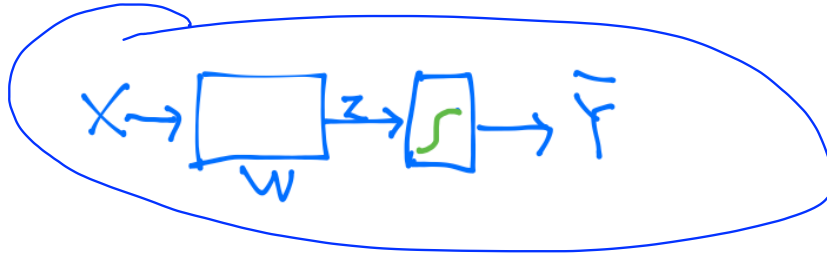


Lecture 9-I

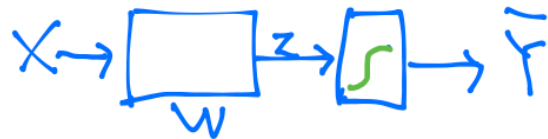
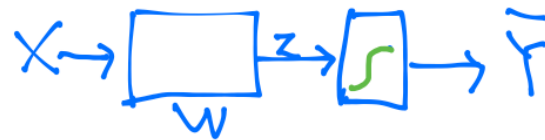
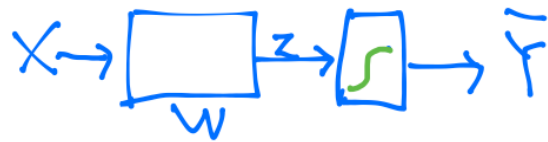
Neural Nets(NN) for XOR

Sung Kim <hunkim+mr@gmail.com>

One logistic regression unit cannot separate XOR

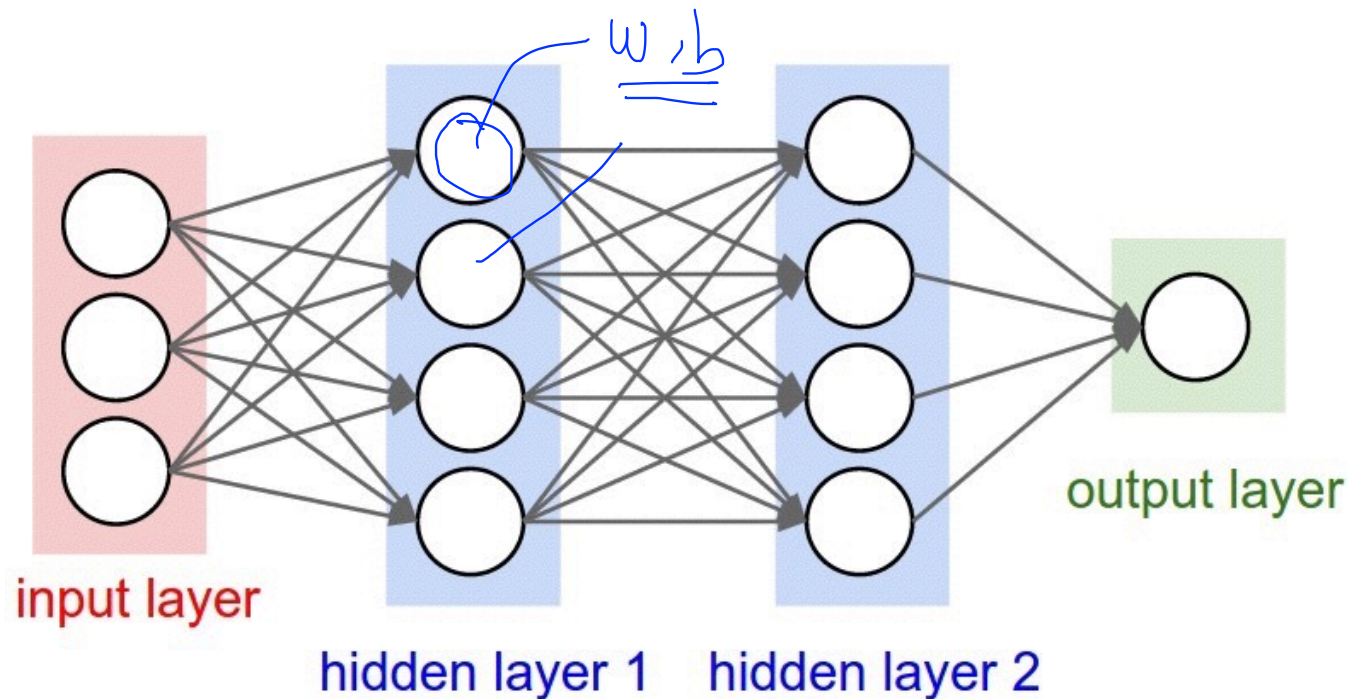


Multiple logistic regression units



Neural Network (NN)

“No one on earth had found a viable way to train*”

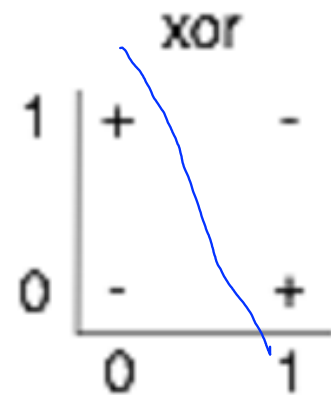


*Marvin Minsky

<http://cs231n.github.io/convolutional-networks/>

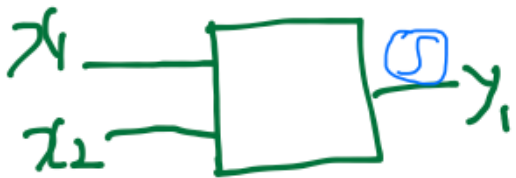
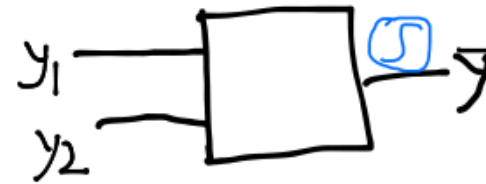
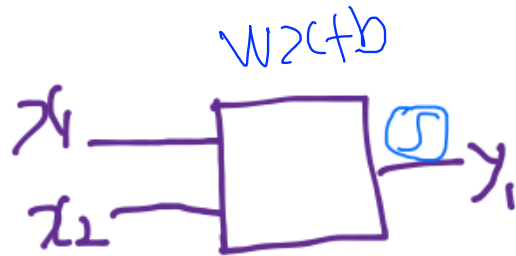
XOR using NN

x_1	x_2	XOR
0	0	0 (-)
0	1	1 (+)
1	0	1 (+)
1	1	0 (-)

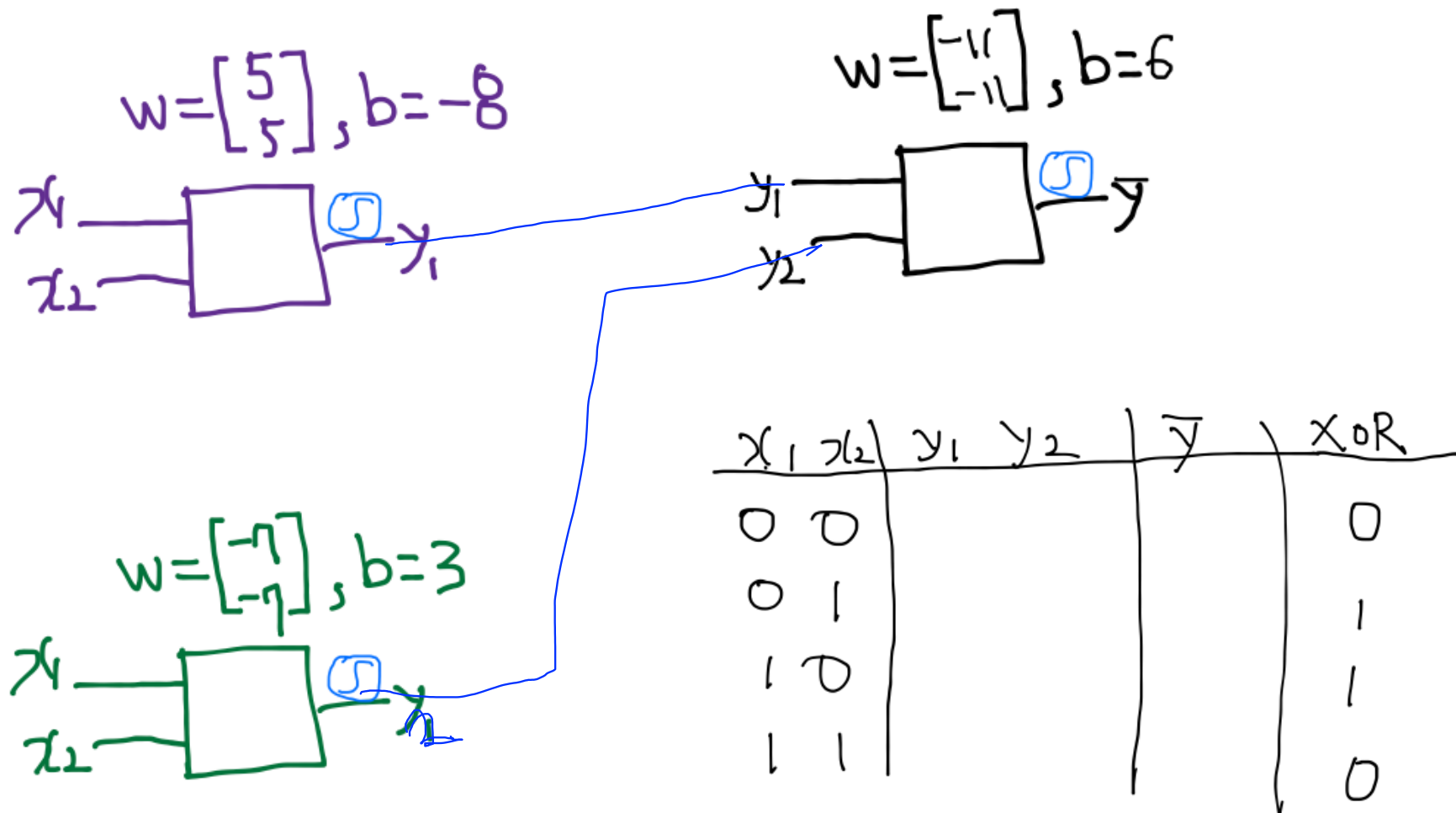


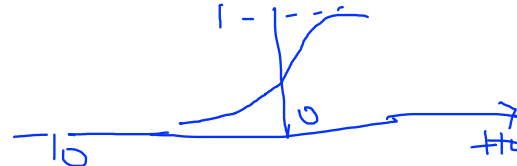
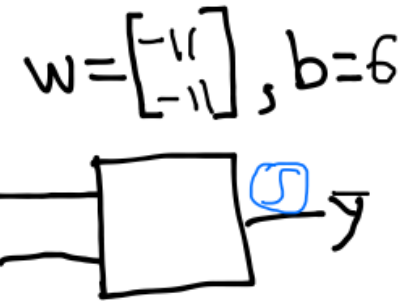
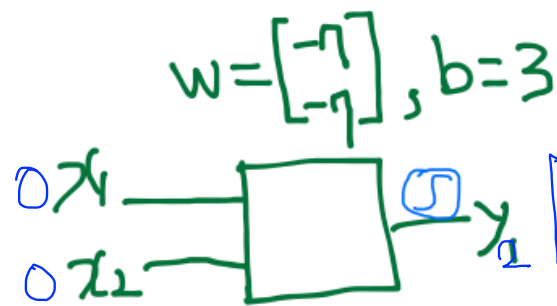
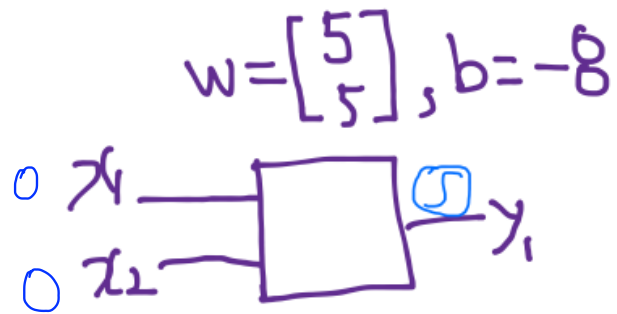
Nope

Neural Net



Neural Net





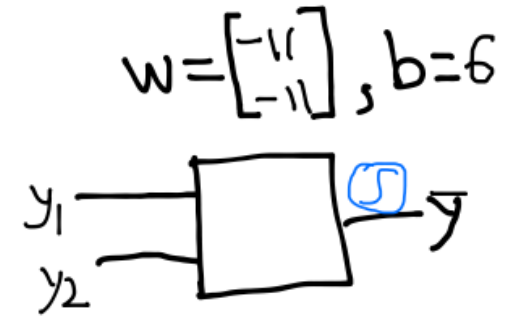
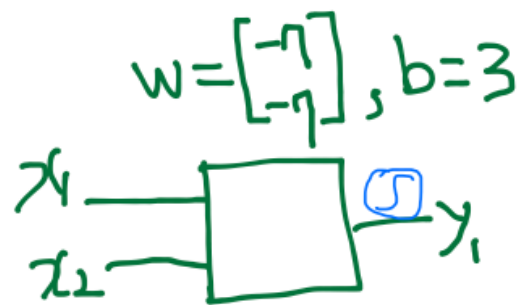
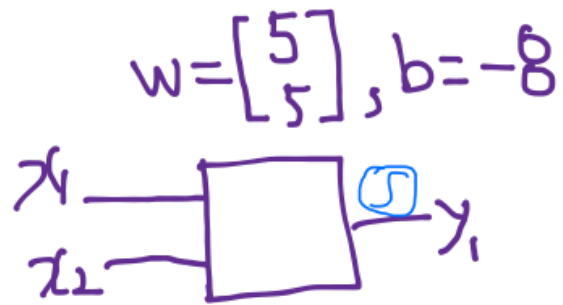
$$[0 \ 0] \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 8 = -8, \quad y_1 = S(-8) = 0$$

$$[0 \ 0] \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 3 = 3, \quad y_2 = S(3) = 1$$

$$[0 \ 1] \begin{bmatrix} -11 \\ -11 \end{bmatrix} + 6 = -11 + 6 = -5$$

$$\bar{y} = S(-5) = 0$$

x_1	x_2	y_1	y_2	\bar{y}	$x \text{ OR } \bar{y}$
0	0	0	1	0	0
0	1				1
1	0				1
1	1				0



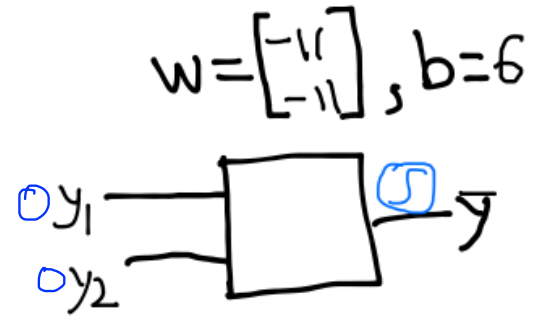
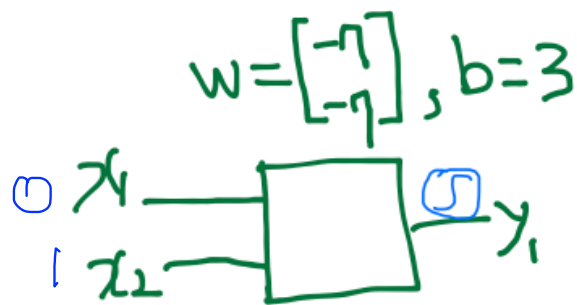
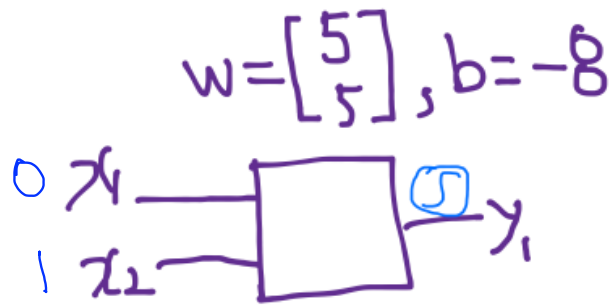
$$[0 \ 0] \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 8 = 0 + 0 - 8 = -8, \text{Sigmoid}(-8) = 0$$

$$[0 \ 0] \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 3 = 0 + 0 + 3 = 3, \text{Sigmoid}(3) = 1$$

$$[0 \ 1] \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 6 = 0 + -1 + 6 = 5$$

$$\text{Sigmoid}(5) = 1$$

x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1				1
1	0				1
1	1				0



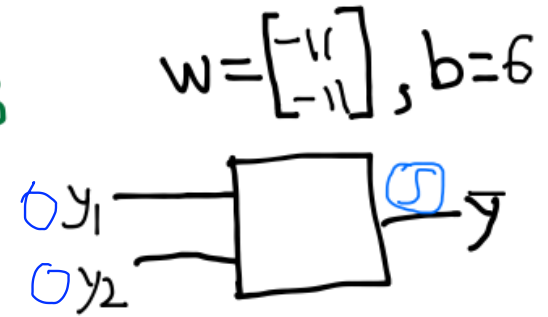
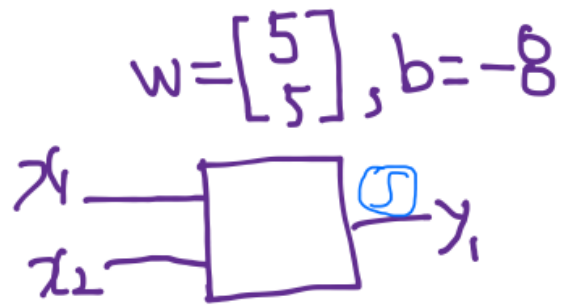
$$\begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 8 = 0 + 5 - 8 = -3, \text{Sigmoid}(-3) = 0$$

$$\begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 3 = 0 + -1 + 3 = -4, \text{Sigmoid}(-4) = 0$$

$$\begin{bmatrix} 0 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 6 = 0 + 0 + 6 = 6$$

$$\text{Sigmoid}(6) = 1$$

x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0	1	1	0	1
1	1	1	1	0	0



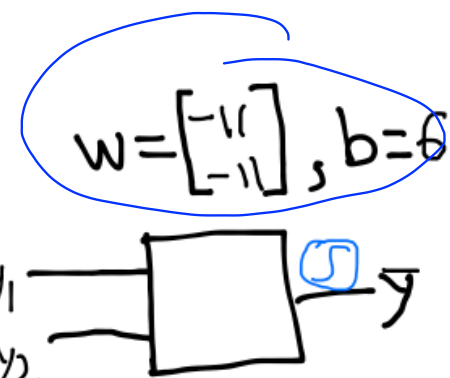
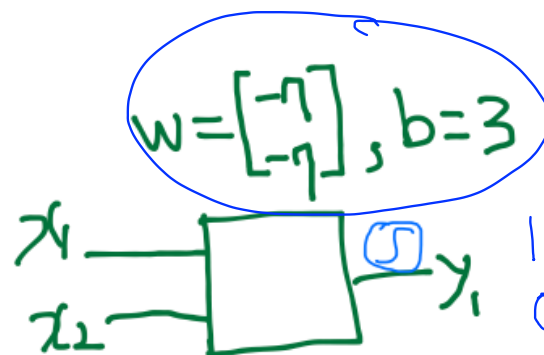
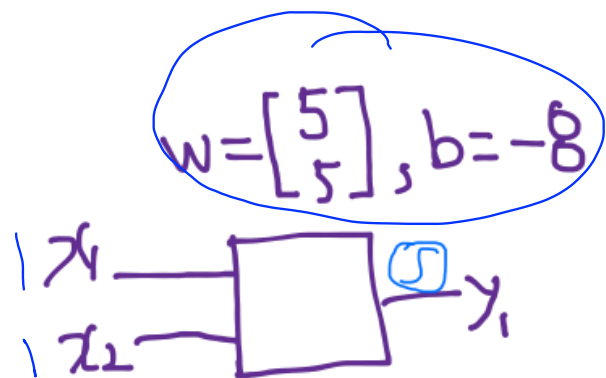
$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 8 = 5 + 0 - 8 = \underline{\underline{-3}}, \text{Sigmoid}(\underline{\underline{-3}}) = \underline{\underline{0}}$$

$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} -7 \\ -7 \end{bmatrix} + 3 = -7 + 0 + 3 = \underline{\underline{-4}}, \text{Sigmoid}(\underline{\underline{-4}}) = \underline{\underline{0}}$$

$$\begin{bmatrix} 0 & 0 \end{bmatrix} \begin{bmatrix} -11 \\ -11 \end{bmatrix} + 6 = 0 + 0 + 6 = 6$$

$$\text{Sigmoid}(6) = 1$$

x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0 ✓
0	1	0	0	1	1 ✓
<u>1</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>1</u>	1 ✓
1	1				0 ?



$$[1 \ 1] \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 8 = \underline{5} + \underline{5} - 8 = \underline{2}, \text{Sigmoid}(2) = \underline{1}$$

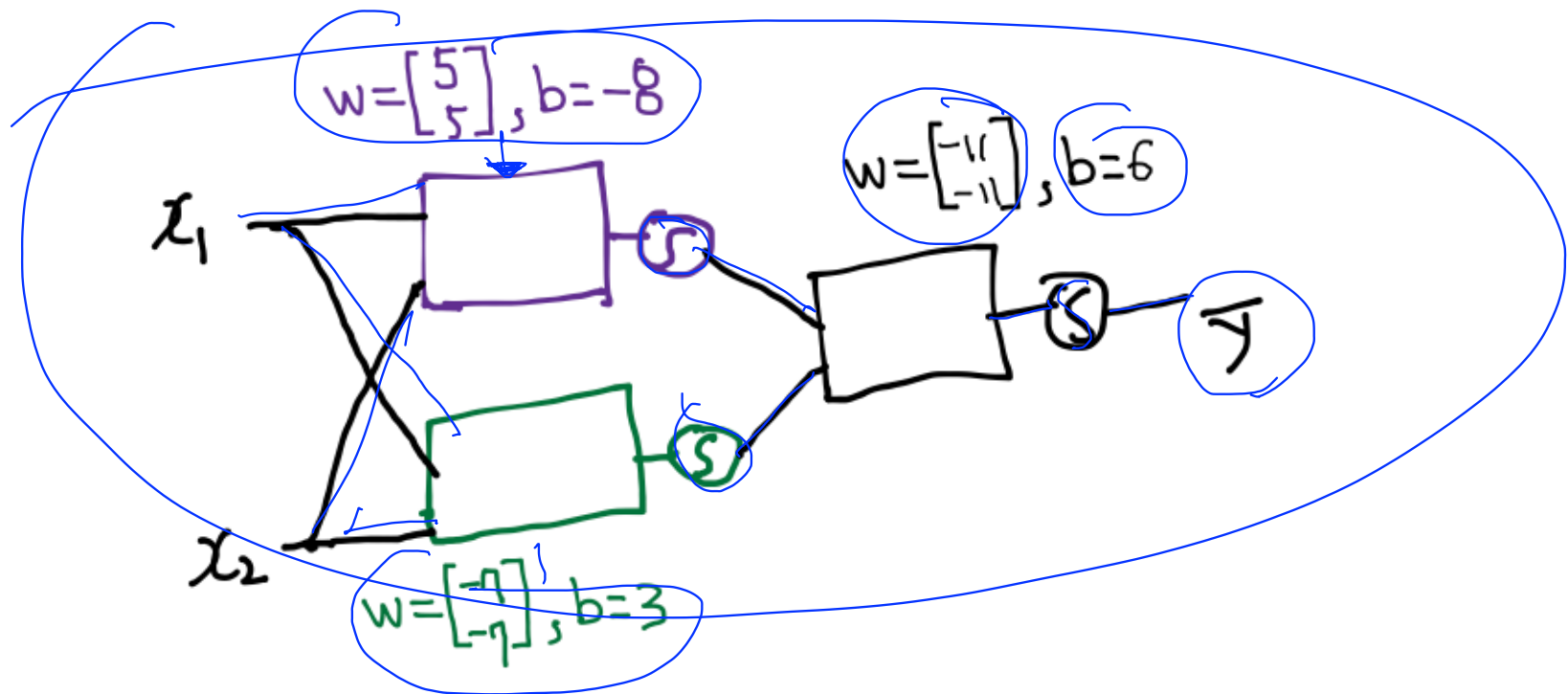
$$[1 \ 1] \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 3 = -1 + -1 + 3 = \underline{-1}, \text{Sigmoid}(\underline{-1}) = \underline{0}$$

$$[1 \ 0] \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 6 = \underline{-1} + 0 + \underline{6} = \underline{-5}$$

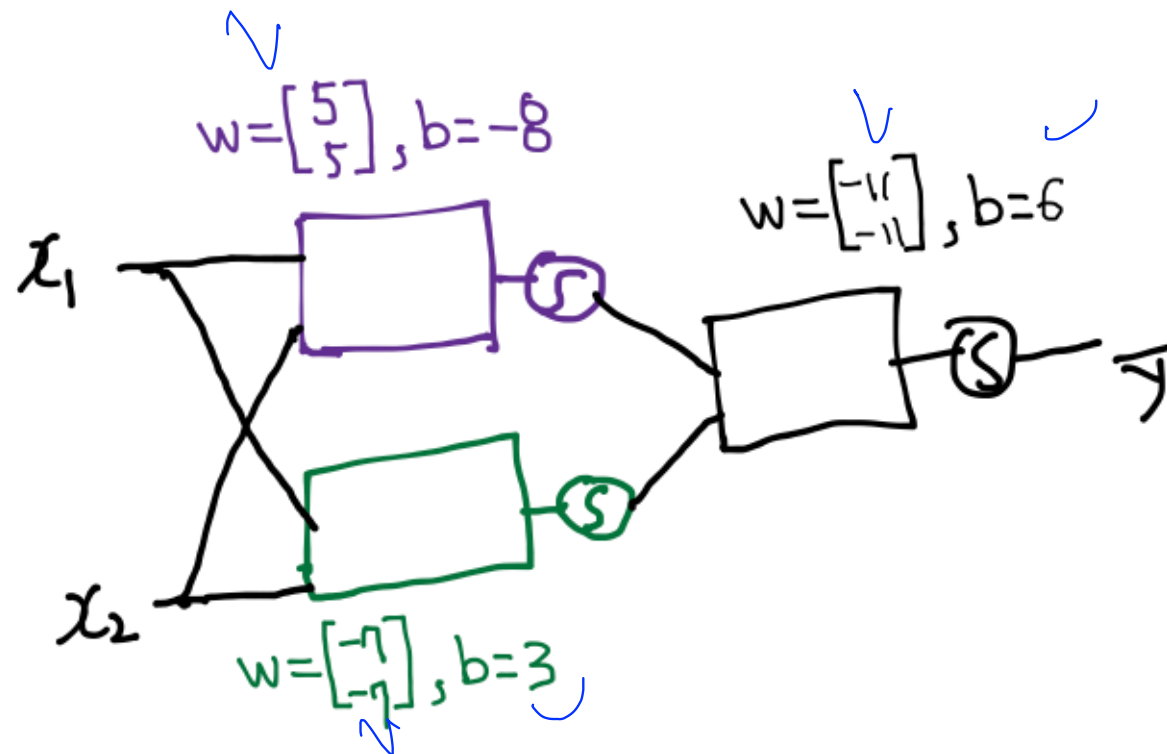
$\text{Sigmoid}(-5) = 0$

x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0 ✓
0	1	0	0	1	1 ✓
1	0	0	0	1	1 ✓
✓ 1	1	1	0	0	0 ✓

Forward propagation

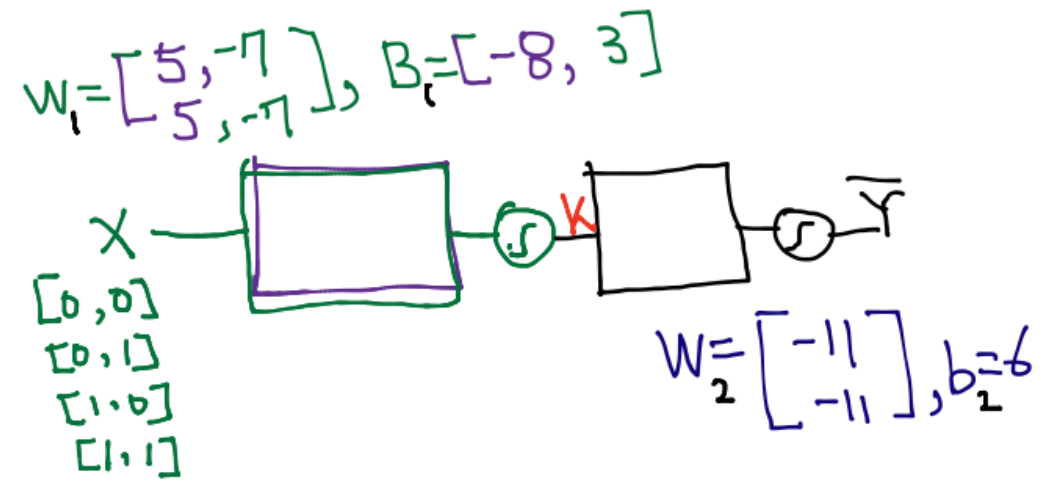
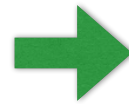
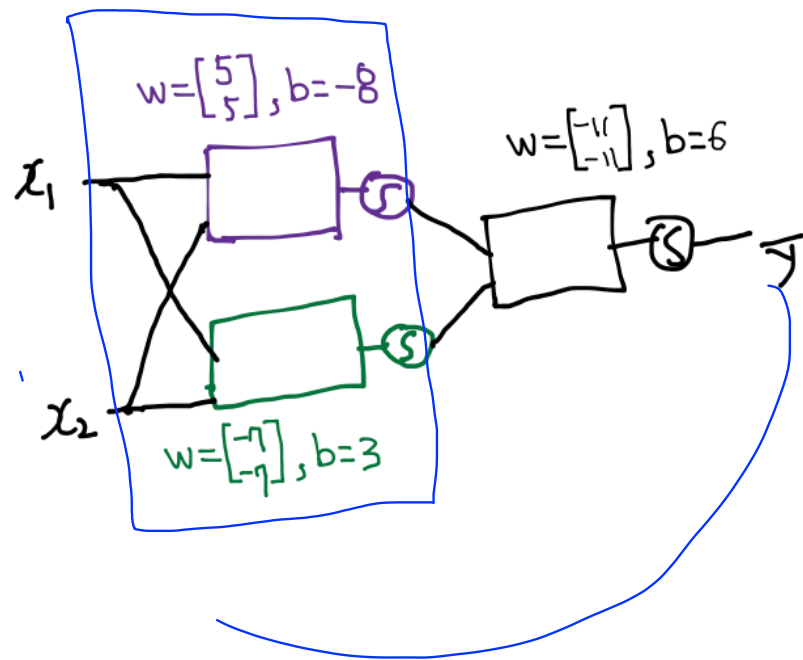


Forward propagation



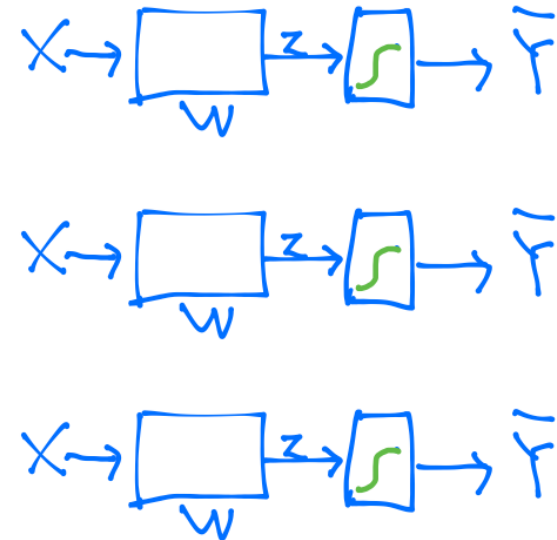
Can you find another W and b for the XOR?

NN

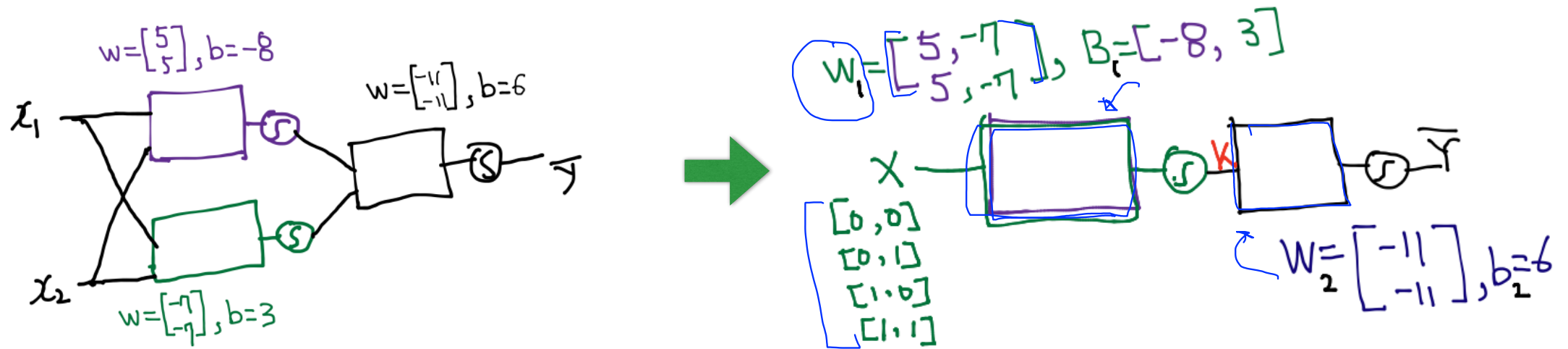


Recap: Lec 6-1 Multinomial classification

$$\begin{bmatrix} w_{A1} & w_{A2} & w_{A3} \\ w_{B1} & w_{B2} & w_{B3} \\ w_{C1} & w_{C2} & w_{C3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{A1}x_1 + w_{A2}x_2 + w_{A3}x_3 \\ w_{B1}x_1 + w_{B2}x_2 + w_{B3}x_3 \\ w_{C1}x_1 + w_{C2}x_2 + w_{C3}x_3 \end{bmatrix} = \begin{bmatrix} \bar{y}_A \\ \bar{y}_B \\ \bar{y}_C \end{bmatrix}$$



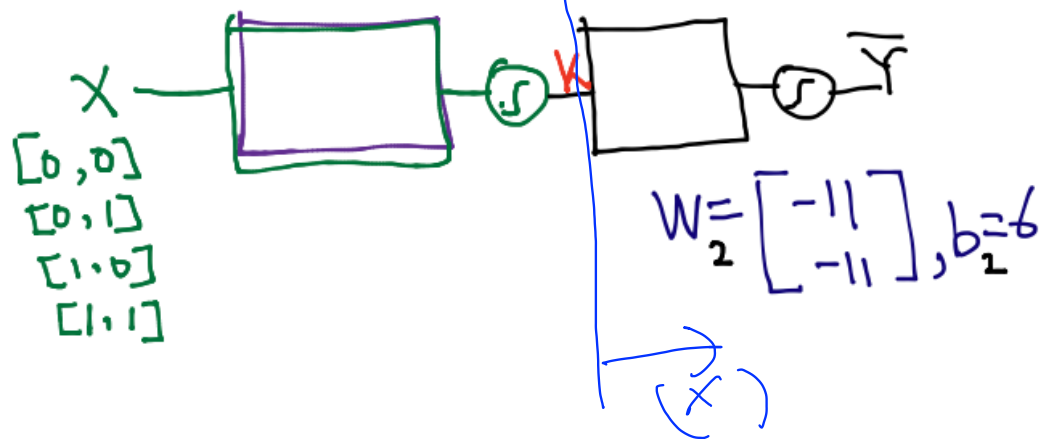
NN



How can we learn W , and b from trading data?

NN

$$W_1 = \begin{bmatrix} 5 & -7 \\ 5 & -7 \end{bmatrix}, B_1 = [-8, 3]$$



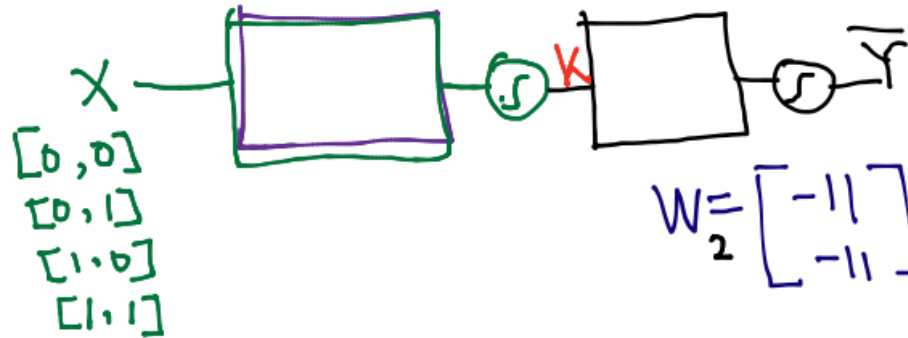
$$W_2 = \begin{bmatrix} -11 \\ -11 \end{bmatrix}, b_2 = 6$$

$$K(X) = \text{sigmoid}(XW_1 + B_1)$$

$$\hat{Y} = H(X) = \text{sigmoid}(K(X)W_2 + b_2)$$

NN

$$W_1 = \begin{bmatrix} 5 & -7 \\ 5 & -7 \end{bmatrix}, B_1 = [-8, 3]$$



$$W_2 = \begin{bmatrix} -11 \\ -11 \end{bmatrix}, b_2 = 6$$

$$K(X) = \text{Sigmoid}(XW_1 + B_1)$$

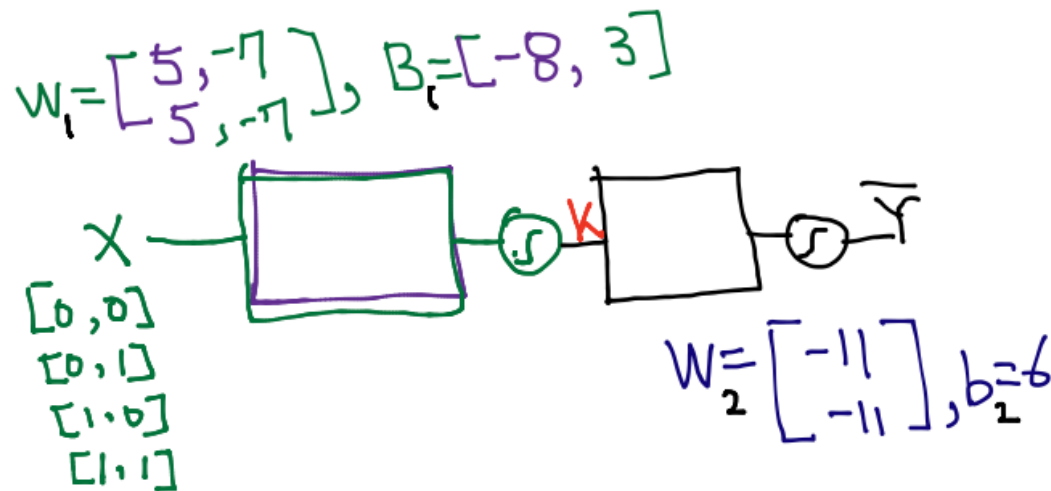
$$\underline{\underline{Y}} = H(X) = \text{Sigmoid}(K(X)W_2 + b_2)$$

#NN

`K = tf.sigmoid(tf.matmul(X, W1) + b1)`

`hypothesis = tf.sigmoid(tf.matmul(K, W2) + b2)`

NN



$$K(X) = \text{sigmoid}(XW_1 + B_1)$$

$$\hat{Y} = H(X) = \text{sigmoid}(K(X)W_2 + b_2)$$

NN

`K = tf.sigmoid(tf.matmul(X, W1) + b1)`

`hypothesis = tf.sigmoid(tf.matmul(K, W2) + b2)`

How can we learn $W1, W2, B1, b2$ from training data?

Next
Backpropagation

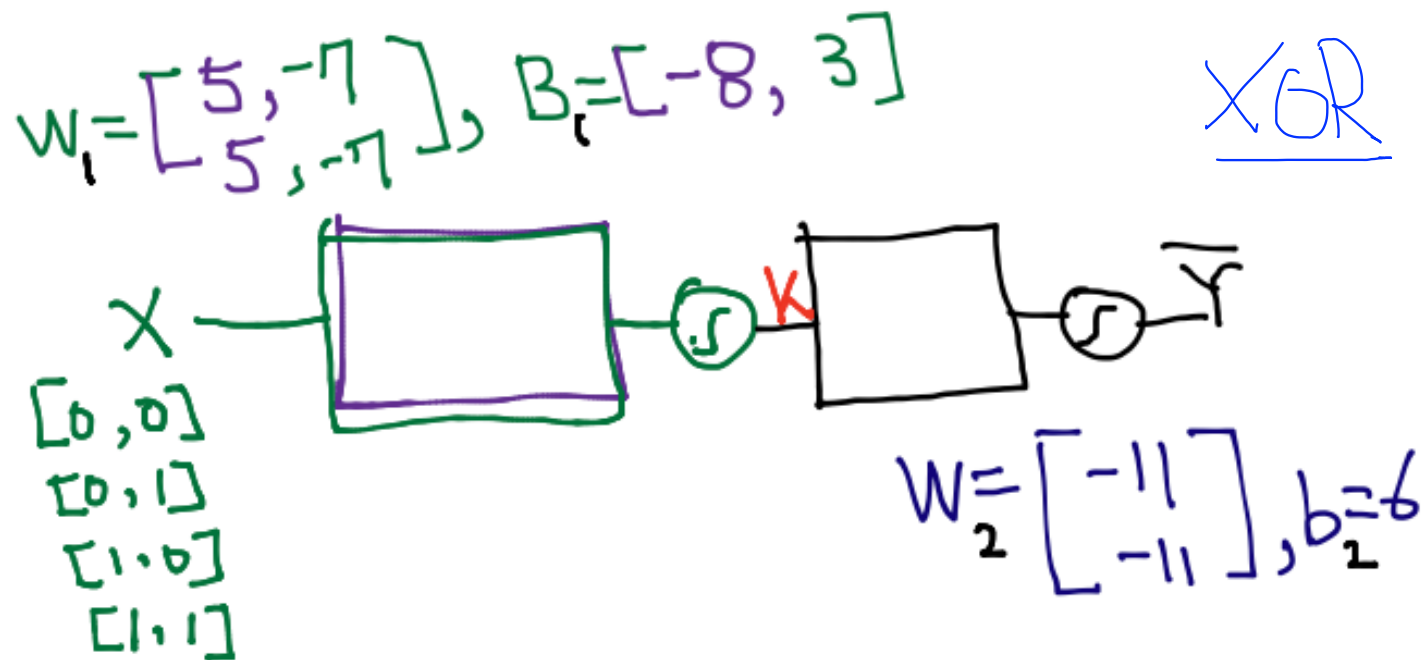


Lecture 9-2

Backpropagation

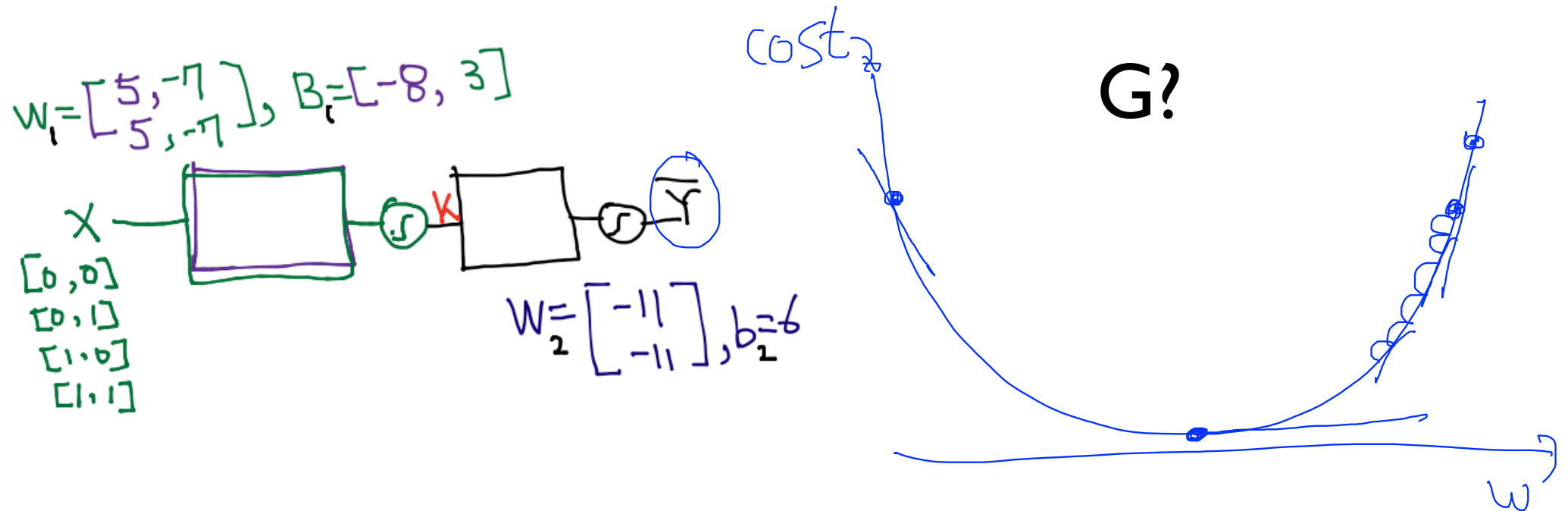
Sung Kim <hunkim+mr@gmail.com>

Neural Network (NN)

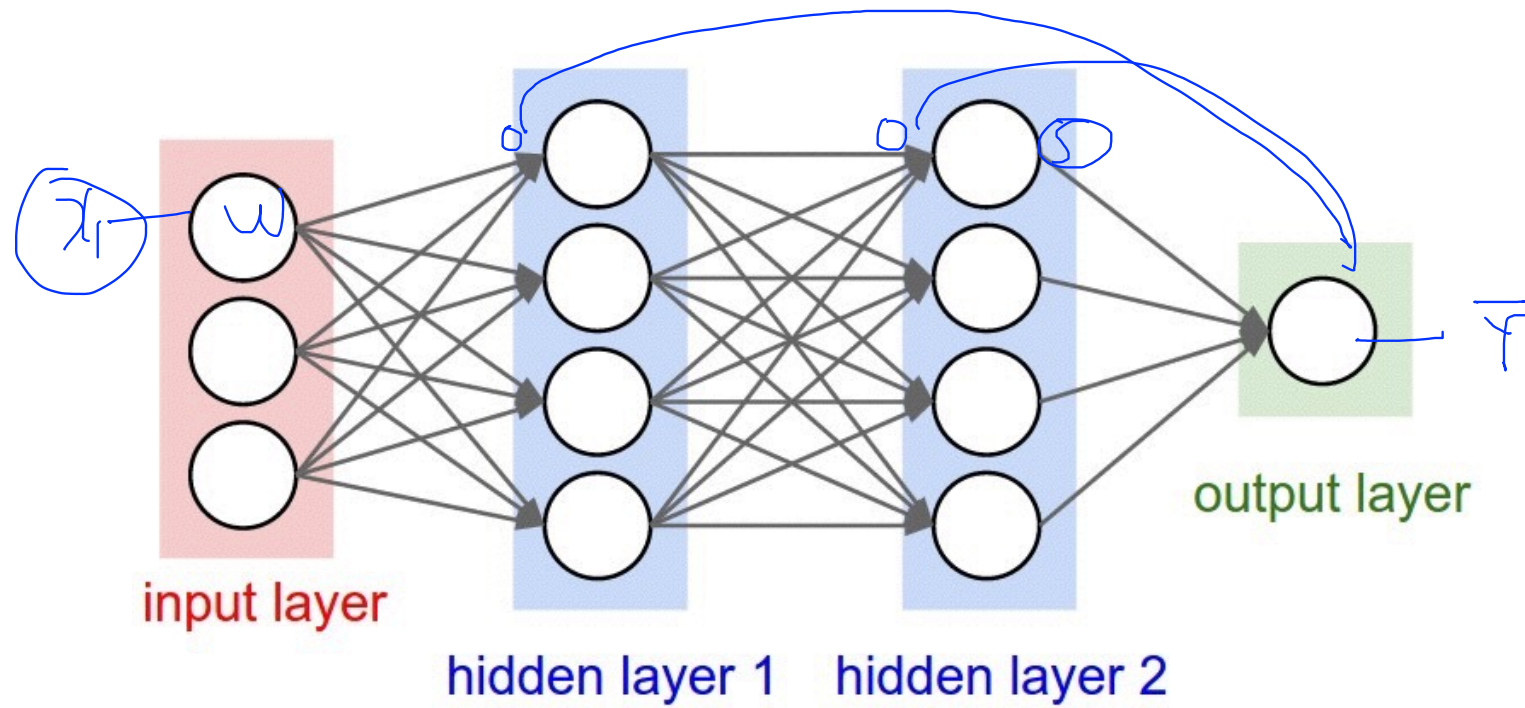


How can we learn W_1, W_2, B_1, b_2 from training data?

How can we learn $W1, W2, B1, b2$ from training data?

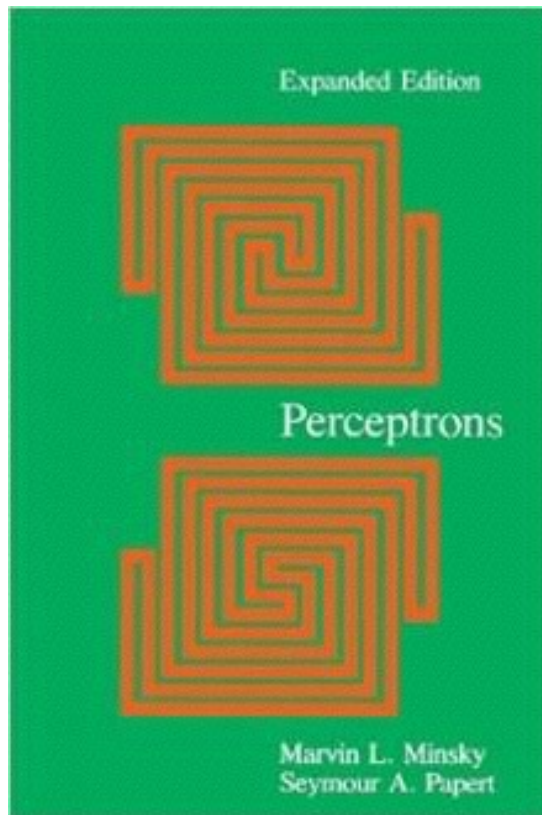


Derivation



Perceptrons (1969)

by Marvin Minsky, founder of the MIT AI Lab

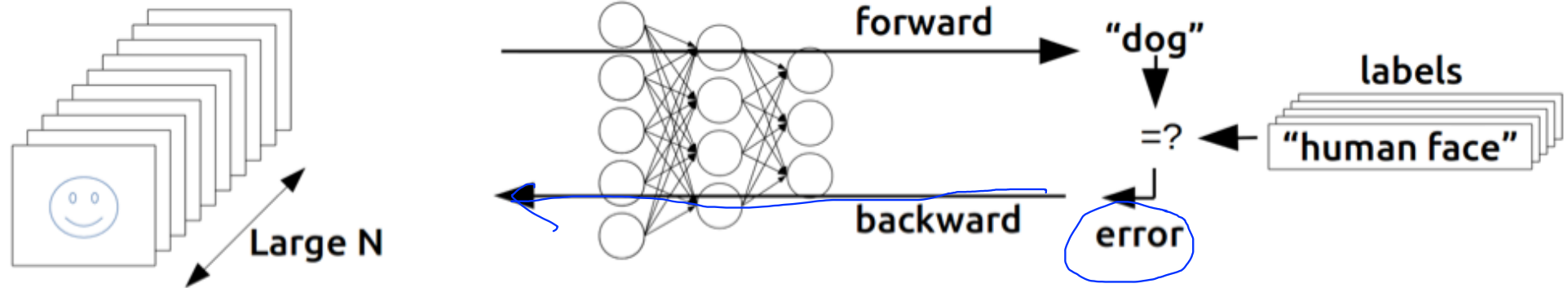


- We need to use MLP, multilayer perceptrons (multilayer neural nets)
- No one on earth had found a viable way to train MLPs good enough to learn such simple functions.

Backpropagation

(1974, 1982 by Paul Werbos, 1986 by Hinton)

Training



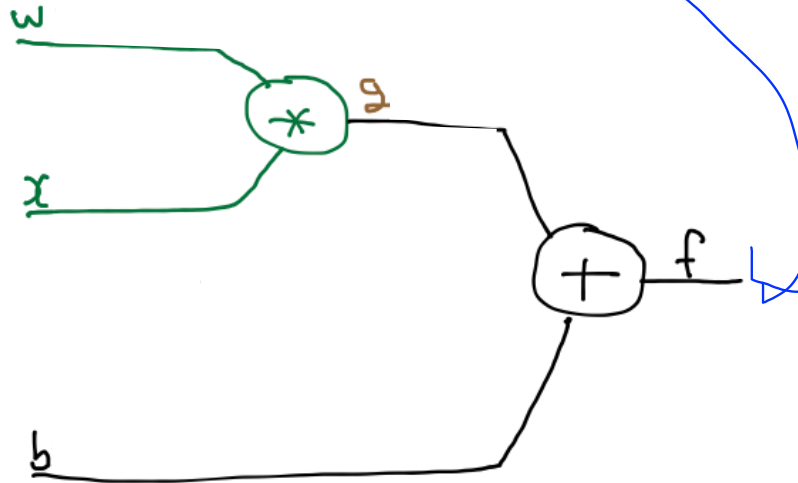
<https://devblogs.nvidia.com/parallelforall/inference-next-step-gpu-accelerated-deep-learning/>

Back propagation (chain rule)

$$\underline{f} = \underbrace{wx + b}_{g}, \quad \underline{g} = wx, \quad \underline{f} = \underline{g} + b$$

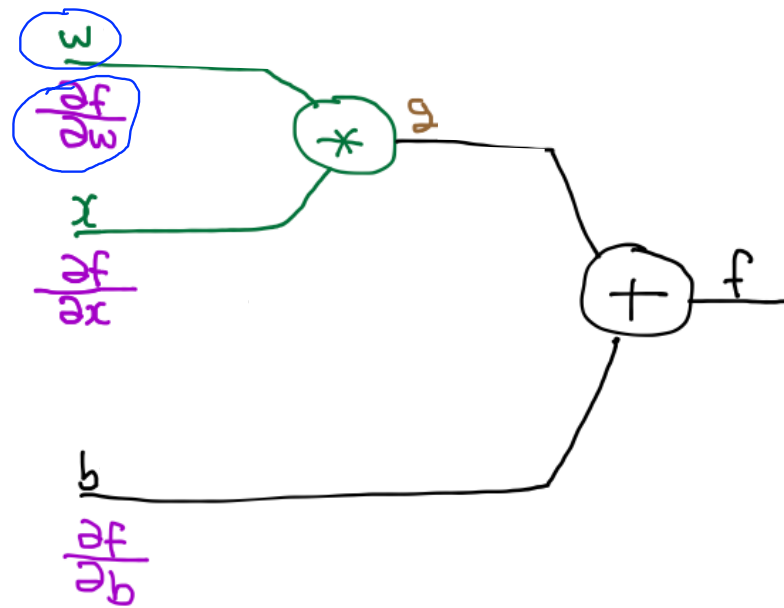
Back propagation (chain rule)

$$\underline{f = wx + b}, \quad g = wx, \quad f = g + b$$



Back propagation (chain rule)

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



Basic derivative

$$\frac{d}{dx}f(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

$$f(x) = 3$$

$$f(x) = x$$

$$f(x) = 2x$$

<https://ko.wikipedia.org/wiki/%EB%AF%B8%EB%B6%84>

Partial derivative: consider other variables as constants

$$f(x) = 2x$$

$$f(x, y) = xy, \frac{\partial f}{\partial x}$$

$$f(x, y) = xy, \frac{\partial f}{\partial y}$$

<https://ko.wikipedia.org/wiki/%EB%AF%B8%EB%B6%84>

Partial derivative: consider other variables as constants

$$f(x) = 3$$

$$f(x) = 2x \quad f(x) = x + x$$

$$f(x) = x + 3$$

$$f(x, y) = x + y, \frac{\partial f}{\partial x}$$

$$f(x, y) = x + y, \frac{\partial f}{\partial y}$$

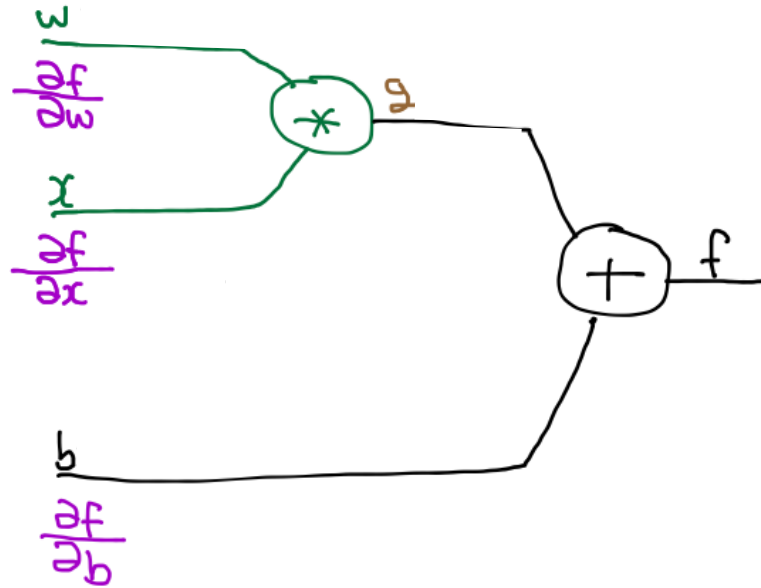
Chain rule

$f(g(x))$

$$\boxed{\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \quad \star}$$

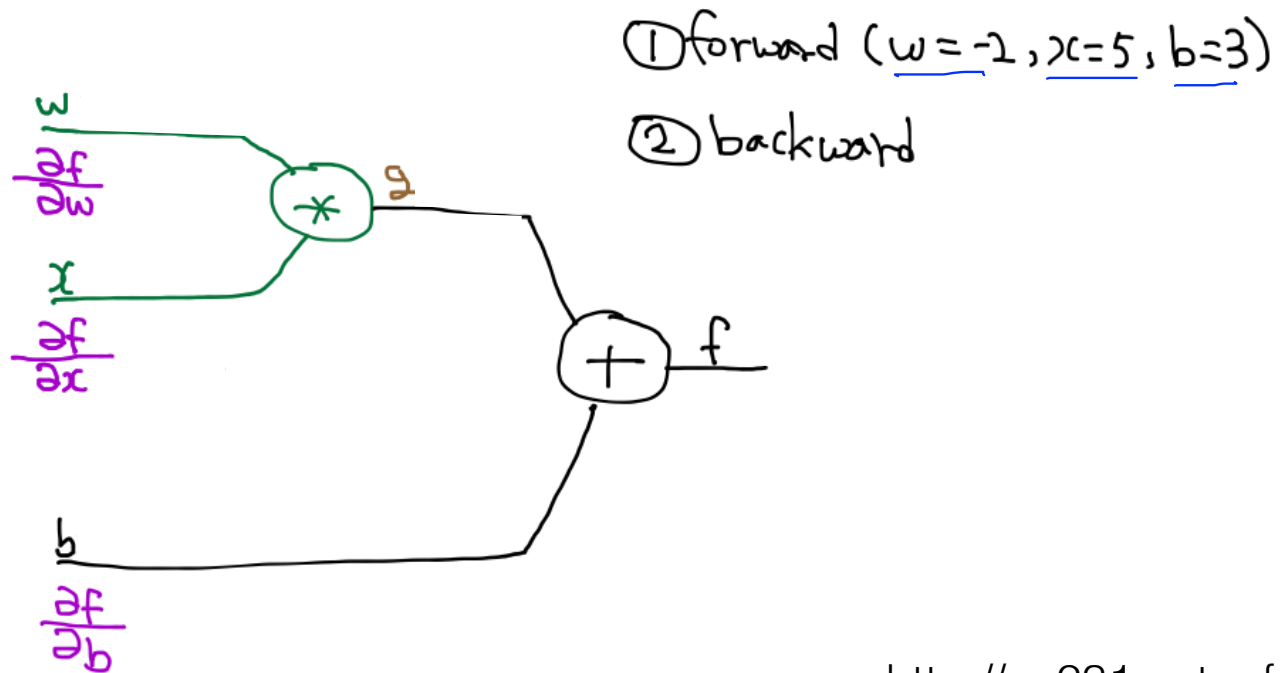
Back propagation (chain rule)

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



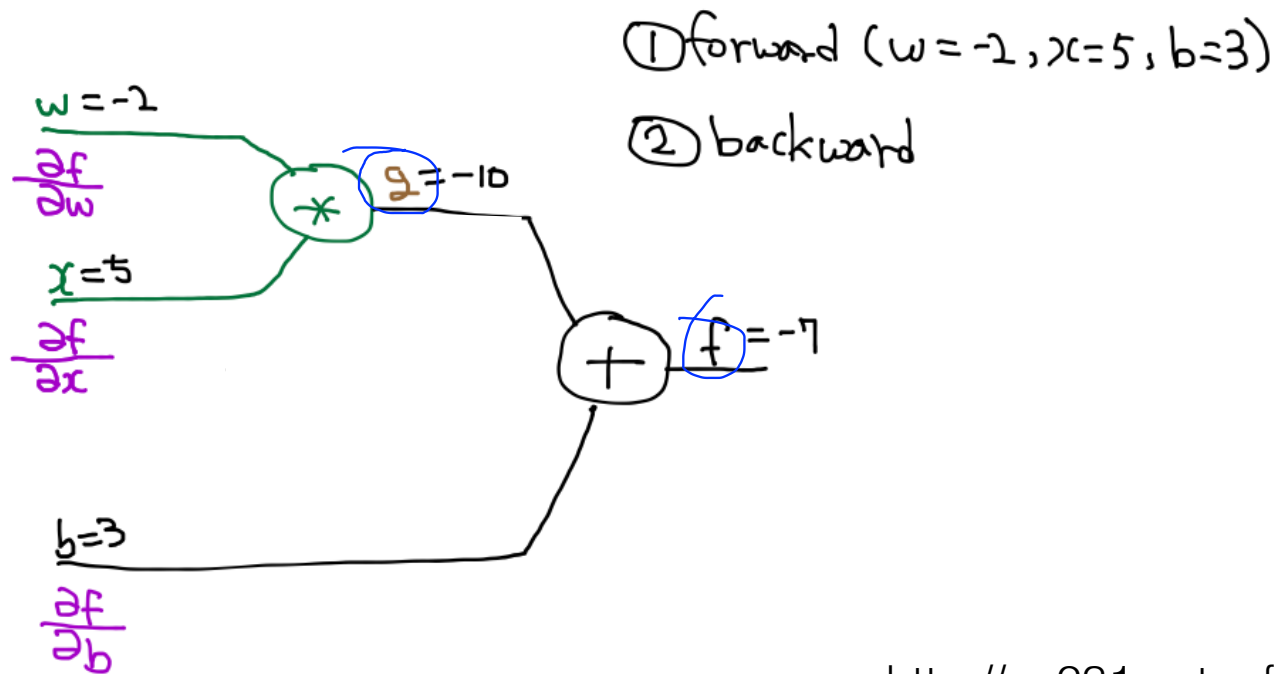
Back propagation (chain rule)

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



Back propagation (chain rule)

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



Back propagation (chain rule)

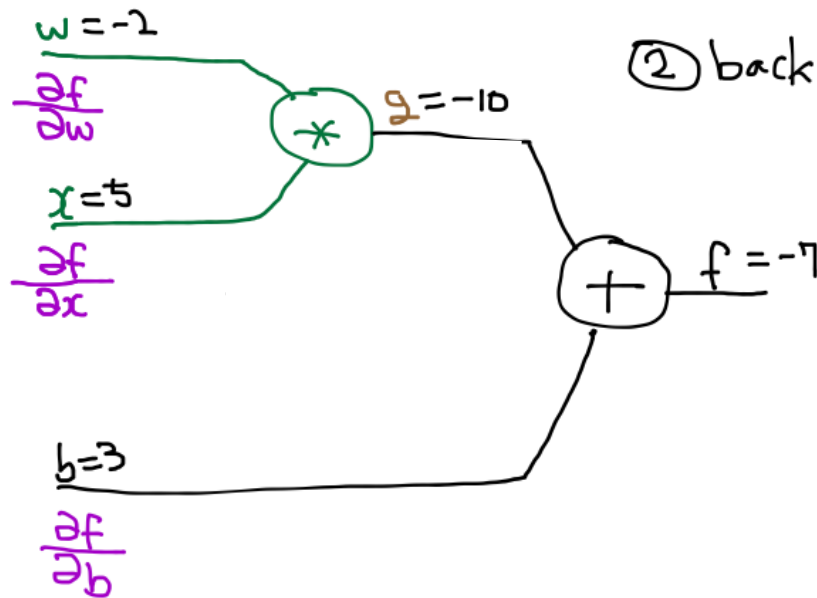
$$f = wx + b, \quad g = wx, \quad f = g + b$$

$$\left(\frac{\partial g}{\partial w} = x, \quad \frac{\partial g}{\partial x} = w \right)$$

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

① forward ($w = -2, x = 5, b = 3$)

② backward



Back propagation (chain rule)

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial w} = 1 * 5 = 5$$

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

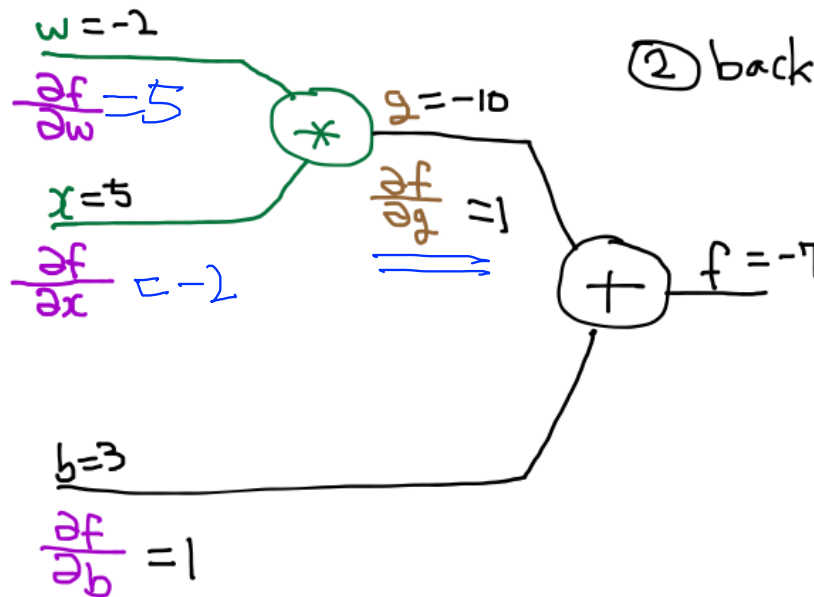
$$f = wx + b, \quad g = wx, \quad f = g + b$$

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

$\hookrightarrow \frac{\partial g}{\partial w} = x, \quad \frac{\partial g}{\partial x} = w$

① forward ($w = -2, x = 5, b = 3$)

② backward



Back propagation (chain rule)

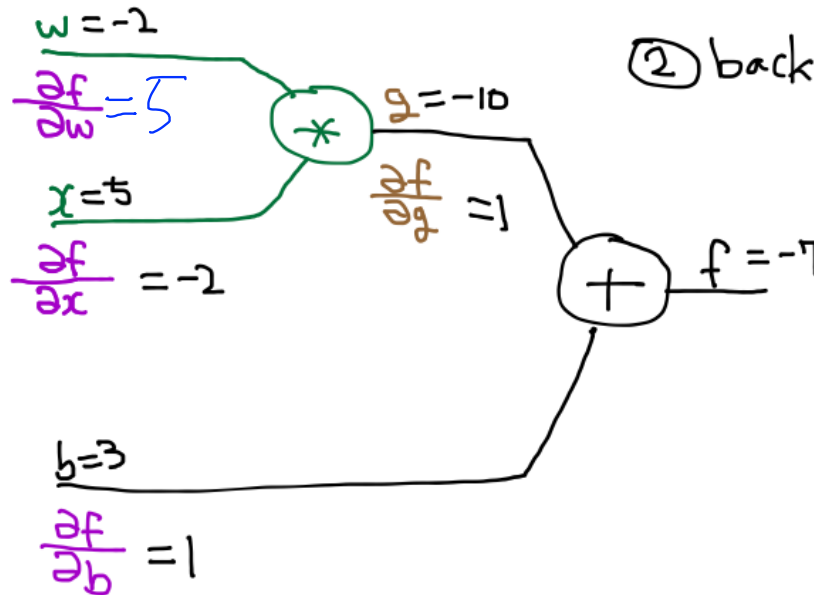
$$f = wx + b, \quad g = wx, \quad f = g + b$$

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

$\hookrightarrow \frac{\partial g}{\partial w} = x, \frac{\partial g}{\partial x} = w$

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

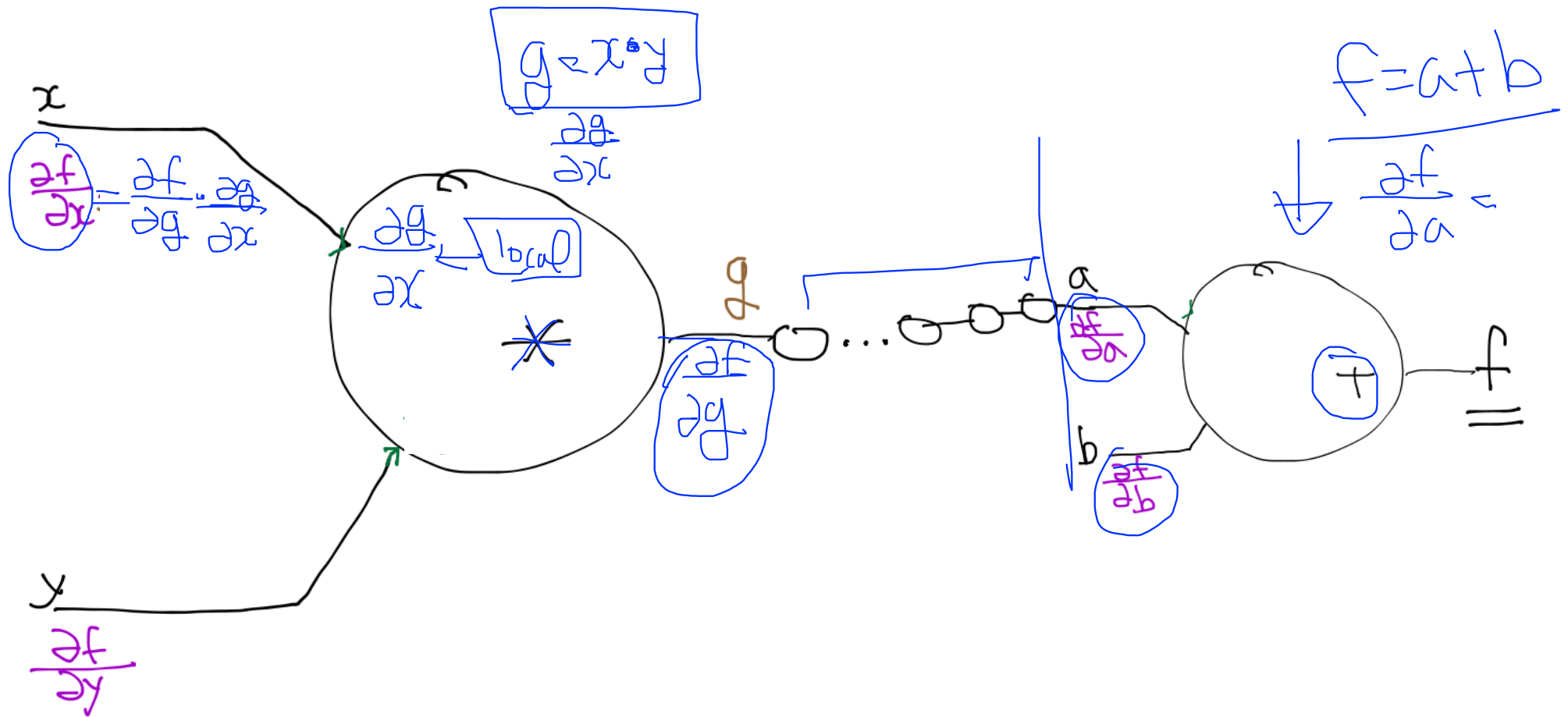
$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w} = 1 * x = 5$$



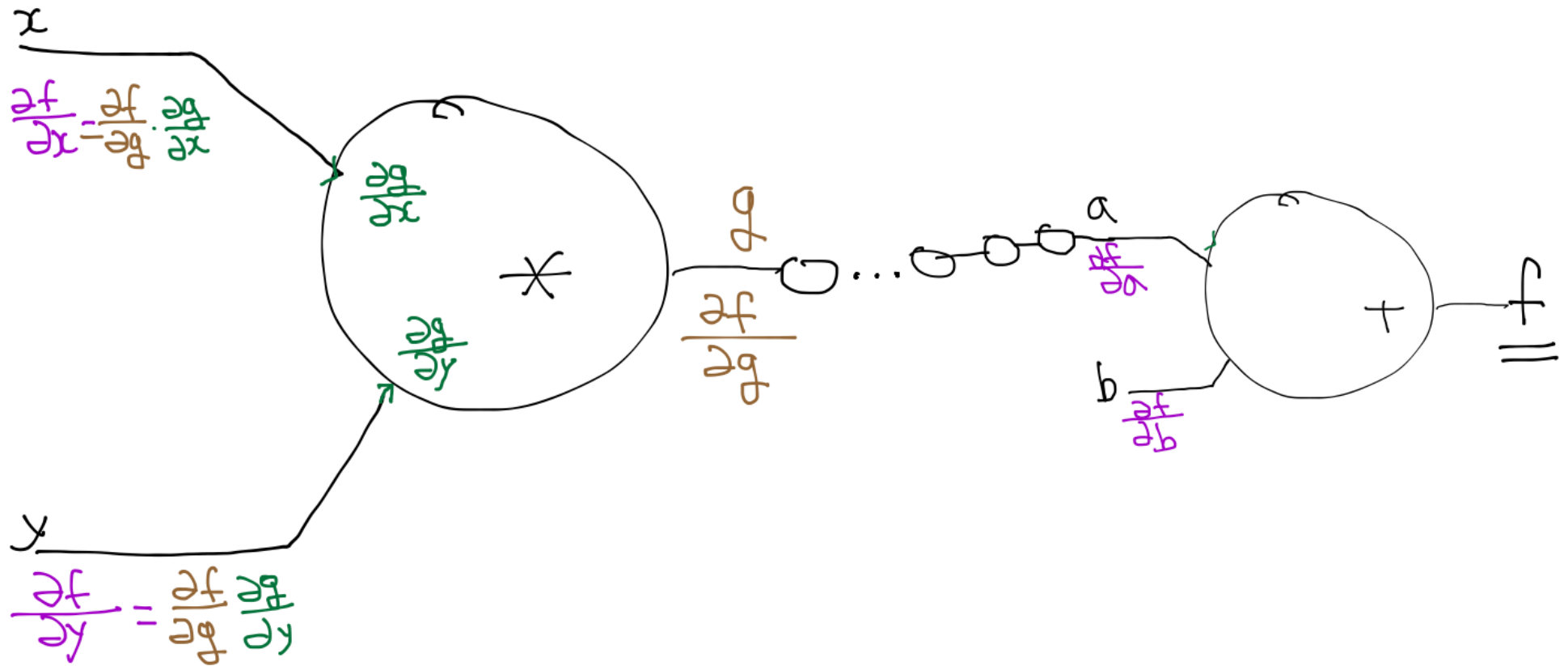
① forward ($w = -2, x = 5, b = 3$)

② backward

Back propagation (chain rule)



Back propagation (chain rule)

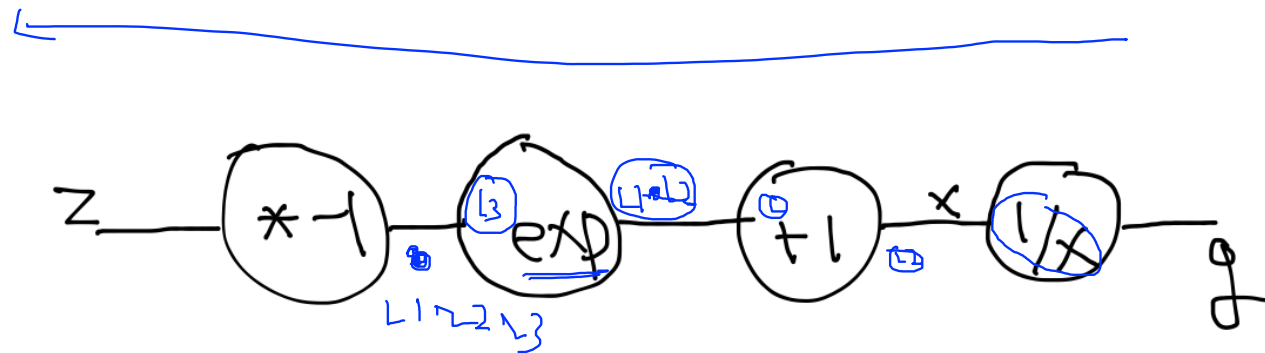


Sigmoid

$$g(z) = \frac{1}{1+e^{-z}} \quad \frac{\partial g}{\partial z}$$

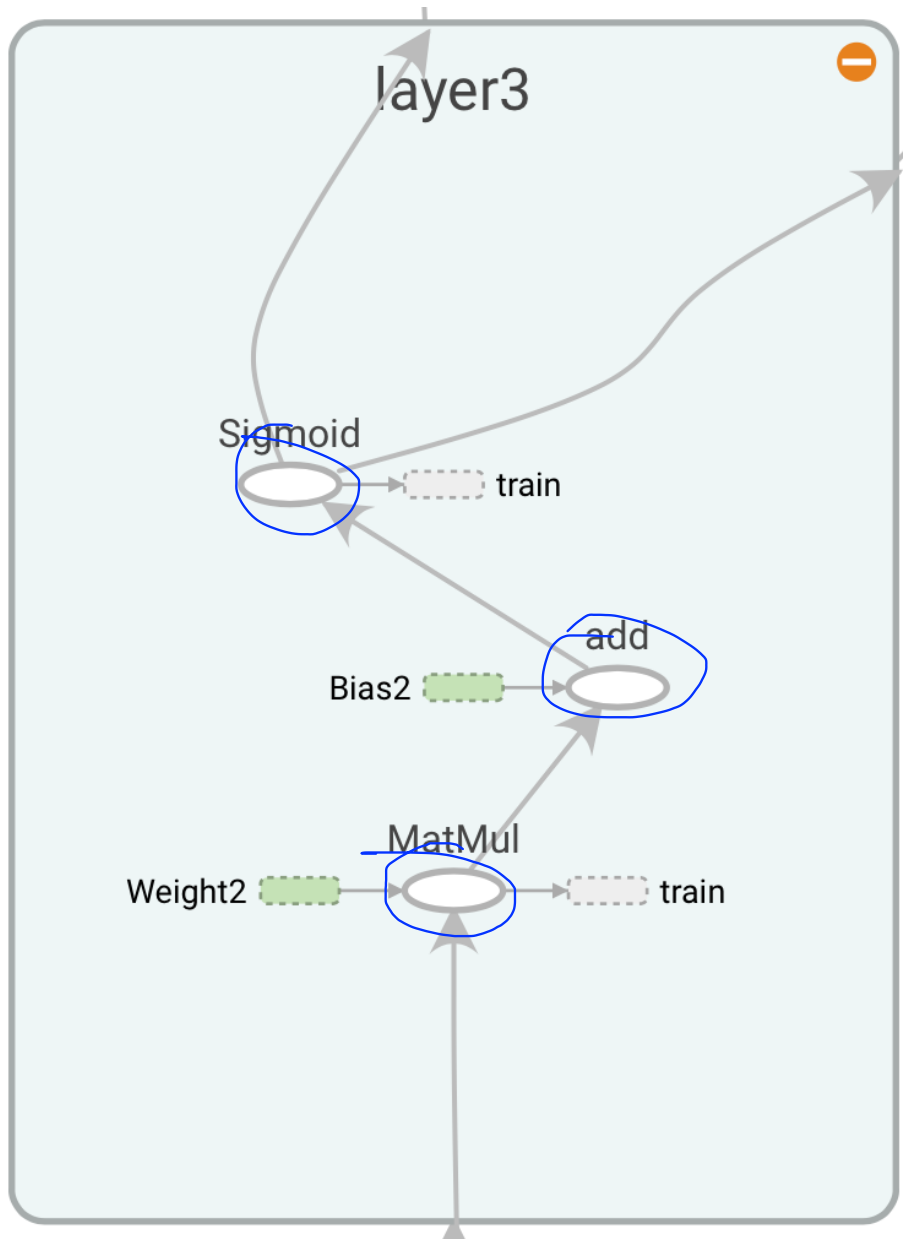
Sigmoid

$$g(z) = \frac{1}{1 + e^{-z}}$$



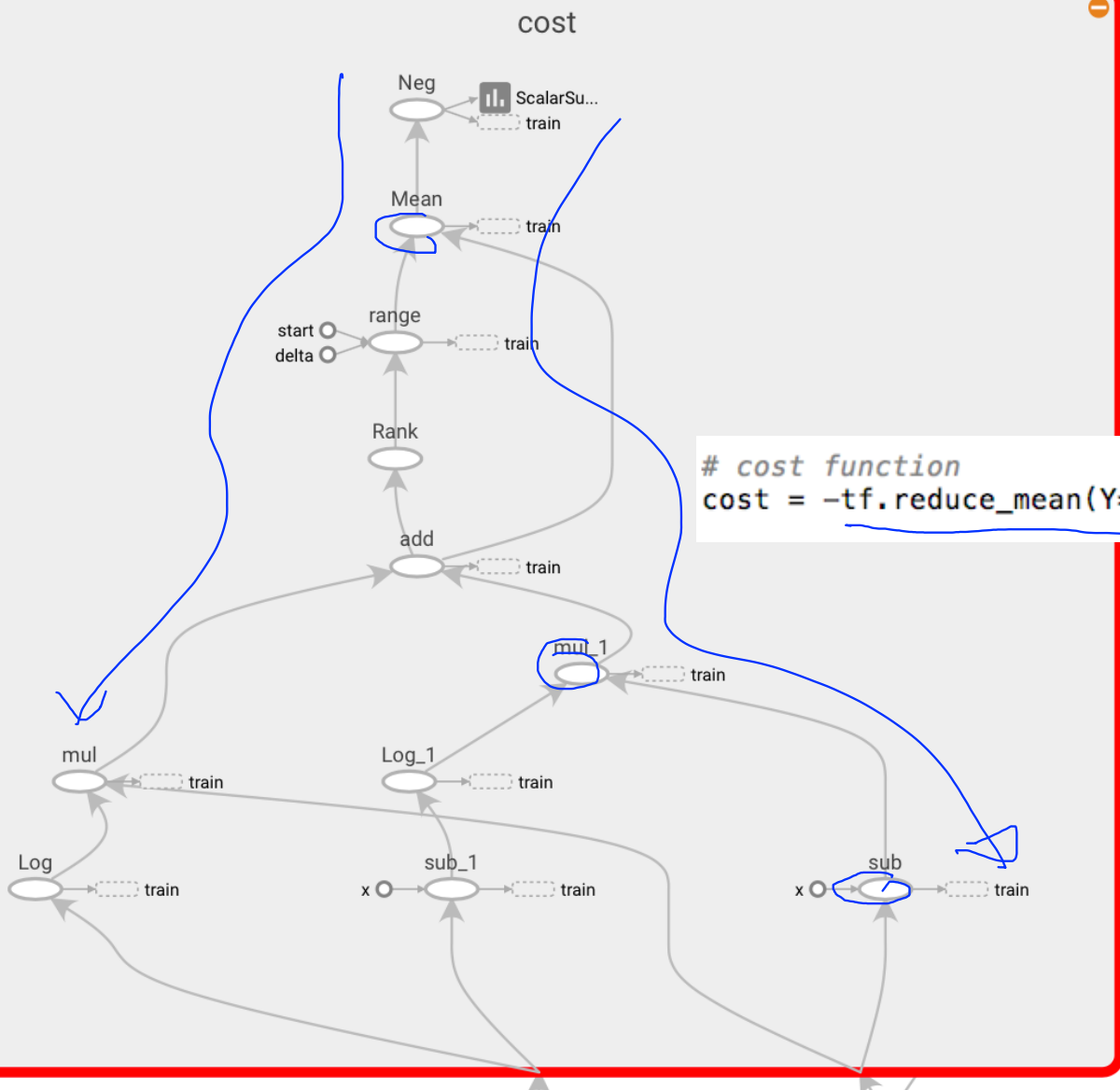
Back propagation in TensorFlow TensorBoard

`hypothesis = tf.sigmoid(tf.matmul(L2, W2) + b2)`

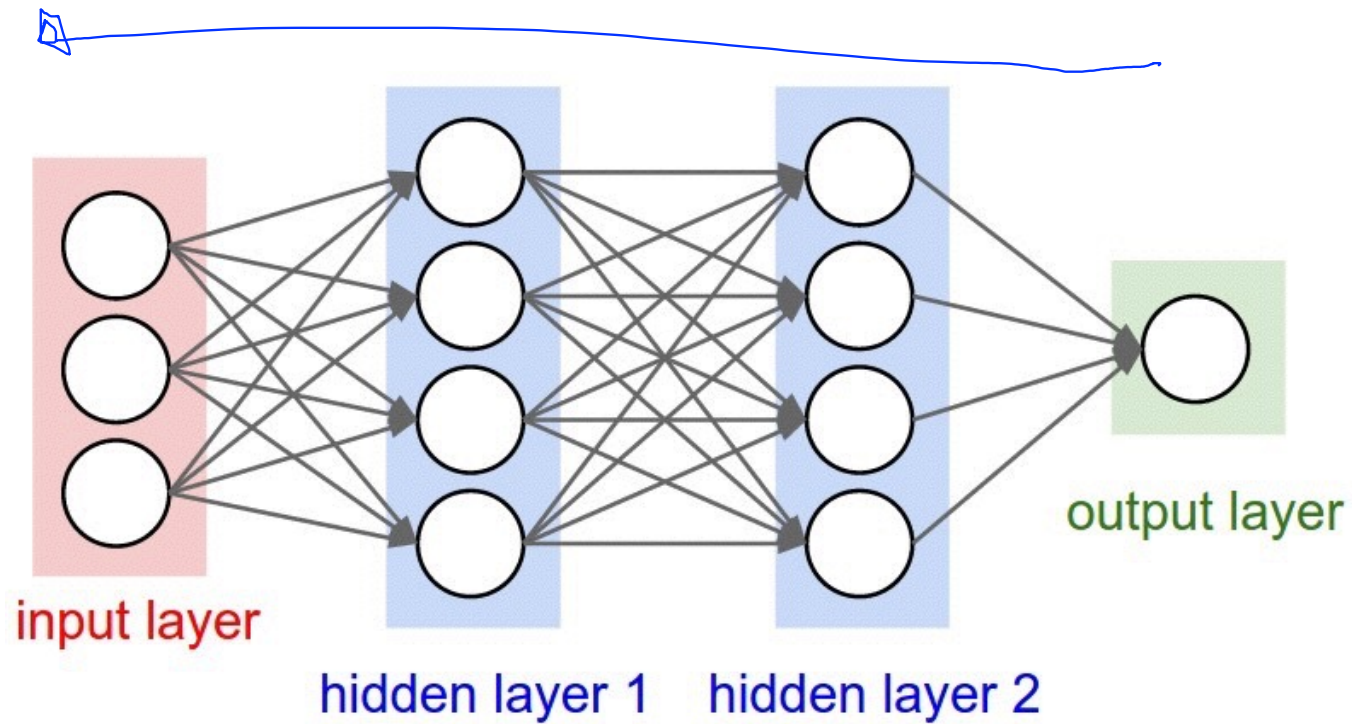


Back propagation in TensorFlow

TensorBoard

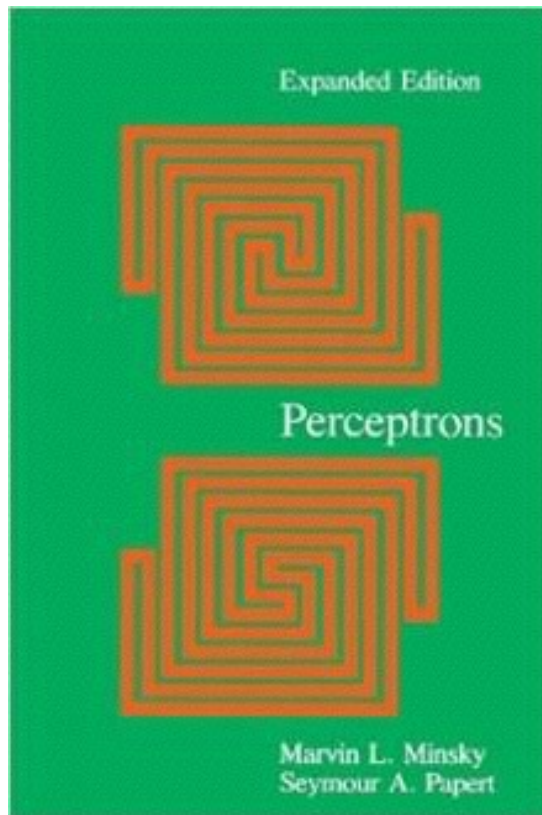


Back propagation



Perceptrons (1969)

by Marvin Minsky, founder of the MIT AI Lab



- We need to use MLP, multilayer perceptrons (multilayer neural nets)
- No one on earth had found a viable way to train MLPs good enough to learn such simple functions.

Next
ReLU

