

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including our cookie policy.

X

# Lecture Notes: Neural Networks and Deep Learning

Learn from my notes for the Deep Learning course on DeepLearning.ai



Daniel Wu

Dec 13, 2019 · 12 min read ★



Photo by Margarida CSilva on Unsplash

The following are my lecture notes for the first of five courses in the Deep Learning Specialization on DeepLearning.ai:

## Week 1 – Introduction

Neural networks are a method to create predictions based on a typically large dataset of inputs with many characteristics called features. The network is made up of units of

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X

until the output node, which presents the outcome.

## Week 2 — Binary Classification

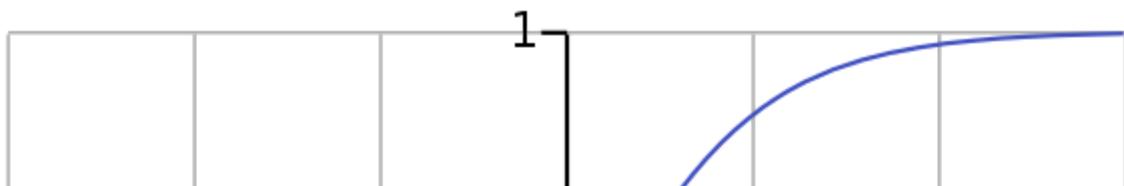
Using logistic regression as an example for neural networks, an image is fed into the network. For example, a color image of a cat is  $20 \times 20$  pixels so it has the image in red, blue, and green channels. That means there are 400 pixels of red, blue, and green each, resulting in 1200 pixels. These 1200 pixels are the number of 'n' features. Input data is unrolled into a column in a matrix of the shape  $(1200,1)$ . If you have a data set of 100 cat images, then the matrix is  $(1200,100)$ . Output data in Y matrix is of shape  $(1200,1)$  since a logistic regression outputs 1 or 0.

Notations:

- Single training example pair is lower case  $(x,y)$
- For a specific training example 'm' is  $\{ (x(i),y(i) \dots x(m),y(m) \}$
- Matrix X has  $x(i)$  to  $x(m)$  stacked in columns
- X shape is  $(n, m)$
- Y shape is  $(1,m)$
- Since I'm not using LaTex or similar mathematical formatting, (i) part of  $x(i)$  or any other variable refers to a superscript of i-th example.

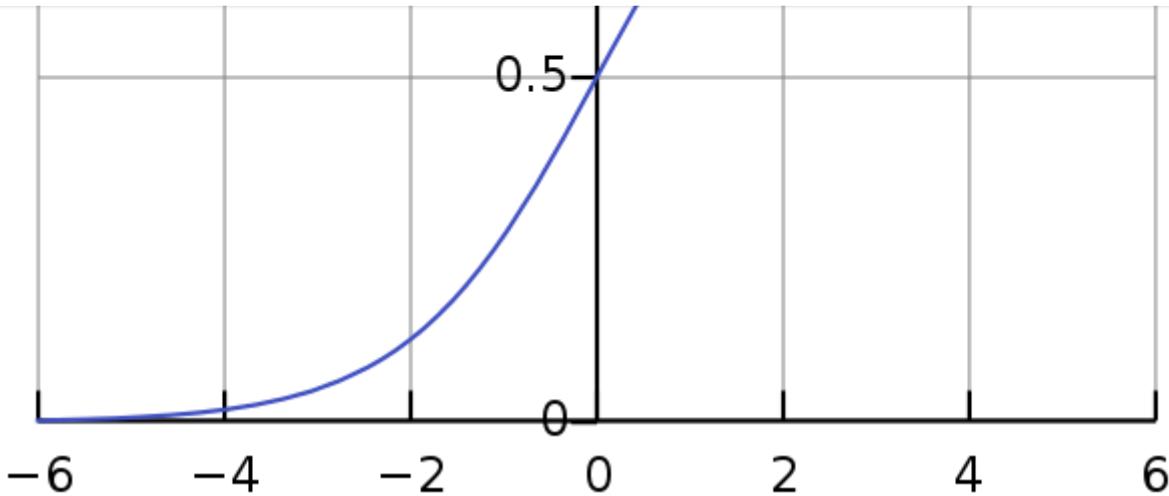
Unlike linear regression, a logistic regression applies a sigmoid function so your output is binary (e.g., 0/1). We're training the parameters for 'w' (weight) and 'b' (bias) so the output 'y' best estimated to be 1. We're learning to get the predicted output  $y\text{-hat}(i)$  approximate as close as possible to actual output  $y(i)$ .

- $z$  is a function =  $w^T x + b$  where  $T$  is represents the transpose of 'w' matrix
- sigmoid function of  $z = 1 / (1 + e^{-z})$



To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X



Sigmoid function outputs values between 0 and 1 (Source: wikipedia)

The loss (error) function measures how well our predicted  $\hat{y}(i)$  compares with actual  $y(i)$ .  $y(i)$  notation designates  $y$  at the  $i$ -th training example. In logistic regression, squared error doesn't do well with gradient descent so the error function  $L(\hat{y}, y) = (\hat{y} - y)^2 / 2$ . The output is convex and easier to optimize for gradient descent.

$$\text{Loss function} = -(y * \log(\hat{y}) + (1-y) * \log(1-\hat{y}))$$

When  $y = 1$ , then the loss function  $L(\hat{y}, y)$  is the first part:  $-\log(\hat{y})$  so that large  $\hat{y}$  also means large  $\log(\hat{y})$ .

When  $y = 0$ , then the loss function  $L(\hat{y}, y)$  is the latter part:  $-(\log(1-\hat{y}))$  so that large  $\log(1-\hat{y})$  means small  $\hat{y}$ .

## Cost Function:

Since above is defining a loss for a single training example, the cost function is the cost of your parameters based on the loss functions from all the training examples ( $i$  to  $m$ ) averaged by  $m$ :

$$J(w, b) = (1/m) * \sum(L(\hat{y}(i), y(i))) \text{ from } i=1 \text{ to } m$$

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. X

$\log(1 - \hat{y}(i))$  from  $i=1$  to  $m$

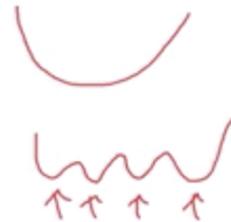
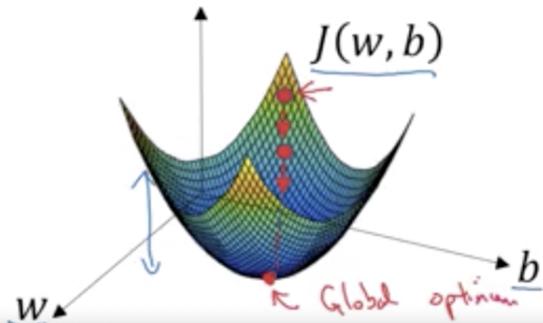
The Gradient Descent goal is to find  $w, b$  that minimize cost function  $J(w, b)$ . Each gradient descent takes steps until the global optimum is reached.

## Gradient Descent

Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$  ↪

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find  $w, b$  that minimize  $J(w, b)$



Andrew Ng

Gradient Descent (Source: DeepLearning.ai / Andrew Ng)

Gradient descent repeats the following where alpha is the learning rate (controls the size of the step of gradient descent) and  $dJ(w, b)/dw$  is the derivative that is the slope and known as “dw” in code.  $w$  “weight” is being updated as  $w$  and derivative changes.  $b$  “bias” is being updated as  $b$  and derivative changes.

$$w = w - \alpha * dJ(w, b) / dw$$

$$b = b - \alpha * dJ(w, b) / db$$

As the gradient descent steps through the updates, global optimum goes “down hill” until the slope approximates to 0 ideally.

Derivatives with a Computation Graph: This section deals with using a computation graph to calculate derivatives. The exercise plugs in values and finds the ratio by chain

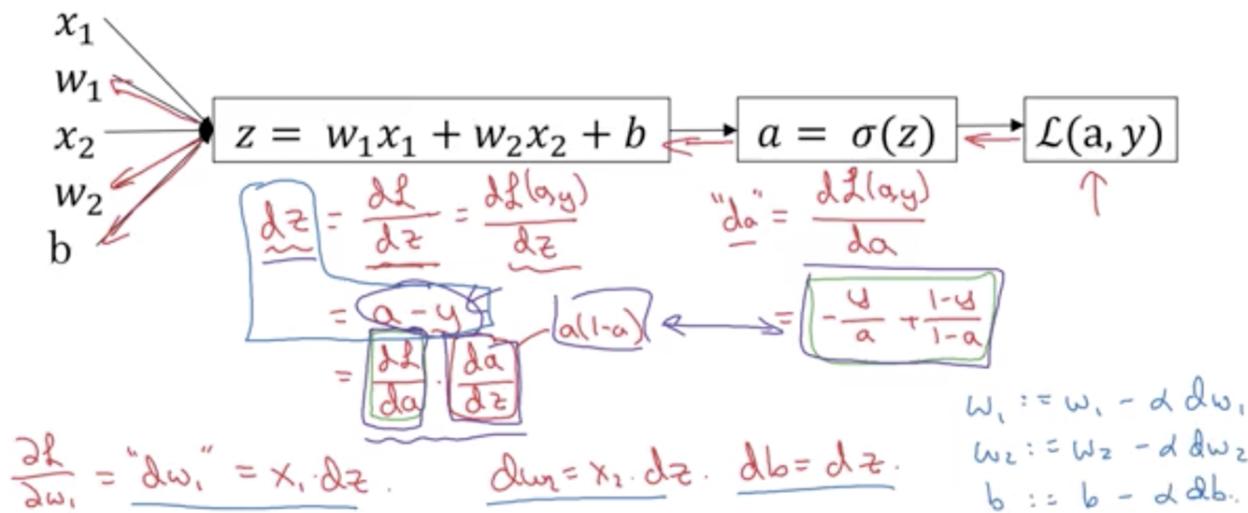
To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including our cookie policy.



## Logistic Regression Gradient Descent:

The following describes the computation graph using logistic regression. Variables  $x_1$ ,  $w_1$ ,  $x_2$ ,  $w_2$ , and  $b$  are inputs into  $z = w_1x_1 + w_2x_2 + b$ .  $Z$  is an input into the sigmoid function, resulting in ‘ $a$ ’ which is the same as  $y\text{-hat}$ . The loss function  $L$  is a function of  $a$  and  $y$ .

## Logistic regression derivatives



Andrew Ng

Logistic regression derivatives (Source: DeepLearning.ai / Andrew Ng)

The goal of the derivatives is to find the  $w_1$ ,  $w_2$ , and  $b$  that have the minimal loss based on  $L(a,y)$ . From  $L(a,y)$ , derivative ‘ $da$ ’ can be calculated; from ‘ $a$ ,’ derivative  $dz$ ; and from  $dz$ ,  $dw_1$ ,  $dw_2$ , and  $db$ . Based on  $dw_1$ ,  $dw_2$ , and  $db$ , the updates to  $w_1$ ,  $w_2$ , and  $b$  can be calculated:

$$da = -(y/a) + (1-y)/(1-a)$$

$$dz = a - y$$

$$dw_1 = x_1 * dz$$

$$dw_2 = x_2 * dz$$

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.



## Gradient Descent on m Examples:

Recap of cost function  $J(w,b)$  which is an average of the sum over all the losses:

$$J(w,b) = \frac{1}{m} * \text{sum}(L(a(i), y(i))) \text{ from } i=1 \text{ to } m \text{ examples.}$$

$$a(i) = y\text{-hat}(i) = \text{sigmoid}(z(i)) = \text{sigmoid}(w^T x(i) + b).$$

The derivative respect to  $dw_1$  of cost function  $J(w,b)$  is similar in approach and is the average of  $dw_1$  across  $i$  to  $m$  examples:

$$\frac{\partial J(w,b)}{\partial w_1} = \frac{1}{m} * \text{sum}(dw_1 * L(a(i), y(i))) = dw_1(i) - (x(i), y(i))$$

In an algorithm for one step of gradient descent,  $J$ ,  $dw_1$ ,  $dw_2$ ,  $db = 0$  on initialization. In the for loop for  $i=1$  to  $m$ , a series of calculation is made to calculate the  $i$ -th example of  $z$ ,  $a$ ,  $J$ ,  $dz$ ,  $dw_1$ ,  $dw_2$ ,  $db$ . After the loop,  $J$ ,  $dw_1$ ,  $dw_2$ , and  $db$  are averaged by  $m$ . Once this loop is complete, we get the derivatives of  $dw_1$ ,  $dw_2$ , and  $db$  — not the same as the  $dw_1(i)$ , etc of a single training example. The derivative is used to update the parameters  $w_1$ ,  $w_2$ , and  $b$ .

## Logistic regression on $m$ examples

$$\begin{aligned} J &= 0; \underline{\delta w_1} = 0; \underline{\delta w_2} = 0; \underline{\delta b} = 0 \\ \text{For } i &= 1 \text{ to } m \\ z^{(i)} &= w^T x^{(i)} + b \\ a^{(i)} &= \sigma(z^{(i)}) \\ J &+= -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})] \\ \underline{\delta z^{(i)}} &= a^{(i)} - y^{(i)} \\ \underline{\delta w_1} &+= x_1^{(i)} \underline{\delta z^{(i)}} \quad \downarrow n=2 \\ \underline{\delta w_2} &+= x_2^{(i)} \underline{\delta z^{(i)}} \\ \underline{\delta b} &+= \underline{\delta z^{(i)}} \end{aligned}$$

$$\left| \begin{aligned} \underline{dw_1} &= \frac{\partial J}{\partial w_1} \\ w_1 &:= w_1 - \alpha \underline{dw_1} \\ w_2 &:= w_2 - \alpha \underline{dw_2} \\ b &:= b - \alpha \underline{db} \end{aligned} \right.$$

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.



Andrew Ng

Logistic regression on m examples (Source: DeepLearning.ai / Andrew Ng)

Note that the above approach shows a single for loop. When there are multiple layers in deep learning, an option is to have a nested for loop within the  $i=1$  to  $m$  for loop at the  $dz(i)$  to  $db$  variables. The better approach with large datasets is to use vectorization which is more efficient.

## Python and Vectorization:

For equation  $z = w^T x + b$  where  $T$  is transpose of  $w$ , the non-vectorized approach looks like the following. In the for loop, we're iterating through  $i$  to  $n$  examples. Then we're multiplying the  $i$ -th ' $w$ ' and ' $x$ ' together, and updating ' $z$ ' with this plus the ' $b$ ' bias:

$$z = 0$$

for  $i$  in range( $n-x$ ):

$$z += w[i] * x[i]$$

$$z += b$$

Instead, the vectorized approach has  $w$  and  $x$  in vector column form where both are  $n$  of  $x$  examples. For example, a ' $w$ ' vector of shape (10,1) and ' $x$ ' vector of shape (1,5) multiplied together using numpy library's dot() method returns a matrix  $z$  of shape (10,5). This eliminates two for-loops and is faster.

$$z = np.dot(w, x)$$

## Logistic regression derivatives

$$\begin{aligned} J &= 0, \quad \cancel{dw1 = 0}, \quad \cancel{dw2 = 0}, \quad db = 0 & dw &= np.zeros((n-x, 1)) \\ \rightarrow \text{for } i &= 1 \text{ to } m: \\ z^{(i)} &= w^T x^{(i)} + b \end{aligned}$$

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X

$$\begin{aligned}
 dz^{(i)} &= a^{(i)}(1 - a^{(i)}) \\
 dw_1 &+= x_1^{(i)} dz^{(i)} \quad n_x = 2 \\
 dw_2 &+= x_2^{(i)} dz^{(i)} \\
 db &+= dz^{(i)} \\
 J &= J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m, \quad db = db/m \\
 dw &/= m
 \end{aligned}$$

Andrew Ng

More Vectorization Examples (Source: DeepLearning.ai / Andrew Ng)

Using vectorization / matrix multiplication, the inside for-loop for  $j=1$  to  $n_x$  features,  $dw_j +=$  updates eliminates the  $dw_1$  and  $dw_2$  updates with a  $dw += x(i)*dz(i)$ . On the top in green, the respective  $dw_1, dw_2 = 0$  is replaced with `np.zeros()`. On the bottom in green, the  $dw_1 = dw_1/m$  and  $dw_2 = dw_2/m$  is simply replaced with  $dw /= m$ .

## Vectorizing Logistic Regression:

In the forward propagation (left to right equations) of the logistic regression, we stack 'm' training examples of each parameter into matrices. For example  $x(1\dots m)$  is in a  $X$  matrix of shape ( $n_x$  rows of features, and  $m$  columns of each  $x$  training examples) or simply  $(n, m)$ . As result,  $W$ ,  $Z$ , and  $A$  are the respective matrices across  $m$  training examples. Note: 'b' bias parameter is of shape  $(1, m)$ .

## Vectorizing Logistic Regression

$$\begin{aligned}
 \rightarrow z^{(1)} &= w^T x^{(1)} + b \\
 \rightarrow a^{(1)} &= \sigma(z^{(1)}) \\
 X &= \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad \frac{(n_x, m)}{\mathbb{R}^{n_x \times m}} \\
 \rightarrow z^{(2)} &= w^T x^{(2)} + b \\
 \rightarrow a^{(2)} &= \sigma(z^{(2)}) \\
 w &\in \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \\
 z^{(3)} &= w^T x^{(3)} + b \\
 a^{(3)} &= \sigma(z^{(3)}) \\
 Z &= \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} w^T x^{(1)} + b \\ w^T x^{(2)} + b \\ \vdots \\ w^T x^{(m)} + b \end{bmatrix}_{1 \times m} \quad "Broadcasting" \\
 A &= [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)
 \end{aligned}$$

Andrew Ng

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X

## Vectorizing Logistic Regression's Gradient Output:

In order to remove the outer for-loop along  $i=1$  to  $m$  examples, vectorization of the gradients is required. The following set of equations addresses a single iteration of gradient descent for logistic regression (Note:  $T$  represents transpose;  $:=$  notation is simply updating the parameter from the prior value):

$$Z = w.T * X + b = \text{np.dot}(w.T, X) + b$$

$$A = \text{sigmoid}(Z)$$

$$dZ = A - Y$$

$$dW = (1/m) * X & dZ.T$$

$$db = (1/m) * \text{np.sum}(dZ)$$

$$w := w - \text{alpha} * dW$$

$$b := b - \text{alpha} * db$$

The above operations depend on a Python concept called broadcasting. For example, when you have a vector of shape  $(4,1)$  and want to add a real number such as 100, then 100 is added to each vector element, resulting in a same shape vector  $(4,1)$ . Python represents 100 as  $(4,1)$ . Another example of a matrix  $(2,3)$  plus vector  $(1,3)$  is interpreted as two rows of  $(1,3)$ , essentially a  $(2,3)$ . The resulting matrix is the element-wise operation.

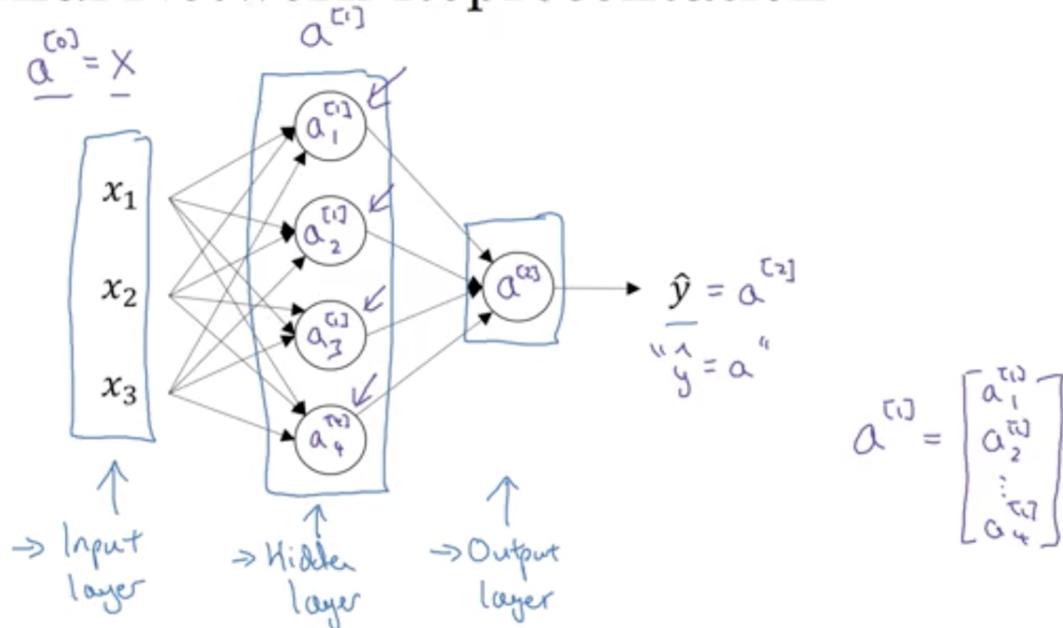
According to Professor Ng, there may be challenging bugs when there are rank 1 arrays (e.g., shape  $(5,)$ ). Recommendation is to reshape into  $(n,1)$  using the shape (e.g., `array.shape == (5,1)`).

## Week 3— Shallow Neural Networks

With this section, a two layer neural network demonstrates how this network works. Notations to be aware of: 1) Previously parentheticals () represented the individual

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

## Neural Network Representation



Andrew Ng

Neural Network Representation (Source: DeepLearning.ai / Andrew Ng)

In the input layer, there are a number of features (e.g.,  $x_1, x_2, x_3$  could be a pixel in the RGB layers of a image). The inputs are passed to the hidden layer represented by a superscript [1]. In this example there are 4 nodes each represented as a subscript (number of node) superscript [layer number].  $a[1]$  is a matrix  $(4,1)$ . In hidden layer 1,  $a[1]$  has  $w, b$  represented by  $w[1]$  and  $b[1]$ . Shape of  $w[1]$  is  $(4,3)$  and  $b[1]$  is  $(4,1)$ . Finally the output is  $y\text{-hat}$  represented by  $a[2]$  and is shape  $(1,1)$ .

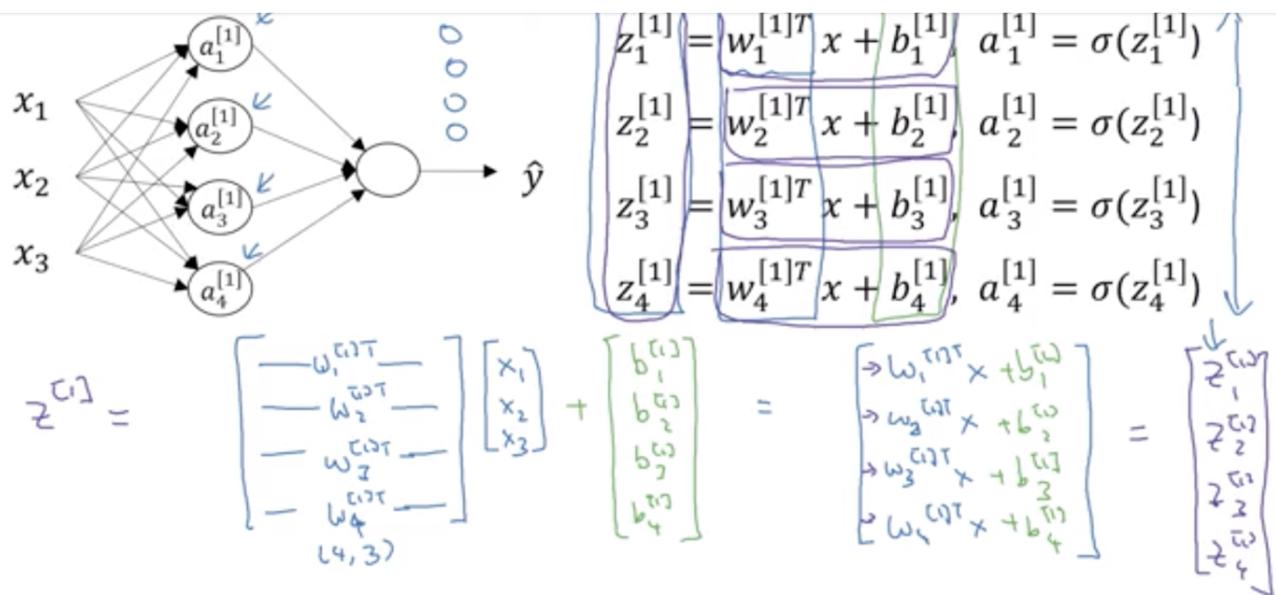
Personally, I keep track of layer matrix shapes with this rule of thumb:

- Input layer is shape (number of features, 1) or (n-x, 1)
  - Hidden layer 1 is shape (number of nodes, n-x)
  - If there is another hidden layer 2, shape is (number of nodes of hidden layer 2, number of nodes of hidden layer 1).
  - Output layer is shape (1, number of nodes of last hidden layer)

Below, the notation  $a$ -subscript 1, superscript[1] means the first node of the first layer. The following shows what is behind each activation node ' $a$ ' and visualizing the shapes of the matrices:

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

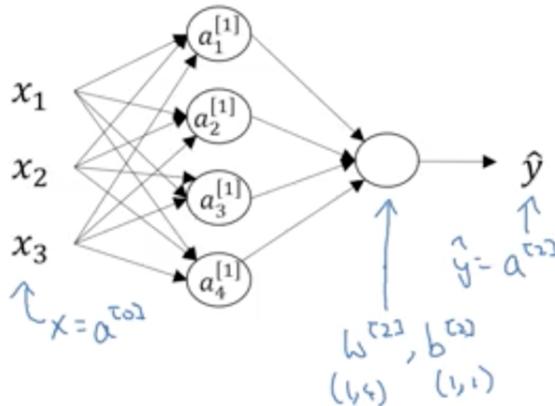
X



Andrew Ng

(Source: DeepLearning.ai / Andrew Ng)

## Neural Network Representation learning



Given input  $x$ :

$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]}_{(4,3)} x^{(3,1)} + b^{[1]}_{(4,1)} \\ \rightarrow a^{[1]} &= \sigma(z^{[1]})_{(4,1)} \\ \rightarrow z^{[2]} &= W^{[2]}_{(1,4)} a^{[1]}_{(4,1)} + b^{[2]}_{(1,1)} \\ \rightarrow a^{[2]} &= \sigma(z^{[2]})_{(1,1)} \end{aligned}$$

Andrew Ng

(Source: DeepLearning.ai / Andrew Ng)

A lot of the debugging in neural networks is ensuring that the shapes of the parameters in the matrices meet expectations. Shape of input, hidden, and output layers are  $(3,1)$ ,  $(4,3)$ , and  $(1,4)$  respectively. Shape of  $z^{[1]}$  is  $(4,1)$  because of the  $W^{[1]} * x$  is  $(4,3)$  multiplied by  $(3,1)$ , resulting in  $(4,1)$ . Then,  $b^{[1]} (4,1)$  is added by to the previous.

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X

## Vectorizing across multiple examples

for  $i = 1$  to  $m$ :

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$X = \begin{bmatrix} & & & \\ \downarrow & \downarrow & \downarrow & \downarrow \\ w & 1 & 1 & 1 \\ X & X & \dots & X \\ \uparrow & | & | & | \\ 1 & (n_x, m) & & 1 \end{bmatrix}$$

↑  
hidden units.  
↑  
training samples

$$\begin{aligned} z^{[1]} &= W^{[1]}X + b^{[1]} \\ \rightarrow A^{[1]} &= \sigma(z^{[1]}) \\ \rightarrow z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ \rightarrow A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

$$\begin{aligned} z^{[1]} &= \begin{bmatrix} 1 & 1 & \dots & 1 \\ z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \\ A^{[1]} &= \begin{bmatrix} \bullet & \bullet & \dots & \bullet \\ a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \end{aligned}$$

↑  
# hidden units  
Andrew Ng

(Source: DeepLearning.ai / Andrew Ng)

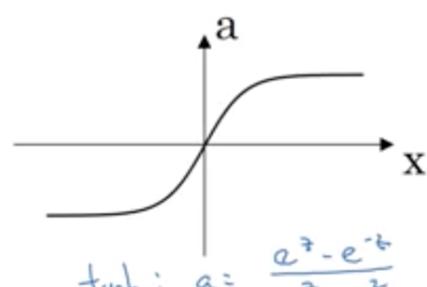
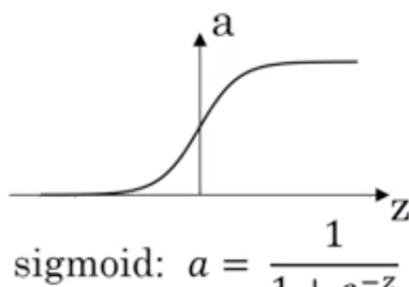
## Activation Functions:

When  $z$  is very large or very small, the slope of the sigmoid or tan-h function is zero, and can slow down gradient descent calculations.

Instead of sigmoid-based activation function in the hidden layers, the tan-h centers the function around zero. Also increasing popular is the rectified linear unit function (ReLU) where a positive value of  $z$  results in a slope of 1; a negative  $z$  value is slope zero.

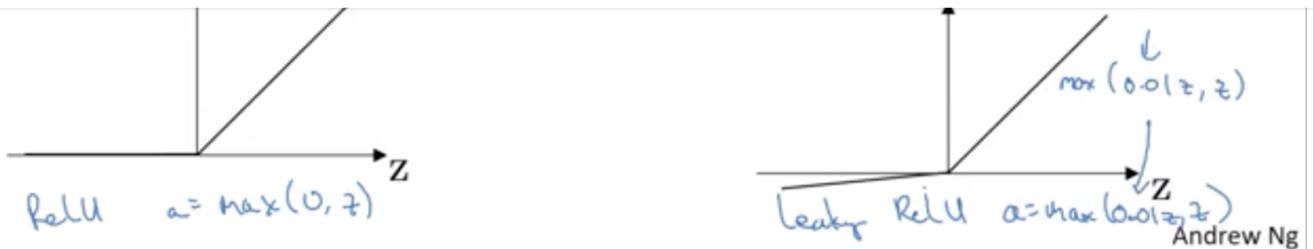
Sometimes a hidden layer would use tan-h or ReLU and the output layer would be sigmoid. Different activation functions by layer are possible.

## Pros and cons of activation functions



To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

X



Activation Functions (Source: DeepLearning.ai / Andrew Ng)

## Derivatives of Activation Functions:

For a sigmoid function  $g(z)$ , the derivative is the slope of  $g(z)$  at  $z$ . When  $z$  is very large at either positive or negative ends, the slope is zero. The derivative function of activation function with 'a' representing  $g(z)$ :

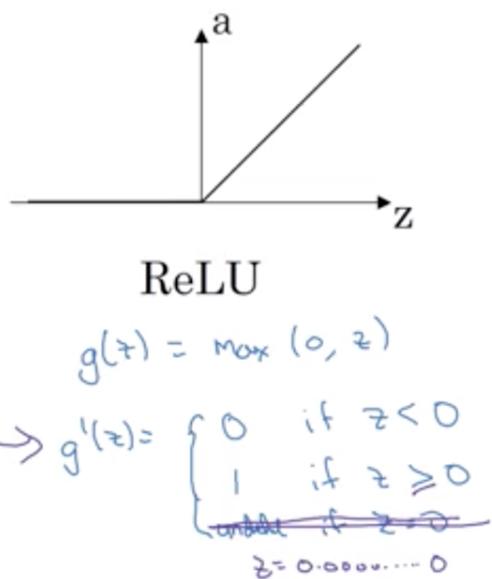
$$g'(z) = g(z) * (1 - g(z)) = a * (1 - a)$$

For tanh function, the derivative of the activation function:

$$g'(z) = 1 - (\tanh(z))^2 = 1 - a^2.$$

For ReLU and Leaky ReLU functions, the derivatives of the activation functions:

## ReLU and Leaky ReLU



Derivatives of Activation Functions (Source: DeepLearning.ai / Andrew Ng)

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. X

forward and back pass being one step of the descent. Initialize parameters randomly and repeat until losses converge.

## Formulas for computing derivatives

Forward propagation:

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]}) \leftarrow$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]}) = \underline{g}(Z^{[2]})$$

Back propagation:

$$dZ^{[2]} = A^{[2]} - Y \leftarrow$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[2]}, m)} \times \underbrace{g^{[2]'}(Z^{[1]})}_{\text{element-wise product}} \quad (n^{[1]}, m)$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Gradient descent for neural network (Source: DeepLearning.ai / Andrew Ng)

In a simple neural network with two layers, the forward propagation equations for Z[1] is calculated based on W[1] and b[1]. A[1] is the sigmoid of the Z[1]. Z[2] is the result of W[2] times the prior A, which is A[1], plus b[2]. A[2] is the sigmoid of Z[2]. Z of current layer is always using the prior layer A and so forth.

In the backpropagation, we start with the last hidden layer. dZ[2] is the difference between A[2] and Y, which represents the predicted and actual results. dW[2] = (1/m) \* dZ[2] \* A[1] transpose. db[2] is the average of the sum of dz[2] where axis=1 and keepdims=True so that the shape of the db[2] matrix will be (n,1) and not (n,), which can result in bugs. dZ[1] is the W[2] transpose \* dZ[2] times elementwise multiplication of sigmoid of Z[1]. dZ is shape (n[1],m). Respective formulas are applied to calculate dW[1] and db[1].

## Random Initialization

If weights W and b are initialized to zero, each node of the neural network does the same calculation and nothing is learned. Except for b, the solution is to initialize weights randomly.

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

number of hidden layers. Essentially, more than single or two hidden layers. The input layer is layer zero.  $n[L]$  is the number of nodes/units in layer L.

## Getting your matrix dimensions right

The key part of making neural network code work properly is ensuring that we have the right shapes. In an individual training example of  $z[1] = w[1]^* x + b[1]$ , these are the following shapes:

$z$ : ( number of nodes for current layer , 1 )

$w$ : ( number of nodes for current layer , number of input nodes from previous layer )

$x$ : ( number of input nodes from previous layer , 1 )

$b$ : ( number of nodes for current layer , 1 )

Shape of  $w = dw$  and  $b = db$

In the vectorized across all training examples for  $Z = (W^*X + B)$  or  $(W^*A + B)$ , these are the shapes of the matrices:

$Z$ : ( nodes current layer , m training examples)

$W$ : ( nodes current layer , nodes previous layer )

$X$  or  $A$ : ( nodes previous layer , m training examples )

$B$ : ( nodes current t layer , m training examples )

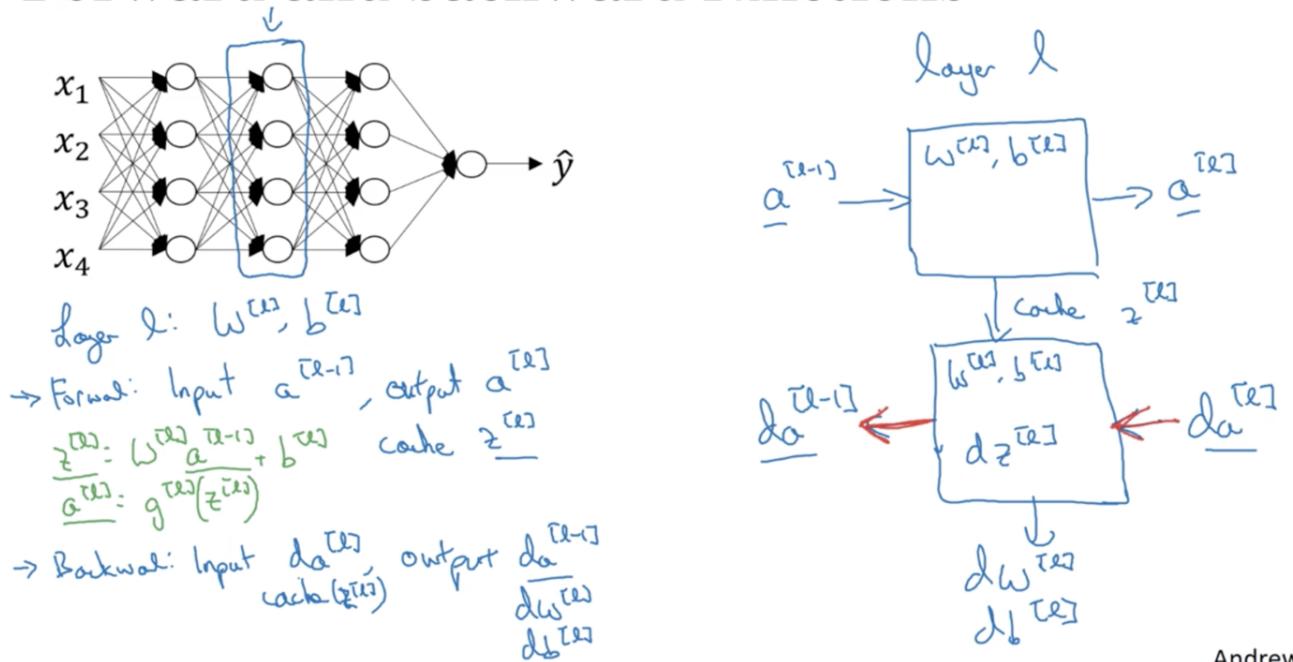
$Z = dZ$  shape and likewise for other derivatives

## Building blocks of deep neural networks

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. X

propagation, the reverse happens where the derivative of  $\hat{y}$  results in  $da[1]$  input to output  $da[l-1]$  and parameters  $dw$  and  $db$ , which are cached like  $z$  in the forward propagation. The  $w$  and  $b$  are updated by the learning rate multiplied by the derivatives. All this accounts for one pass of gradient descent to calculate the optimal  $w$  and  $b$ . The cache is used to keep the parameters for accessibility.

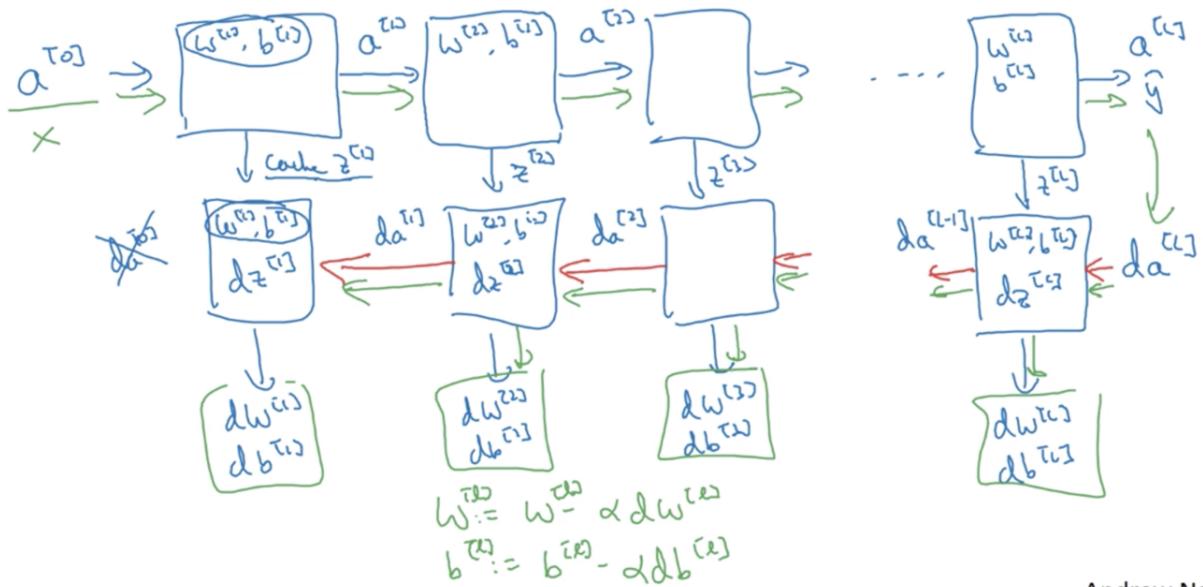
## Forward and backward functions



Andrew Ng

Building blocks of deep neural networks (Source: DeepLearning.ai / Andrew Ng)

## Forward and backward functions



Andrew Ng

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.



## Forward and backward propagation

The following shows the calculations for backward propagation using vectorization:

### Backward propagation for layer $l$

→ Input  $\underline{da}^{[l]}$

→ Output  $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$\begin{cases} \underline{dz}^{[l]} = \underline{da}^{[l]} * g^{[l]}(z^{[l]}) \\ dW^{[l]} = dz^{[l]} \cdot a^{[l-1]} \\ db^{[l]} = dz^{[l]} \\ da^{[l-1]} = W^{[l]} \cdot dz^{[l]} \\ dz^{[l]} = W^{[l+1]} dz^{[l+1]} + g^{[l]}(z^{[l]}) \end{cases}$$

$$\begin{cases} dz^{[l]} = dA^{[l]} * g^{[l]}(z^{[l]}) \\ dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} np.sum(dz^{[l]}), axis=1, keepdims=True \\ dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]} \end{cases}$$

Andrew Ng

(Source: DeepLearning.ai / Andrew Ng)

• • •

Learn together with me and follow my blog as I share more about machine learning, data science, and product management.

Machine Learning

Data Science

Deep Learning

Neural Networks

About Help Legal