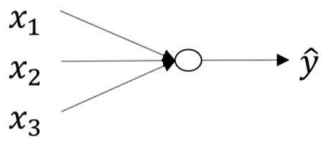
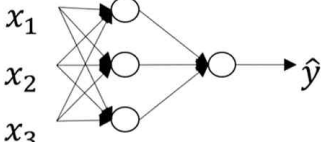
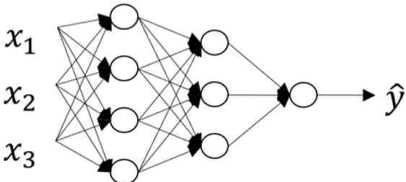
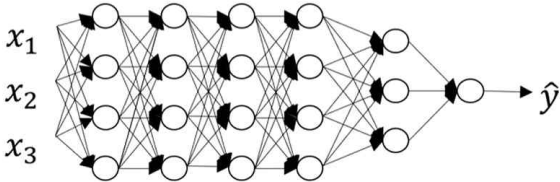
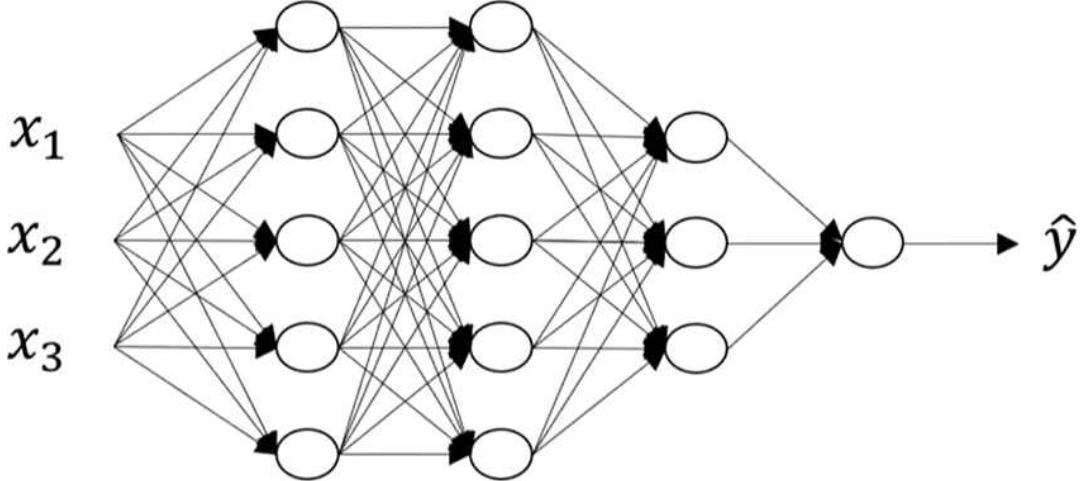


Deep Neural Networks

Deep L-layer Neural Networks

	
logistic regression	1 hidden layer
	
2 hidden layers	5 hidden layers

	
$L = \text{\#layers}$ $n^{[l]} = \text{\#units in layer } l$ $a^{[l]} = \text{activation value in layer } l$ $W^{[l]}, b^{[l]} = \text{weights for } z^{[l]} \text{ in layer } l$ $a^{[l]} = g^{[l]}(z^{[l]})$	$n^{[0]} = n_x = 3,$ $n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3,$ $n^{[4]} = n^{[L]} = 1$

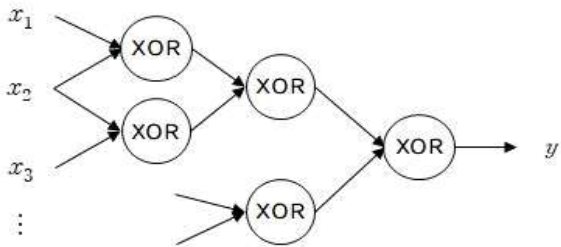
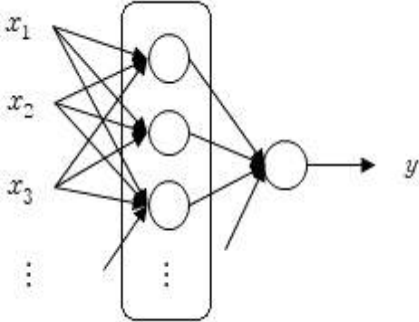
Forward Propagation in a Deep Network

$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$ $A^{[l]} = g^{[l]}(Z^{[l]})$	$Z^{[1]} = W^{[1]} X + b^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$ $A^{[1]} = g^{[1]}(Z^{[1]})$ $Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$ $A^{[2]} = g^{[2]}(Z^{[2]}) = \sigma(Z^{[2]})$ \vdots $\hat{Y} = g^{[4]}(Z^{[4]}) = A^{[4]}$
--	---

Getting matrix dimensions right

Parameters $W^{[l]}$ and $b^{[l]}$
$W^{[l]} : (n^{[l]}, n^{[l-1]})$ $b^{[l]} : (n^{[l]}, 1)$ $dW^{[l]} : (n^{[l]}, n^{[l-1]})$ $db^{[l]} : (n^{[l]}, 1)$
$l = 0, A^{[0]} = X : (n^{[0]}, m)$ $A^{[l-1]} : (n^{[l-1]}, m)$ $Z^{[l]} = A^{[l]} : (n^{[l]}, m)$ $dZ^{[l]} = dA^{[l]} : (n^{[l]}, m)$
$Z^{[l]}_{(n^{[l]}, m)} = W^{[l]}_{(n^{[l]}, n^{[l-1]})} A^{[l-1]}_{(n^{[l-1]}, m)} + b^{[l]}_{(n^{[l]}, 1) \rightarrow (n^{[l]}, m)}$

Why deep representations?

circuit theory and deep learning	
Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.	
$y = x_1 \text{ xor } x_2 \text{ xor } \dots \text{ xor } x_n$  <p>$\Rightarrow O(\log n)$</p>	 <p>$\Rightarrow O(2^n) = 2^{n-1}$</p>
The deeper layers of a neural network are typically computing more complex features of the input than the earlier layers.	
(i) To compute the function using a shallow network circuit, you will need a large network (where we measure size by the number of logic gates in the network), but (ii) To compute it using a deep network circuit, you need only an exponentially smaller network.	

"Applied deep learning is a very empirical process."

최근에는 데이터가 새로 들어오면서 가장 좋은 hyperparameter 값이 쉽게 변할 수 있다. 오늘 가장 좋았던 learning rate 값이 한 달 뒤 데이터가 변하거나 새로 추가되면서 가장 좋은 learning rate 값이 변할 수 있다. 그 때문에 경험적 과정인 것이다.

Building blocks of deep neural networks

Forward and Backward functions

<p>layer l: $W^{[l]}, b^{[l]}$</p> <p>Forward:</p> <p>input $A^{[l-1]}$, output $A^{[l]}$</p> <p>$Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]} \rightarrow \text{cache } Z^{[l]}$</p> <p>$A^{[l]} = g^{[l]}(Z^{[l]})$</p>	<p>$da^{[0]}$은 어느 곳에도 사용되지 않기 때문에 구하지 않아도 괜찮다.</p>
<p>Backward:</p> <p>input $da^{[l]}$, cache $Z^{[l]}$</p> <p>output $da^{[l-1]}, dw^{[l]}, db^{[l]}$</p>	

layer l

$A^{[0]} \dots A^{[l-1]} \rightarrow \begin{matrix} W^{[l]}, b^{[l]} \\ g^{[l]}(Z^{[l]}) \end{matrix} \rightarrow A^{[l]} \dots A^{[L]} = \hat{y}$

$\text{cache } Z^{[l]}$

$\dots da^{[l-1]} \leftarrow \begin{matrix} W^{[l]}, b^{[l]} \\ dZ^{[l]} \end{matrix} \leftarrow da^{[l]} \dots da^{[L]} = \frac{\partial L(\hat{Y}, Y)}{\partial A}$

$dW^{[l]}, db^{[l]}$

$L(\hat{Y}, Y)$

Parameters vs Hyperparameters

Parameters: weights $W^{[l]}, b^{[l]}$

Hyperparameters: parameters를 결정하는 요소

learning rate

iterations

hidden layers

hidden units

choice of activation function

momentum, minibatch size, regularization