabilities  $(z, 1-z)$ :

$$y_i = \frac{e^z}{e^z + e^{1-z}}$$

- With  $z'=(z+1)/2$ , we recover the Sigmoid function.

- As we found with the sigmoid,

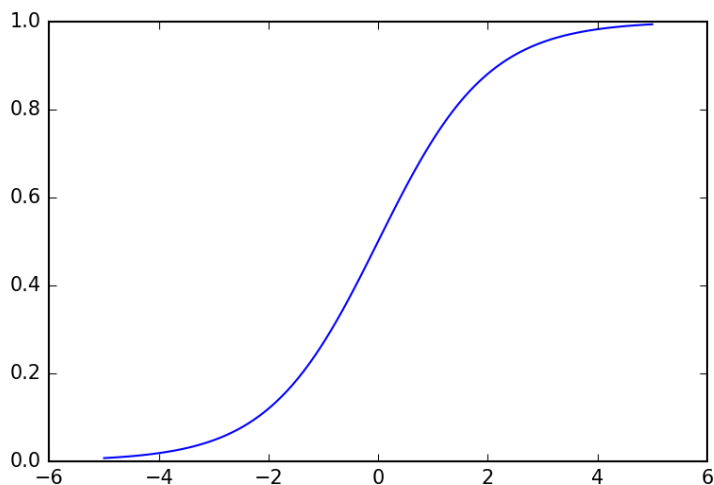
$$y'_i = y_i(1 - y_i)$$



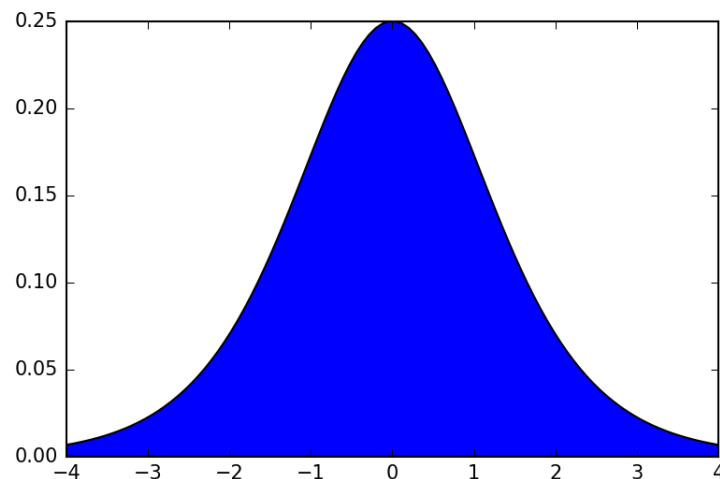
## Recall: Trouble in Sigmoid land

- Sigmoids are relatively effective in active region

$$\sigma(s)$$



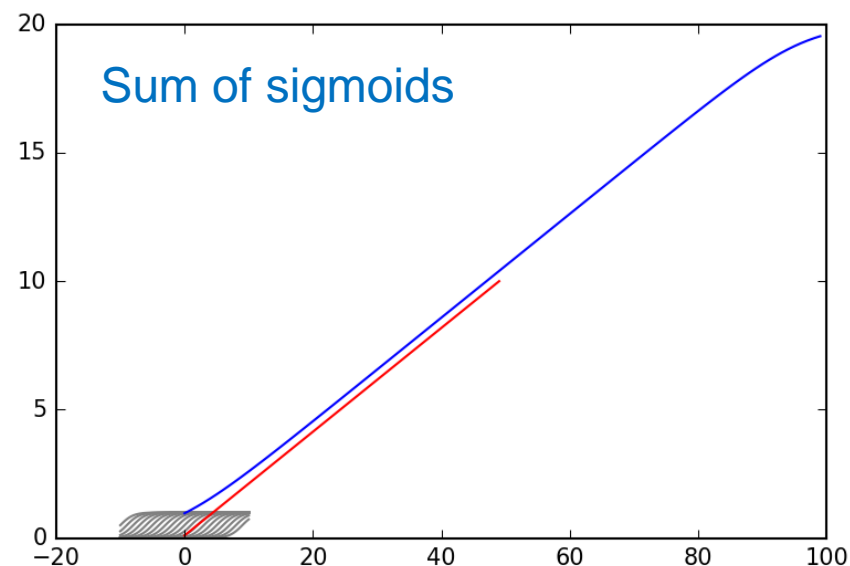
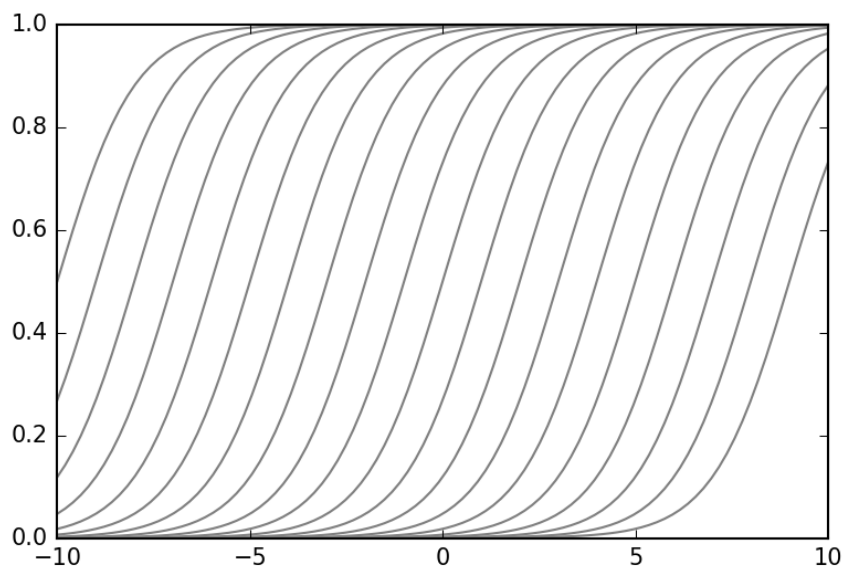
$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$





## Recall: Trouble in Sigmoid land

- Sigmoids are relatively effective in active region
- What if we replicate the sigmoids with different biases?  
$$\sigma(s) + \sigma(s - 1) + \sigma(s - i) + \cdots + \sigma(s - n)$$
- And share their weights  $w_i$ ?



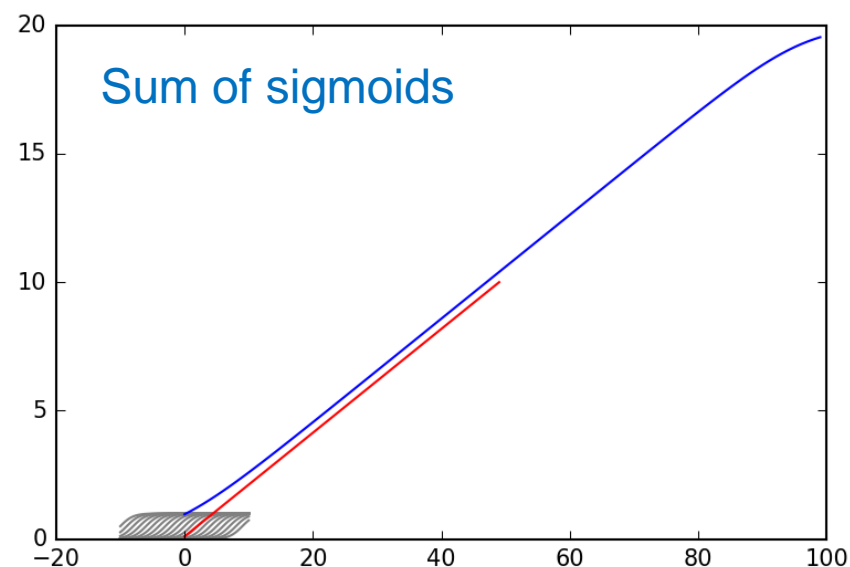
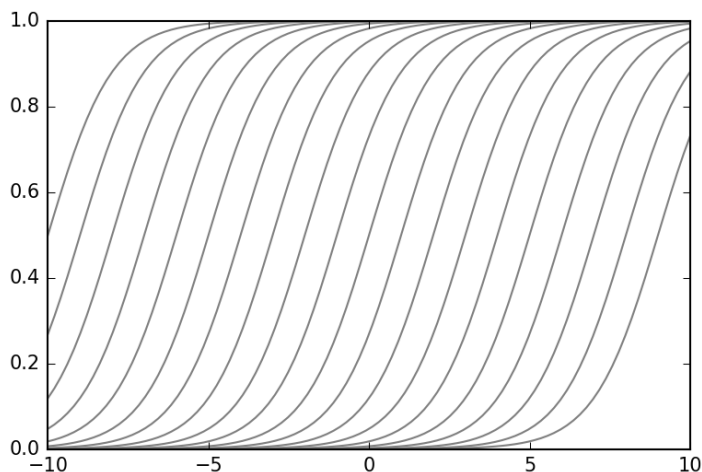


## Recall: Trouble in Sigmoid land

- What if we replicate the sigmoids with different biases?

$$\sum_{i=1}^{\infty} \sigma(s + 0.5 - i) \approx \log(1 + e^s) \approx \max(0, s + \text{noise})$$

- Called rectified linear unit: ReLU
- Easy to approximate

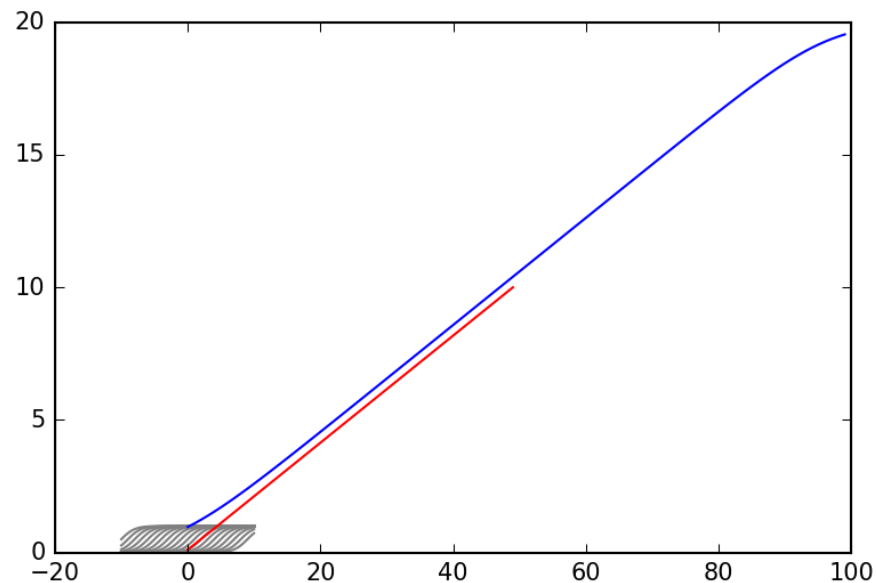






Nice try, but....

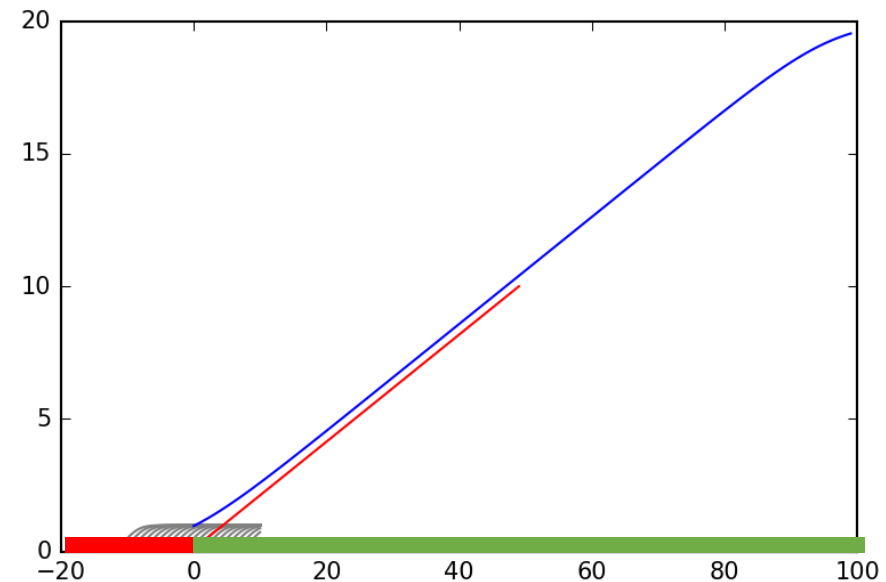
- $f(s) = \max(0, s)$
- $\frac{\partial}{\partial s} f = \begin{cases} 0, & s \leq 0 \\ 1 * ds, & s > 0 \end{cases}$





Nice try, but....

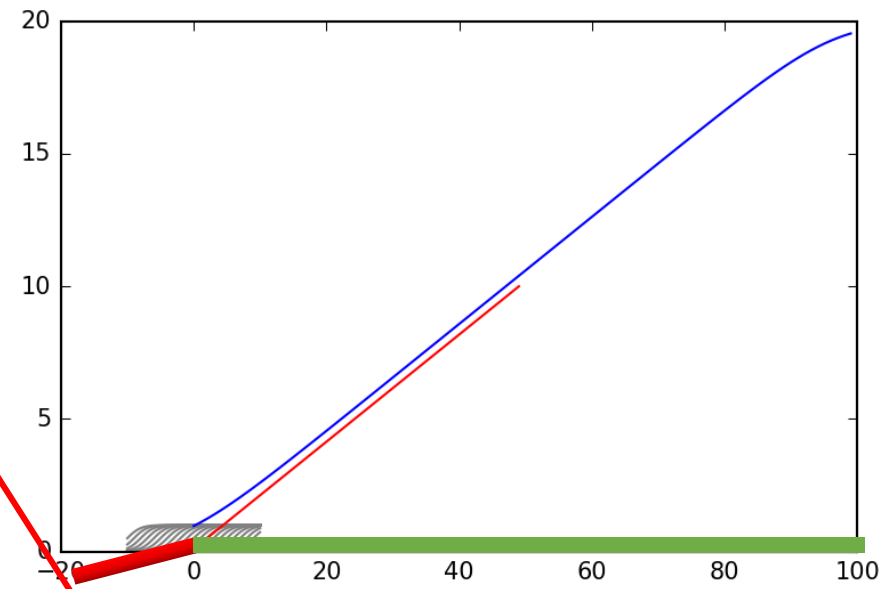
- $f(s) = \max(0, s)$
- $\frac{\partial}{\partial s} f = \begin{cases} 0, & s \leq 0 \\ 1 * ds, & s > 0 \end{cases}$
- ReLU neurons also suffer from vanishing gradients
- Neurons can “die” if they get a negative activation and never recover





Nice try, but....

- $f(s) = \max(0, s)$
- $\frac{\partial}{\partial s} f = \begin{cases} 0.01s, & s \leq 0 \\ 1 * ds, & s > 0 \end{cases}$
- ReLU neurons also suffer from vanishing gradients
- Neurons can “die” if they get a negative activation and never recover

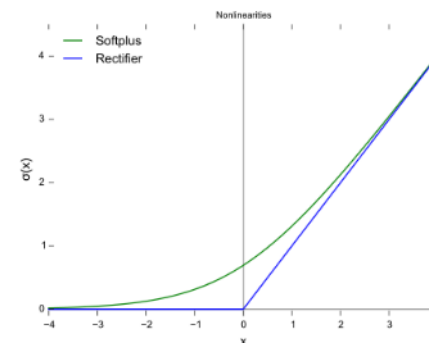
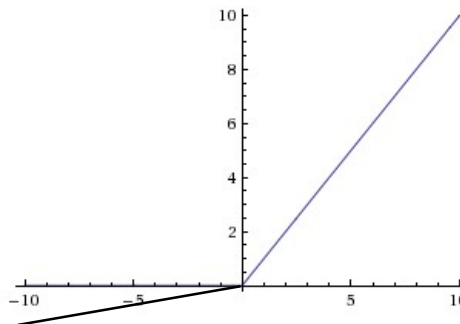
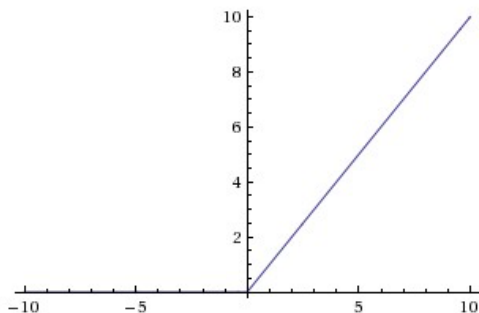
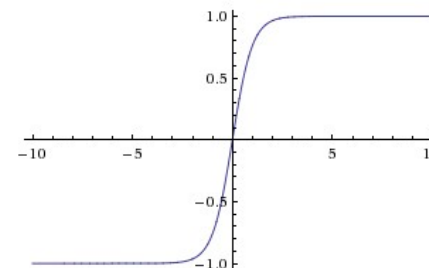
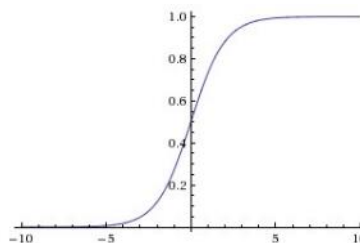


Fix:  
Clip 0 gradient to something small so neuron can “recover”  
Called: Leaky ReLU



## Activation Functions

- Sigmoid
- Hyperbolic tangent
- ReLU =  $\max(0, x)$ 
  - Leaky ReLU
  - Soft ReLU –  $\ln(1 + e^x)$



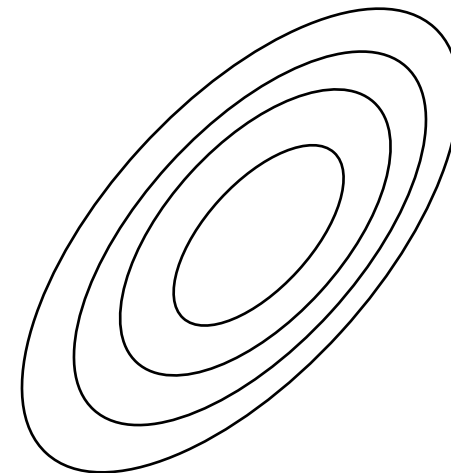


NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization

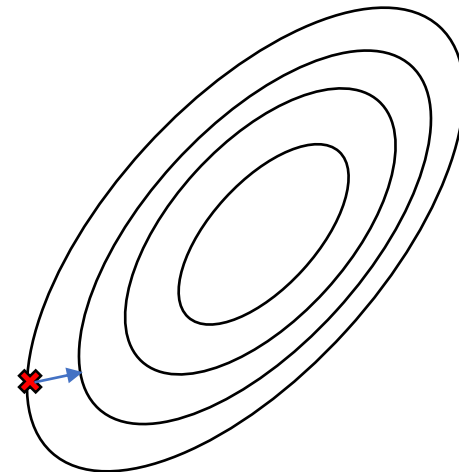


Cost function is usually  
MSE or cross entropy



## Optimization Visualized

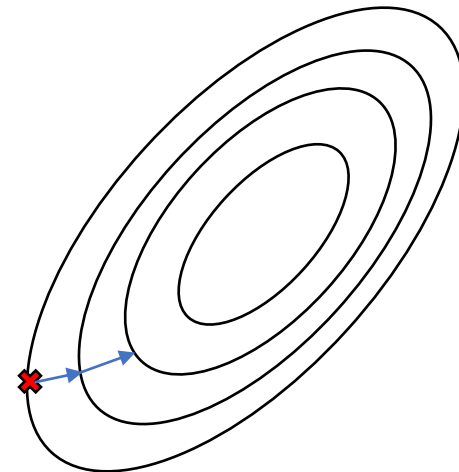
- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point





## Optimization Visualized

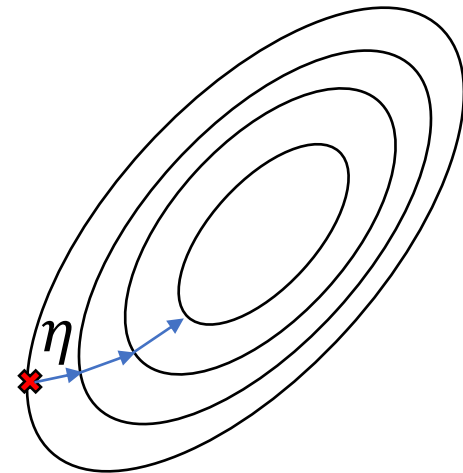
- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima





## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient



$\eta$  is known as the  
learning rate

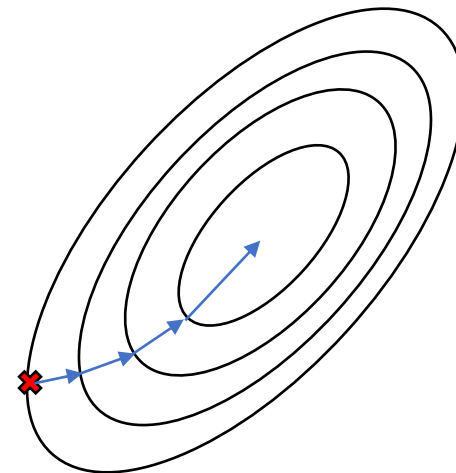
$$w_{i+1} = w_i - \eta \nabla w_i$$





## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient
- This leads to a path perpendicular to the contours

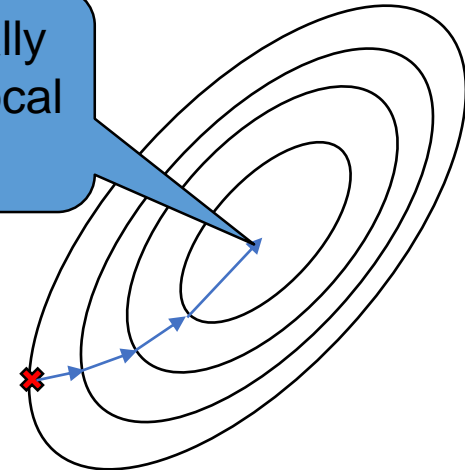




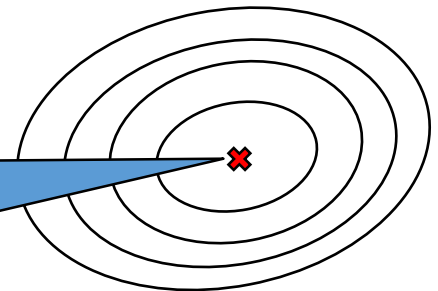
## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient
- This leads to a path perpendicular to the contours
- However, optimization can have some pitfalls...
  - Local minima

Accidentally  
found a local  
minimum



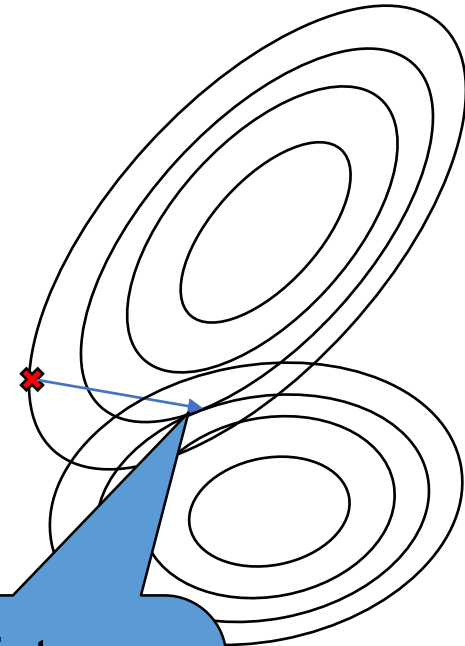
This is the  
global  
minimum





## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient
- This leads to a path perpendicular to the contours
- However, optimization can have some pitfalls...
  - Local minima
  - Saddle points



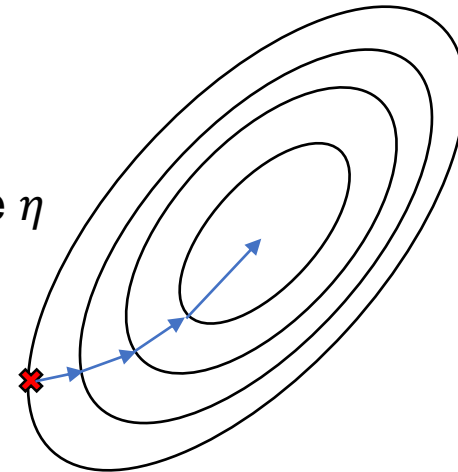
Saddle point:  
Stationary point  
that's not an  
extrema where we  
can get stuck



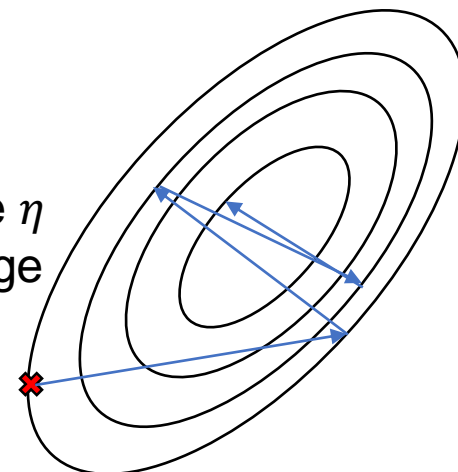
## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient
- This leads to a path perpendicular to the contours
- However, optimization can have some pitfalls...
  - Local minima
  - Saddle points
  - Bad learning rates  $\eta$

Step size  $\eta$   
is good



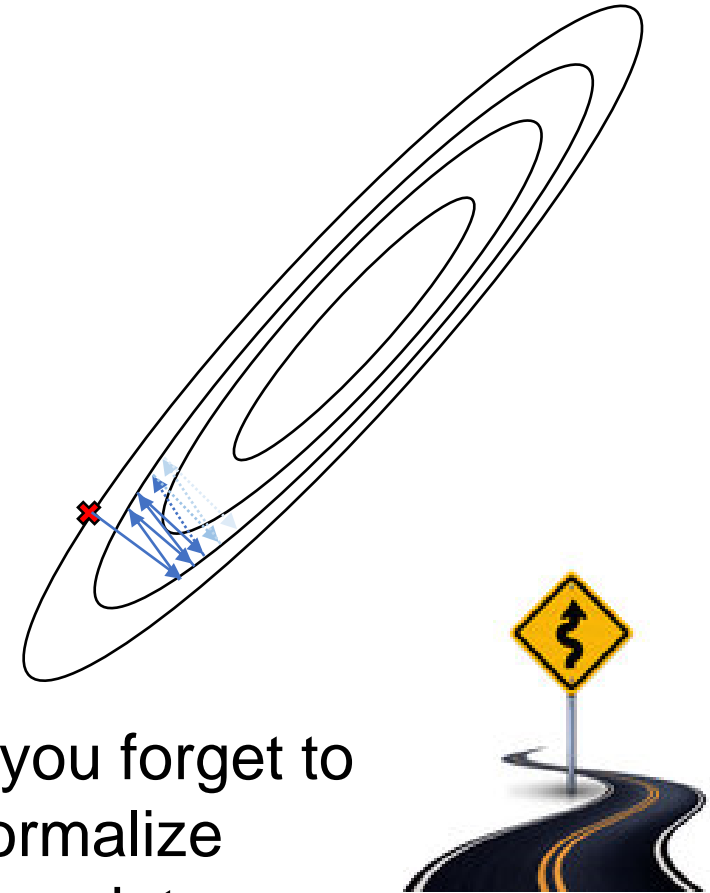
Step size  $\eta$   
is too large





## Optimization Visualized

- To find the best weights, we minimize the cost function using optimization
- Optimization starts from an initial point
- Each iteration attempts to make progress towards global minima
- Usually the fastest way is to take a step  $\eta$  in the direction of the steepest descent, called the gradient
- This leads to a path perpendicular to the contours
- However, optimization can have some pitfalls...
  - Local minima
  - Saddle points
  - Bad learning rates  $\eta$

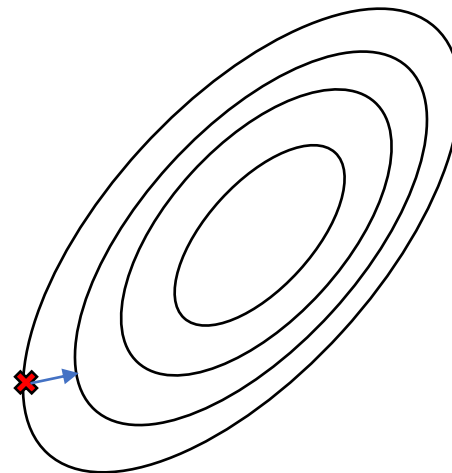


If you forget to  
normalize  
your data,  
weights or  
activations ;)



## Estimating gradients

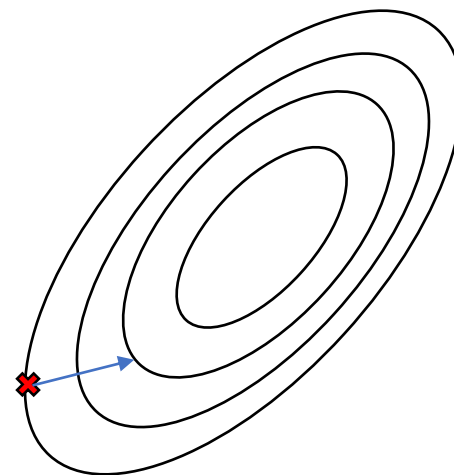
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$





## Estimating gradients

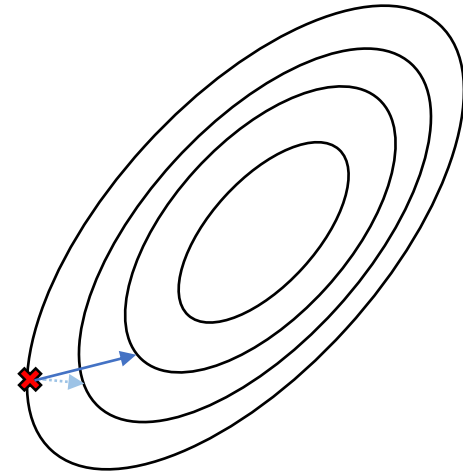
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
- Full batch means running through all the training data
  - For MNIST, this mean 60,000 samples
  - For some datasets, this could mean terabytes or petabytes of data
  - ... if we do that, perhaps we should take a larger step  $\eta$  to justify the cost of computing over all the data
  - Full gradient does the best job in finding the true gradient





## Estimating gradients

- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
- Mini batch means running through a smaller batch of the training data
  - This is much cheaper than running through all the data. Each batch is called an epoch.
  - However, the gradient estimate is more noisy, so it makes sense to take a smaller step  $\eta$
  - Depending on the conditioning of the problem, the added cost of small step sizes may outweigh the savings of running a small sample



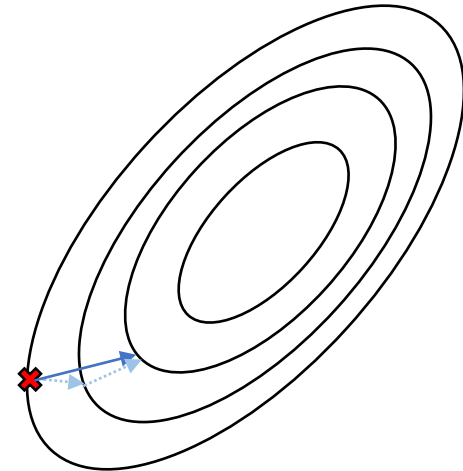
**NB:** Statisticians often take a smaller sample of the data to build a model. Doing a mini-batch is similar.





## Estimating gradients

- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
- Mini batch means running through a smaller batch of the training data
  - This is much cheaper than running through all the data. Each batch is called an epoch.
  - However, the gradient estimate is more noisy, so it makes sense to take a smaller step  $\eta$
  - Depending on the conditioning of the problem, the added cost of small step sizes may outweigh the savings of running a small sample



**NB:** Statisticians often take a smaller sample of the data to build a model. Doing a mini-batch is similar.

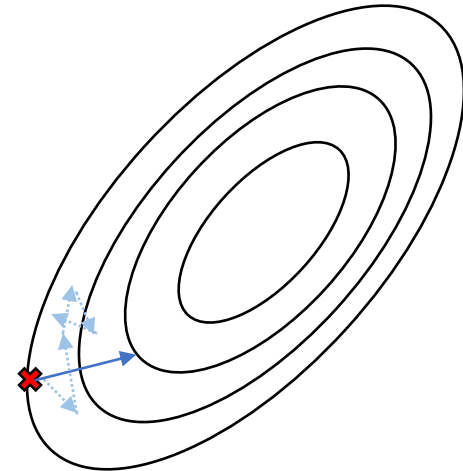
Also, we could take several mini-batches and average over them to get a better estimate

Batch size can make a difference, so choose judiciously



## Estimating gradients

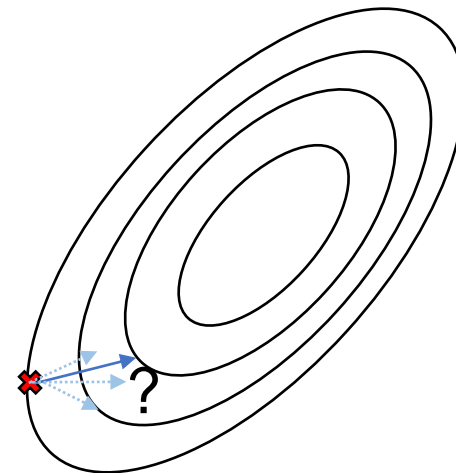
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
  - Online
- Online methods update after each training case
  - The gradient can be very noisy due to the lack of averaging
  - Smaller steps are advised
  - This could work well for streaming data sets where it's impractical to store the incoming data





## Estimating gradients

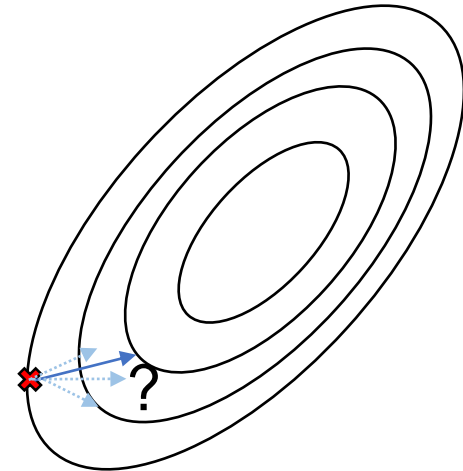
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
  - Online
- How big should the step size  $\eta$  be?





## Estimating gradients

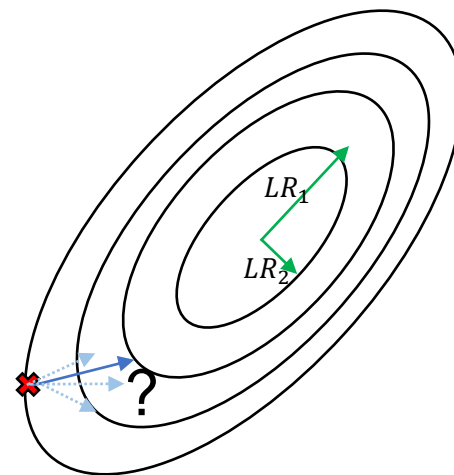
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
  - Online
- How big should the step size  $\eta$  be?
  - Fixed global learning rate?
  - Adaptive global learning rate?
  - Learning rate per direction?





## Estimating gradients

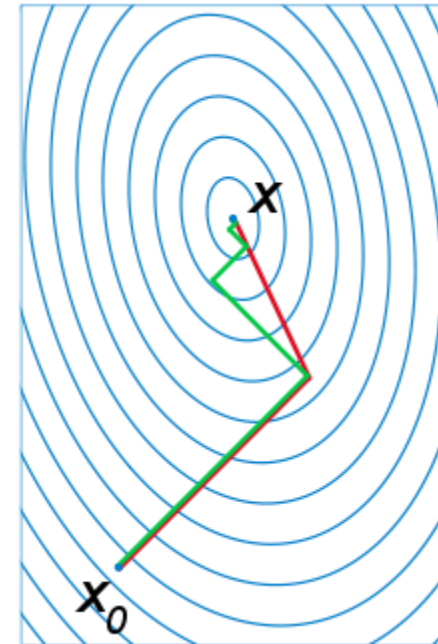
- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
  - Online
- How big should the step size  $\eta$  be?
  - Fixed global learning rate?
  - Adaptive global learning rate?
  - Learning rate per direction?





## Estimating gradients

- At each step, the optimizer computes the location of the next point estimate: direction  $\nabla w$  and step size  $\eta$
- This can be done several ways:
  - Full batch
  - Mini-batch
  - Online
- How big should the step size  $\eta$  be?
  - Fixed global learning rate?
  - Adaptive global learning rate?
  - Learning rate per direction?
  - Using something other than steepest descent?
    - E.g. Conjugate gradient



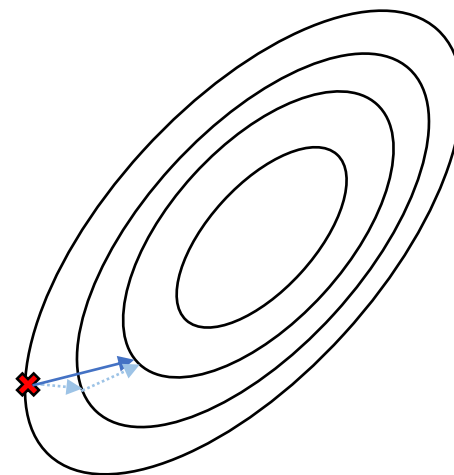
Conjugate gradient [Wikipedia]  
Green – gradient descent  
Red - Conjugate gradient

Conjugate gradient means choosing a new direction that's conjugate to the previous directions to not "mess up" the work of previous iterations.



## Speeding up mini batches

- Each mini-batch epoch estimates gradient  $\nabla w$  and step size  $\eta$  using a random subset of the data
- Since the gradient is noisy, it's usually necessary to take smaller steps



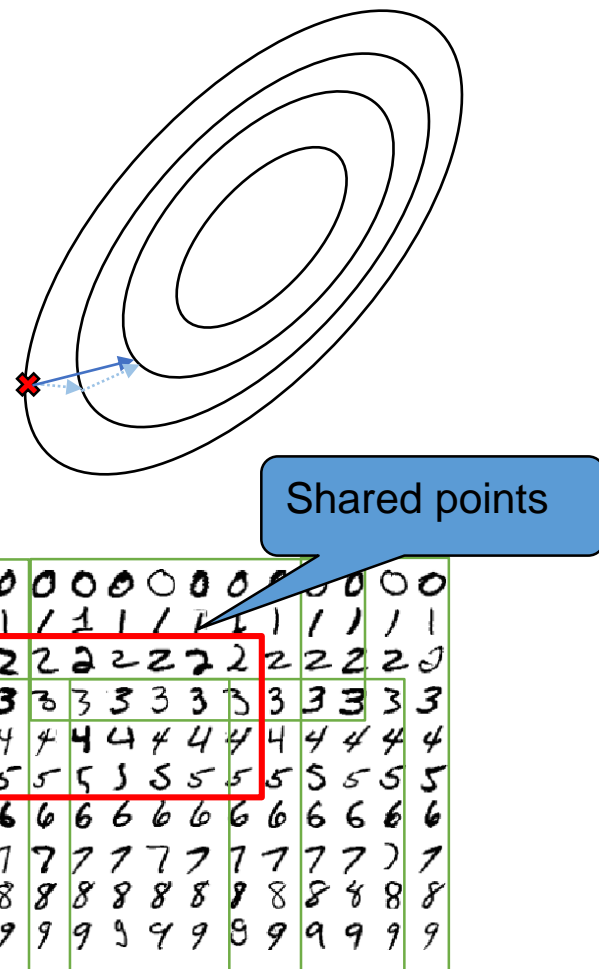
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

Various ways to slice data for minibatching



## Speeding up mini batches

- Each mini-batch epoch estimates gradient  $\nabla w$  and step size  $\eta$  using a random subset of the data
- Since the gradient is noisy, it's usually necessary to take smaller steps
- However, some of the minibatches share training examples, and hence share some part of their gradient

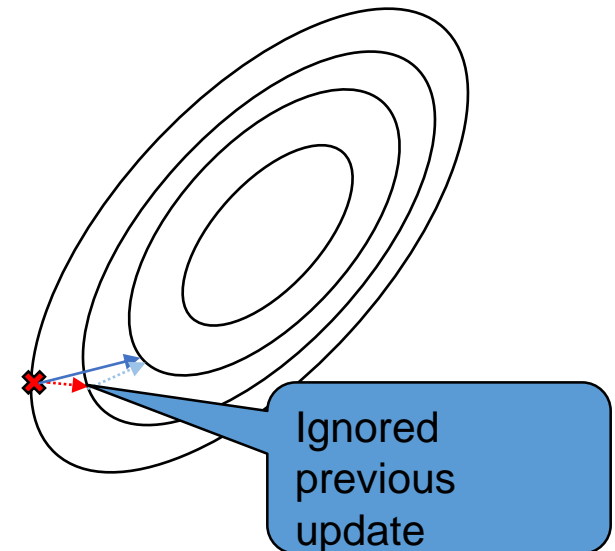






## Speeding up mini batches

- Each mini-batch epoch estimates gradient  $\nabla w$  and step size  $\eta$  using a random subset of the data
- Since the gradient is noisy, it's usually necessary to take smaller steps
- However, some of the minibatches share training examples, and hence share some part of their gradient
- Also, information about previous directions of travel are ignored...

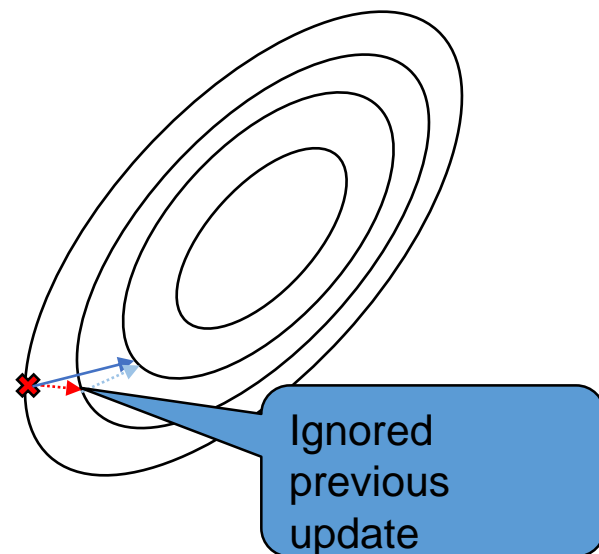


0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9



## Speeding up mini batches

- Each mini-batch epoch estimates gradient  $\nabla w$  and step size  $\eta$  using a random subset of the data
- Since the gradient is noisy, it's usually necessary to take smaller steps
- However, some of the minibatches share training examples, and hence share some part of their gradient
- Also, information about previous directions of travel are ignored...
- We can exploit these properties to speed up mini-batches by making a better estimate of the gradient to take larger steps with more confidence



0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9

Various ways to slice data for minibatching

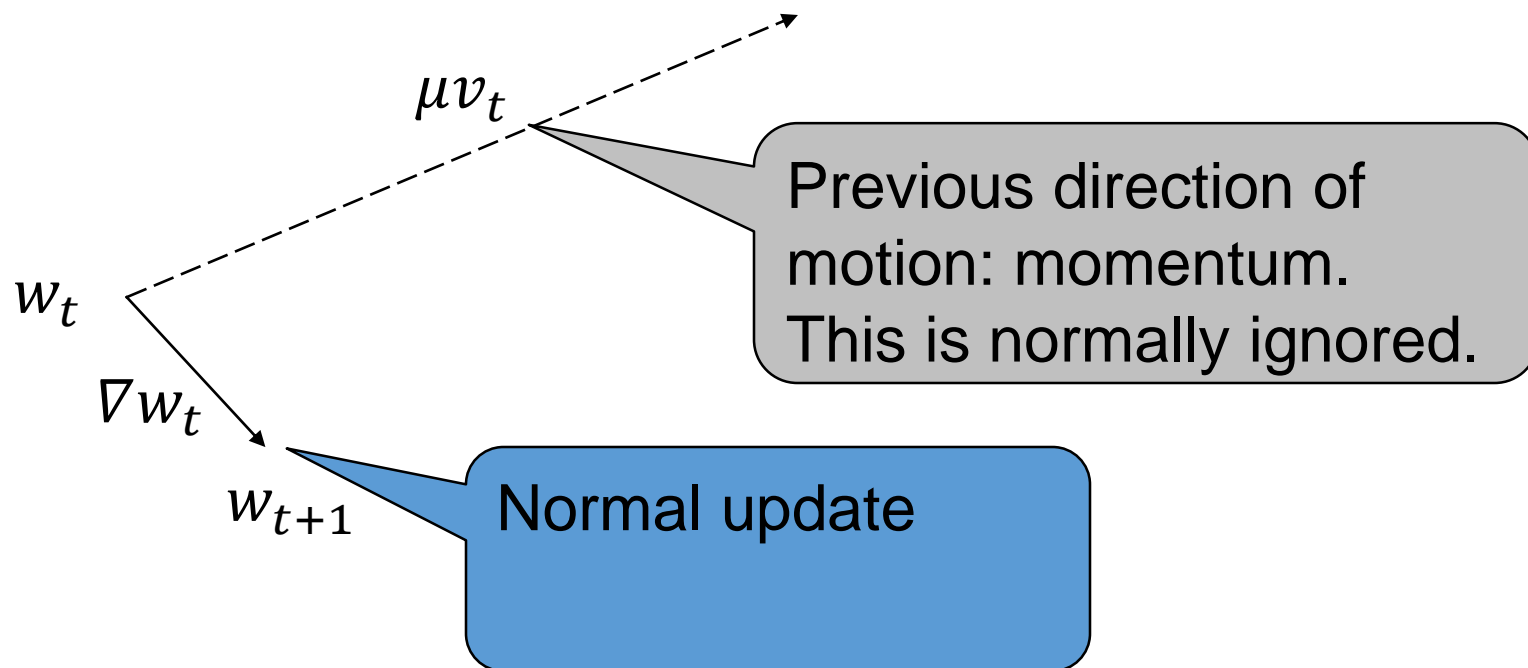


# Momentum

## Momentum

$$v_{t+1} = \mu v_t - \eta \nabla w_t$$

$$w_{t+1} = w_t + v_{t+1}$$



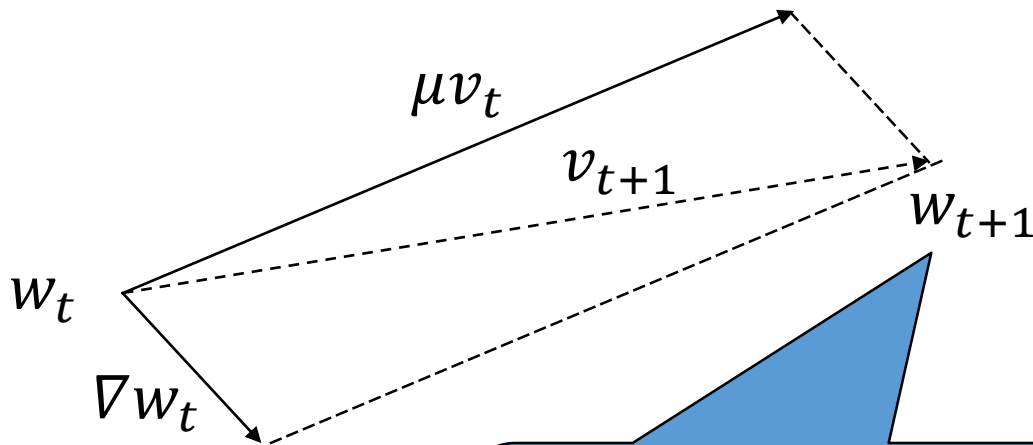


## Momentum

### Momentum

$$v_{t+1} = \mu v_t - \eta \nabla w_t$$

$$w_{t+1} = w_t + v_{t+1}$$



Momentum-adjusted update.  
Note how momentum is favored  
over noisy gradient estimate.

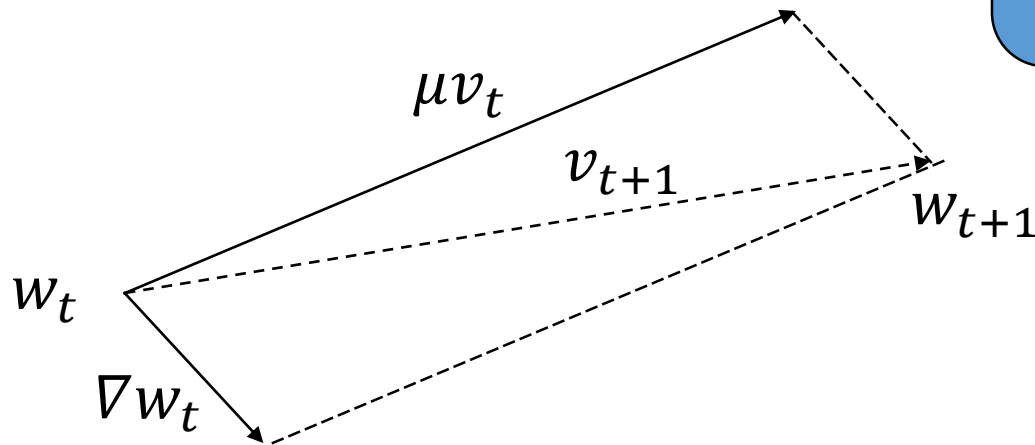


## Momentum and acceleration

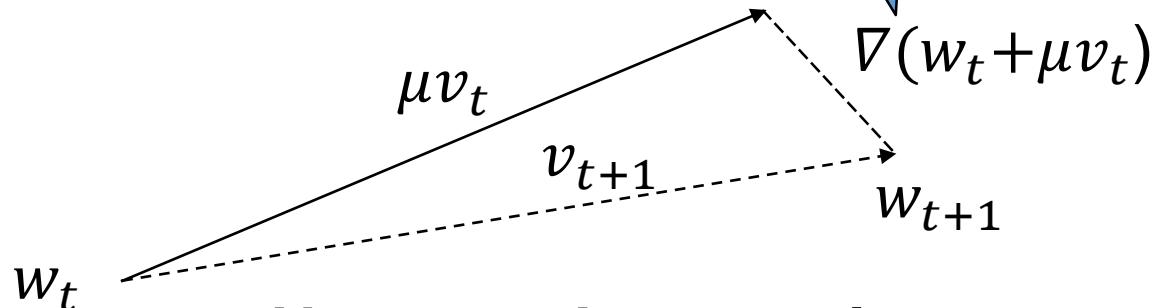
### Momentum

$$v_{t+1} = \mu v_t - \eta \nabla w_t$$

$$w_{t+1} = w_t + v_{t+1}$$



Acceleration (2<sup>nd</sup> derivative) is now accounted for



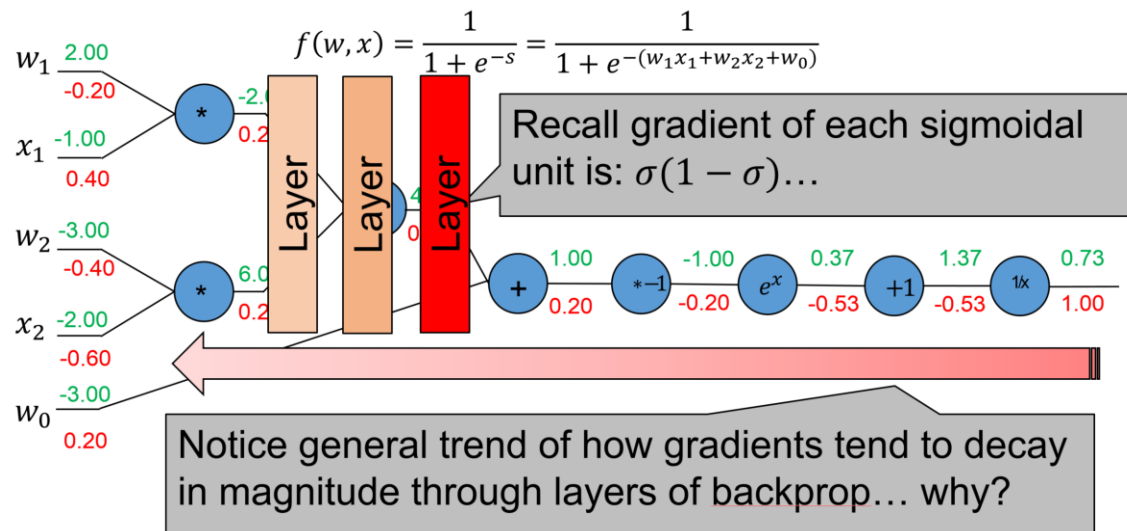
### Nesterov Acceleration

Sutskever, ICML 2013



## Other ways to speed up mini batches

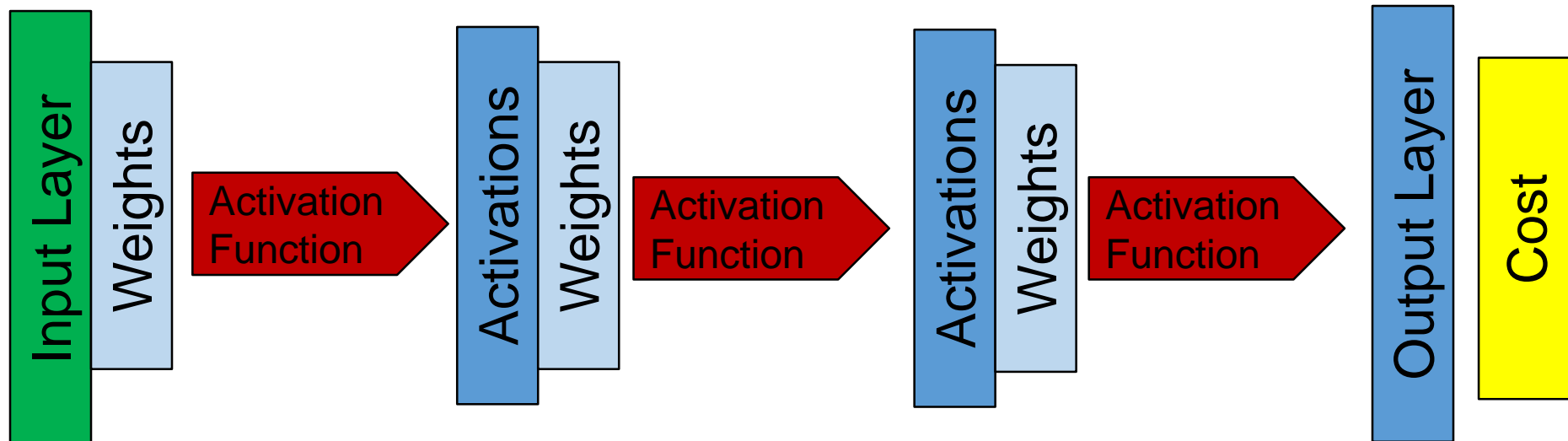
- Adaptive learning rates
  - Remember that earlier layers tend to have smaller gradients



- Use an adaptive multiplier to boost/attenuate gradients for each layer/weight
- RMSprop: Keep a moving average of the squared gradient of each weight, divide gradient by this to normalize



## Brute force

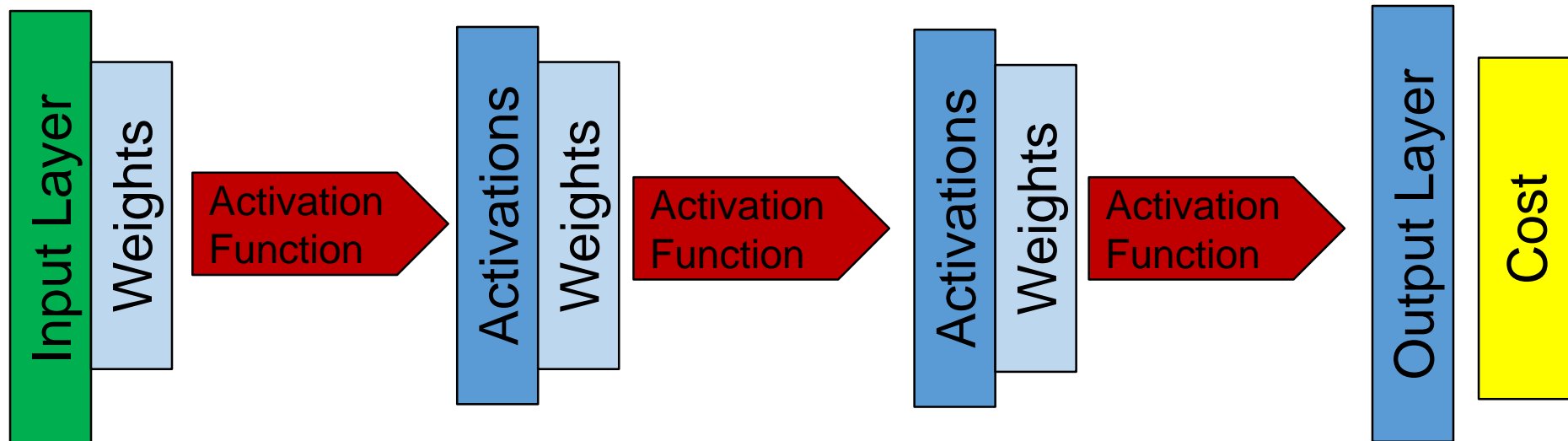


- Brute force methods





## Brute force



- Brute force methods
- Commercial networks can have billions of parameters, > 30 levels, can take weeks to months to train
- MSR: 152 layers
- <http://www.wired.com/2016/01/microsoft-neural-net-shows-deep-learning-can-get-way-deeper/>







NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

## GPUs for Deep Learning

VGG net takes ~2-3  
weeks to train on 4 GPUs



NVIDIA Digits  
Devbox ~\$10k



~\$1k per Titan X

Four TITAN X GPUs with 7 TFlops of single precision,  
336.5 GB/s of memory bandwidth, and 12 GB of  
memory per board.

Source: NVIDIA





## GPUs

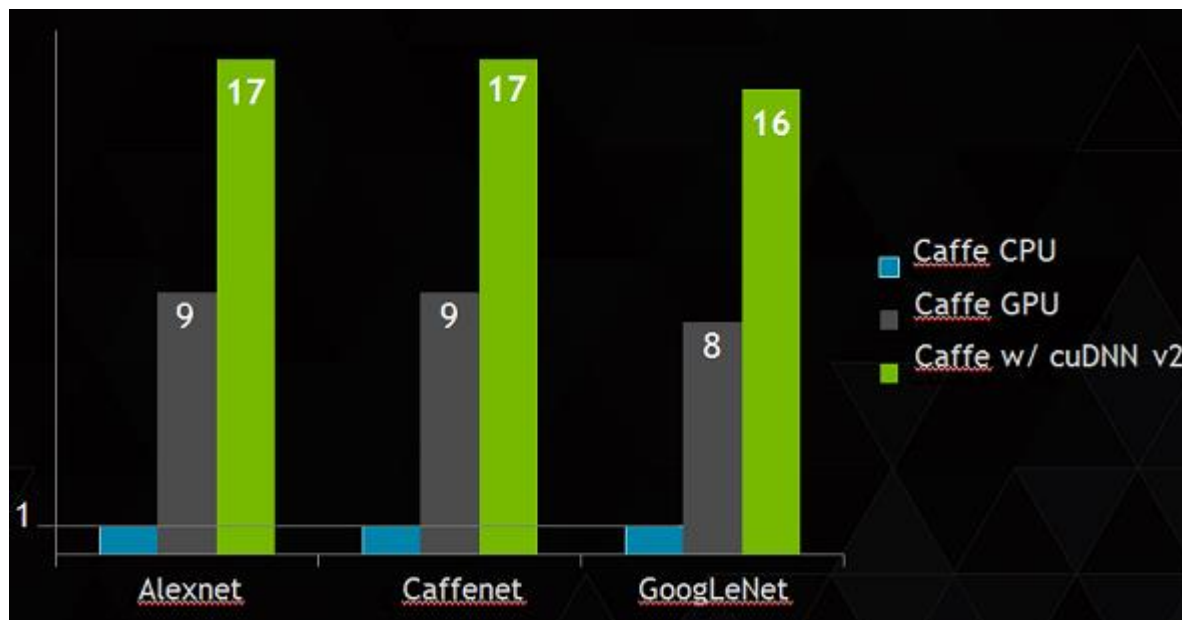
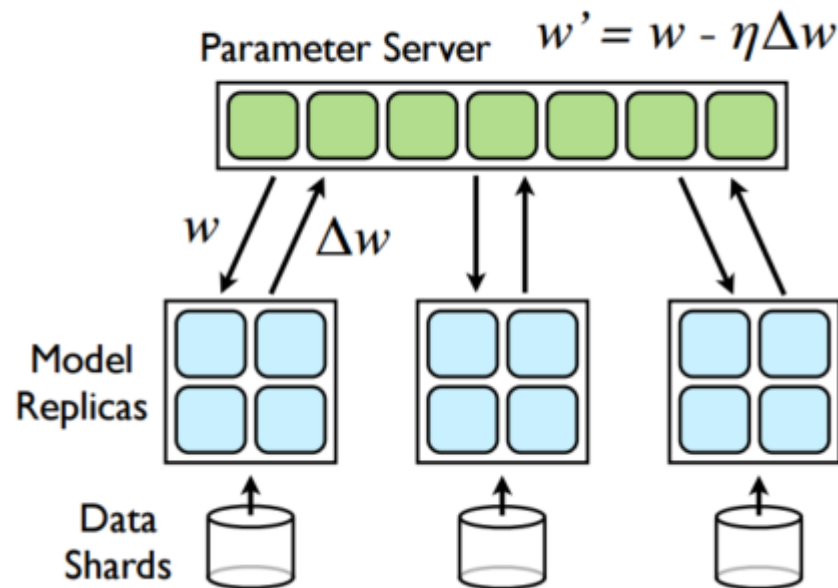
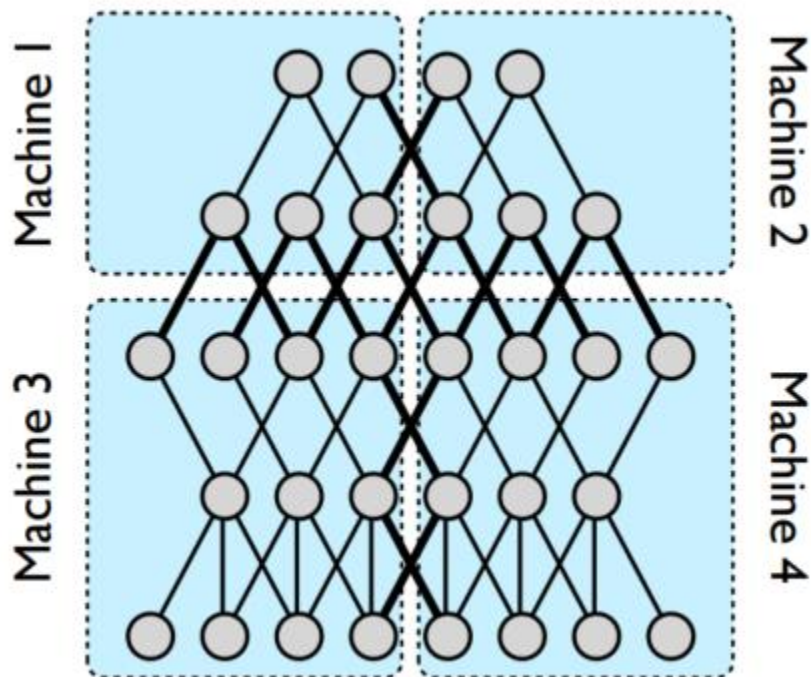


Figure 1: cuDNN performance comparison with CAFFE, using several well known networks. CPU is 16-core Intel Haswell E5-2698 2.3 GHz with 3.6 GHz Turbo. GPU is NVIDIA GeForce GTX TITAN X.



## Deep learning parallelized

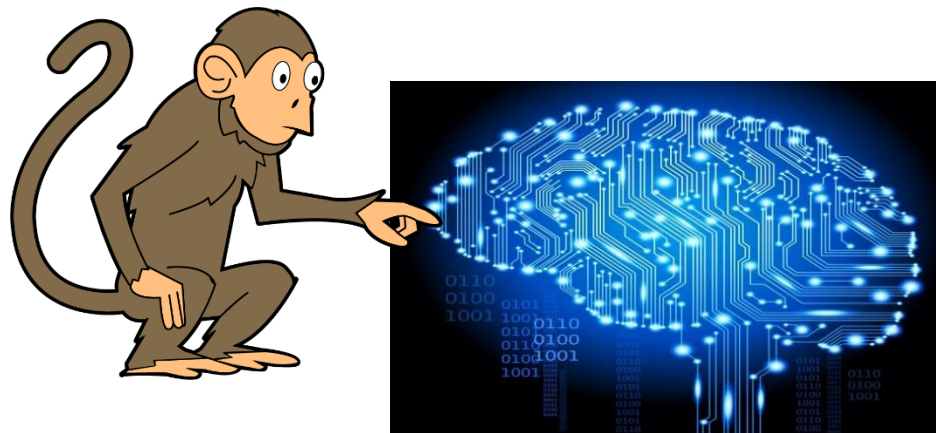


- Divide model across machines
- Sync model periodically
- Out of sync states act as noise  $\rightarrow$  SGD can handle this



## In practice (for this class)

- Training models takes a lot of GPU/CPU and memory
- Luckily, for students like us, there's:
  - A set of pretrained Keras/Tensorflow models we can play with





## Summary

- **Data prep:**

- Remember to clean data by subtracting the mean and normalizing inputs
- If possible, try to decorrelate inputs



- **Network architecture:**

- Use ReLU activations instead of sigmoid or tanh
- Use softmax or cross entropy if modeling mutually exclusive classes: e.g. digits

- **Training:**

- Try using minibatch training with momentum
- Be careful about learning rates



NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

## Assignment

- Assignment
- Quizzes
- Project





NORTHWESTERN  
UNIVERSITY

MSIA 490-29: Deep Learning. Spring 2017.  
Instructor: Dr. Ellick Chan. TA: Mark Harmon.

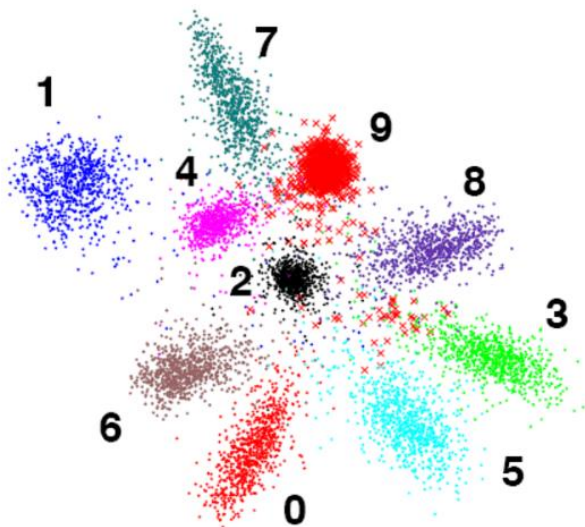
## Assignment

- Discuss assignment
- Random labels & NN capacity



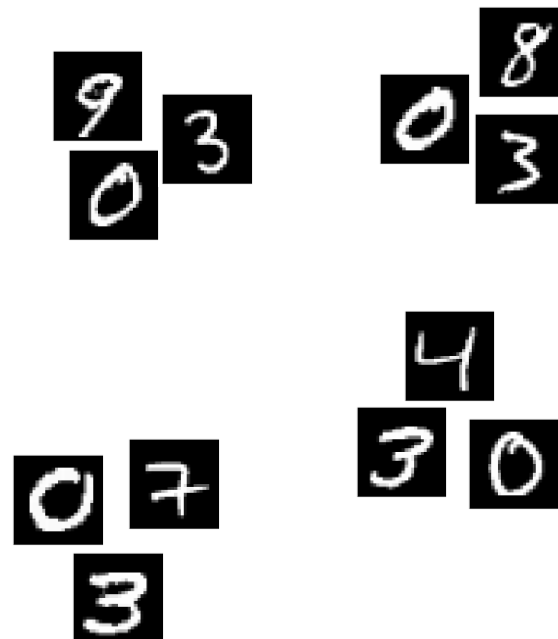


## Random labels



### Non-random labels

Show some spatial consistency in high dimensional space (clustering)



### Random labels

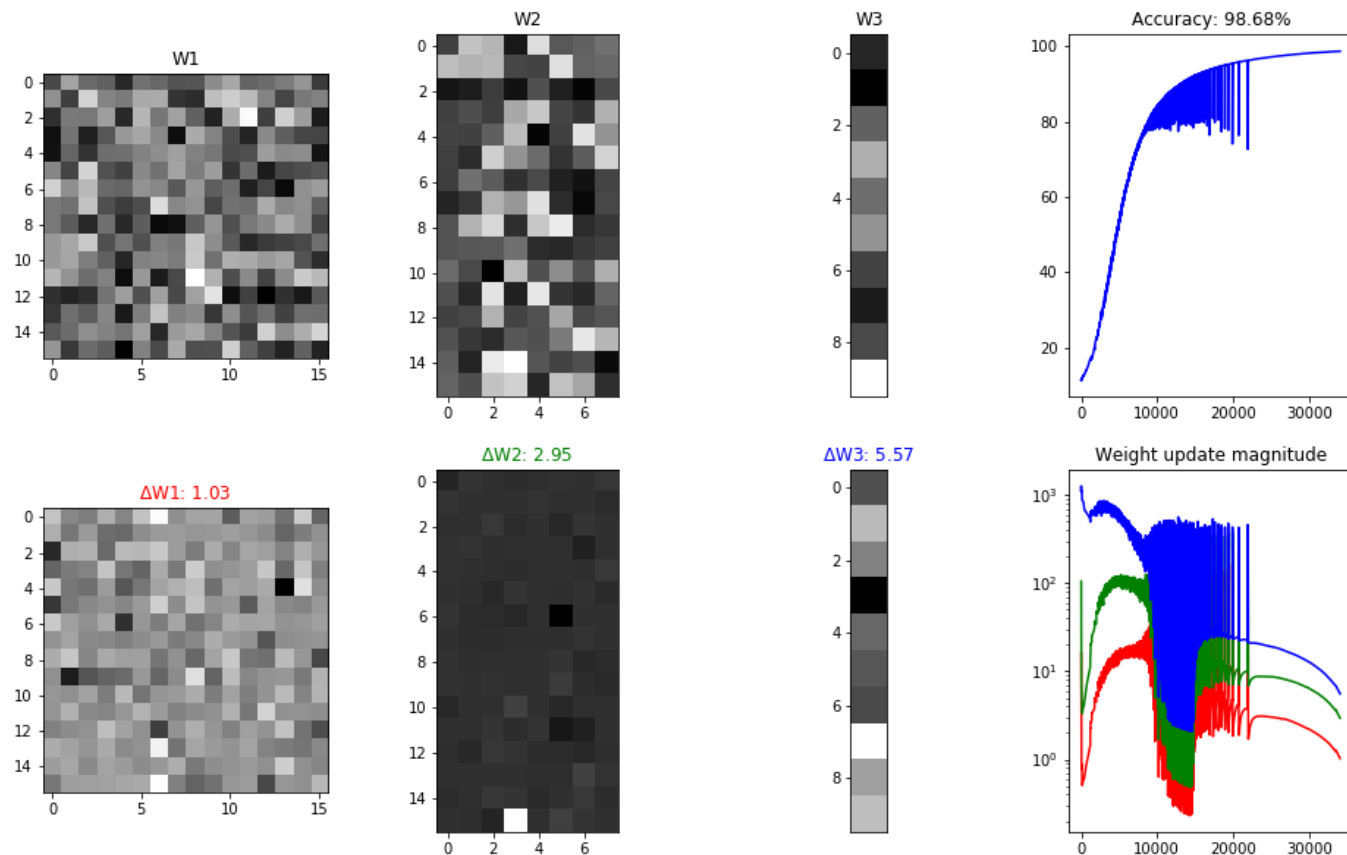
Show no spatial consistency in high dimensional space. Relies on model overfitting.





## Random labels

Weight and update visualization ACC: 98.68%





## Further reading

- Batch Normalization:  
<http://arxiv.org/pdf/1502.03167.pdf>
- 1-bit gradient updates:  
<http://arxiv.org/pdf/1602.02830v1.pdf>
- Drop-in for training deep nets:  
<http://arxiv.org/pdf/1511.06951v1.pdf>
- Weight normalization  
<https://arxiv.org/abs/1602.07868>





## Optional Reading

- Learning to learn gradient descent by gradient descent:  
<https://arxiv.org/abs/1606.04474>
- Synthetic gradient: <https://arxiv.org/abs/1608.05343>
- Learning to reinforcement learn:  
<https://arxiv.org/abs/1611.05763>
- Batch Normalization:  
<https://arxiv.org/pdf/1502.03167.pdf>
- Data-dependent init:  
<https://arxiv.org/pdf/1511.06856.pdf>



## Optional Reading

- Continuum of activations:  
<http://uaf46365.ddns.uark.edu/continuum.pdf>
- Designing NN using RL: <https://arxiv.org/abs/1611.02167>
- Neural arch search using RL:  
<https://openreview.net/forum?id=r1Ue8Hcxg>
- Sacred metaoptimizer: <https://github.com/IDSIA/sacred>
- What is softmax? <http://www.kdnuggets.com/2016/07/softmax-regression-related-logistic-regression.html>
- Overcoming forgetting:  
<https://arxiv.org/abs/1612.00796>
- Local minima: <https://arxiv.org/abs/1611.06310>



## Optional Reading

- Hyperopt:
  - <https://github.com/hyperopt/hyperopt-sklearn>
  - <http://optunity.readthedocs.io/en/latest/>
- Tanh paper: [2010]  
<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf> - ReLU invented in 2011 (Wikipedia)
- Overview of gradient descent algs:  
<http://sebastianruder.com/optimizing-gradient-descent/>



## Optional Reading

- More on training and tanh:  
<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>
- More on ReLU:  
<http://www.jmlr.org/proceedings/papers/v15/glorot11a/glorot11a.pdf>