

Building your Deep Neural Network: Step by Step

Notation:

- Superscript $[l]$ denotes a quantity associated with the l^{th} layer.
 - Example: $a^{[L]}$ is the L^{th} layer activation. $W^{[L]}$ and $b^{[L]}$ are the L^{th} layer parameters.
- Superscript (i) denotes a quantity associated with the i^{th} example.
 - Example: $x^{(i)}$ is the i^{th} training example.
- Lowercase i denotes the i^{th} entry of a vector.
 - Example: $a_i^{[l]}$ denotes the i^{th} entry of the l^{th} layer's activations).

1 - Packages

Let's first import all the packages that you will need during this assignment.

- [numpy](http://www.numpy.org) (www.numpy.org) is the main package for scientific computing with Python.
- [matplotlib](http://matplotlib.org) (<http://matplotlib.org>) is a library to plot graphs in Python.
- dnn_utils provides some necessary functions for this notebook.
- testCases provides some test cases to assess the correctness of your functions
- np.random.seed(1) is used to keep all the random function calls consistent. It will help us grade your work. Please don't change the seed.

```
In [4]: import numpy as np
import h5py
import matplotlib.pyplot as plt
from testCases_v4 import *
from dnn_utils_v2 import sigmoid, sigmoid_backward, relu, relu_backward

%matplotlib inline
plt.rcParams['figure.figsize'] = (5.0, 4.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

%load_ext autoreload
%autoreload 2

np.random.seed(1)
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

2 - Outline of the Assignment

To build your neural network, you will be implementing several "helper functions". These helper functions will be used in the next assignment to build a two-layer neural network and an L-layer neural network. Each small helper function you will implement will have detailed instructions that will walk you through the necessary steps. Here is an outline of this assignment, you will:

- Initialize the parameters for a two-layer network and for an L -layer neural network.
- Implement the forward propagation module (shown in purple in the figure below).
 - Complete the LINEAR part of a layer's forward propagation step (resulting in $Z^{[l]}$).
 - We give you the ACTIVATION function (relu/sigmoid).
 - Combine the previous two steps into a new [LINEAR->ACTIVATION] forward function.
 - Stack the [LINEAR->RELU] forward function $L-1$ time (for layers 1 through $L-1$) and add a [LINEAR->SIGMOID] at the end (for the final layer L). This gives you a new `L_model_forward` function.
- Compute the loss.
- Implement the backward propagation module (denoted in red in the figure below).
 - Complete the LINEAR part of a layer's backward propagation step.
 - We give you the gradient of the ACTIVATE function (`relu_backward`/`sigmoid_backward`)
 - Combine the previous two steps into a new [LINEAR->ACTIVATION] backward function.
 - Stack [LINEAR->RELU] backward $L-1$ times and add [LINEAR->SIGMOID] backward in a new `L_model_backward` function
- Finally update the parameters.

3 - Initialization

3.1 - 2-layer Neural Network

Create and initialize the parameters of the 2-layer neural network.

- The model's structure is: *LINEAR -> RELU -> LINEAR -> SIGMOID*.
- Use random initialization for the weight matrices. Use `np.random.randn(shape)*0.01` with the correct shape.
- Use zero initialization for the biases. Use `np.zeros(shape)`.

```
In [5]: def initialize_parameters(n_x, n_h, n_y):
    """
        Argument:
        n_x -- size of the input layer
        n_h -- size of the hidden layer
        n_y -- size of the output layer

        Returns:
        parameters -- python dictionary containing your parameters:
            W1 -- weight matrix of shape (n_h, n_x)
            b1 -- bias vector of shape (n_h, 1)
            W2 -- weight matrix of shape (n_y, n_h)
            b2 -- bias vector of shape (n_y, 1)
    """
    np.random.seed(1)

    W1 = np.random.randn(n_h,n_x) * 0.01
    b1 = np.zeros((n_h,1))
    W2 = np.random.randn(n_y,n_h) * 0.01
    b2 = np.zeros((n_y,1))

    assert(W1.shape == (n_h, n_x))
    assert(b1.shape == (n_h, 1))
    assert(W2.shape == (n_y, n_h))
    assert(b2.shape == (n_y, 1))

    parameters = {"W1": W1,
                  "b1": b1,
                  "W2": W2,
                  "b2": b2}

    return parameters
```

```
In [6]: parameters = initialize_parameters(3,2,1)
print("W1 = " + str(parameters["W1"]))
print("b1 = " + str(parameters["b1"]))
print("W2 = " + str(parameters["W2"]))
print("b2 = " + str(parameters["b2"]))
```

```
W1 = [[ 0.01624345 -0.00611756 -0.00528172]
      [-0.01072969  0.00865408 -0.02301539]]
b1 = [[ 0.]
      [ 0.]]
W2 = [[ 0.01744812 -0.00761207]]
b2 = [[ 0.]]
```

Expected output:

W1	[[0.01624345 -0.00611756 -0.00528172] [-0.01072969 0.00865408 -0.02301539]]
b1	[[0.] [0.]]
W2	[[0.017444812 -0.00761207]]
b2	[[0.]]

3.2 - L-layer Neural Network

The initialization for a deeper L-layer neural network is more complicated because there are many more weight matrices and bias vectors. When completing the `initialize_parameters_deep`, you should make sure that your dimensions match between each layer. Recall that $n^{[l]}$ is the number of units in layer l . Thus for example if the size of our input X is (12288, 209) (with $m = 209$ examples) then:

	Shape of W	**Shape of b**	**Activation**	**Shape of Activation**
Layer 1	($n^{[1]}$, 12288)	($n^{[1]}$, 1)	$Z^{[1]} = W^{[1]}X + b^{[1]}$	($n^{[1]}$, 209)
Layer 2	($n^{[2]}$, $n^{[1]}$)	($n^{[2]}$, 1)	$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$	($n^{[2]}$, 209)
:	:	:	:	:
Layer L-1	($n^{[L-1]}$, $n^{[L-2]}$)	($n^{[L-1]}$, 1)	$Z^{[L-1]} = W^{[L-1]}A^{[L-2]} + b^{[L-1]}$	($n^{[L-1]}$, 209)
Layer L	($n^{[L]}$, $n^{[L-1]}$)	($n^{[L]}$, 1)	$Z^{[L]} = W^{[L]}A^{[L-1]} + b^{[L]}$	($n^{[L]}$, 209)

Remember that when we compute $WX + b$ in python, it carries out broadcasting. For example, if:

$$W = \begin{bmatrix} j & k & l \\ m & n & o \\ p & q & r \end{bmatrix} \quad X = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \quad b = \begin{bmatrix} s \\ t \\ u \end{bmatrix} \quad (2)$$

Then $WX + b$ will be:

$$WX + b = \begin{bmatrix} (ja + kd + lg) + s & (jb + ke + lh) + s & (jc + kf + li) + s \\ (ma + nd + og) + t & (mb + ne + oh) + t & (mc + nf + oi) + t \\ (pa + qd + rg) + u & (pb +qe + rh) + u & (pc + qf + ri) + u \end{bmatrix} \quad (3)$$

Implement initialization for an L-layer Neural Network.

- The model's structure is $[LINEAR \rightarrow RELU] \times (L-1) \rightarrow LINEAR \rightarrow SIGMOID$. I.e., it has $L - 1$ layers using a ReLU activation function followed by an output layer with a sigmoid activation function.
- Use random initialization for the weight matrices. Use `np.random.randn(shape) * 0.01`.
- Use zeros initialization for the biases. Use `np.zeros(shape)`.
- We will store $n^{[l]}$, the number of units in different layers, in a variable `layer_dims`. For example, the `layer_dims` for the "Planar Data classification model" from last week would have been [2,4,1]: There were two inputs, one hidden layer with 4 hidden units, and an output layer with 1 output unit. Thus means $W1$'s shape was (4,2), $b1$ was (4,1), $W2$ was (1,4) and $b2$ was (1,1). Now you will generalize this to L layers!
- Here is the implementation for $L = 1$ (one layer neural network). It should inspire you to implement the general case (L-layer neural network).

```
if L == 1:
    parameters["W" + str(L)] = np.random.randn(layer_dims[1], layer_dims[0]
) * 0.01
    parameters["b" + str(L)] = np.zeros((layer_dims[1], 1))
```

```
In [7]: def initialize_parameters_deep(layer_dims):
    """
    Arguments:
    layer_dims -- python array (List) containing the dimensions of each layer
    in our network

    Returns:
    parameters -- python dictionary containing your parameters "W1", "b1",
    ... , "WL", "bL":
        WL -- weight matrix of shape (Layer_dims[L], Layer_dims[L-
1])
        bl -- bias vector of shape (Layer_dims[L], 1)
    """
    np.random.seed(3)
    parameters = {}
    L = len(layer_dims) # number of layers in the network

    for l in range(1, L):
        parameters['W' + str(l)] = np.random.randn(layer_dims[l], layer_dims[l-
1]) * 0.01
        parameters['b' + str(l)] = np.zeros((layer_dims[l], 1))

        assert(parameters['W' + str(l)].shape == (layer_dims[l], layer_dims[l-
1]))
        assert(parameters['b' + str(l)].shape == (layer_dims[l], 1))

    return parameters
```

```
In [8]: parameters = initialize_parameters_deep([5,4,3])
print("W1 = " + str(parameters["W1"]))
print("b1 = " + str(parameters["b1"]))
print("W2 = " + str(parameters["W2"]))
print("b2 = " + str(parameters["b2"]))
```

```
W1 = [[ 0.01788628  0.0043651   0.00096497 -0.01863493 -0.00277388]
[-0.00354759 -0.00082741 -0.00627001 -0.00043818 -0.00477218]
[-0.01313865  0.00884622  0.00881318  0.01709573  0.00050034]
[-0.00404677 -0.0054536   -0.01546477  0.00982367 -0.01101068]]
b1 = [[ 0.]
[ 0.]
[ 0.]
[ 0.]]
W2 = [[-0.01185047 -0.0020565   0.01486148  0.00236716]
[-0.01023785 -0.00712993  0.00625245 -0.00160513]
[-0.00768836 -0.00230031  0.00745056  0.01976111]]
b2 = [[ 0.]
[ 0.]
[ 0.]]
```

Expected output:

W1	[[0.01788628 0.0043651 0.00096497 -0.01863493 -0.00277388] [-0.00354759 -0.00082741 -0.00627001 -0.00043818 -0.00477218] [-0.01313865 0.00884622 0.00881318 0.01709573 0.00050034] [-0.00404677 -0.0054536 -0.01546477 0.00982367 -0.01101068]]
b1	[[0.] [0.] [0.] [0.]]
W2	[[-0.01185047 -0.0020565 0.01486148 0.00236716] [-0.01023785 -0.00712993 0.00625245 -0.00160513] [-0.00768836 -0.00230031 0.00745056 0.01976111]]
b2	[[0.] [0.] [0.]]

4 - Forward propagation module

4.1 - Linear Forward

Now that you have initialized your parameters, you will do the forward propagation module. You will start by implementing some basic functions that you will use later when implementing the model. You will complete three functions in this order:

- LINEAR
- LINEAR -> ACTIVATION where ACTIVATION will be either ReLU or Sigmoid.
- [LINEAR -> RELU] \times (L-1) -> LINEAR -> SIGMOID (whole model)

The linear forward module (vectorized over all the examples) computes the following equations:

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]} \quad (4)$$

where $A^{[0]} = X$.

Build the linear part of forward propagation.

```
In [9]: def linear_forward(A, W, b):
    """
    Implement the Linear part of a Layer's forward propagation.

    Arguments:
    A -- activations from previous layer (or input data): (size of previous layer, number of examples)
    W -- weights matrix: numpy array of shape (size of current layer, size of previous layer)
    b -- bias vector, numpy array of shape (size of the current layer, 1)

    Returns:
    Z -- the input of the activation function, also called pre-activation parameter
    cache -- a python dictionary containing "A", "W" and "b" ; stored for computing the backward pass efficiently
    """
    Z = np.dot(W,A) + b

    assert(Z.shape == (W.shape[0], A.shape[1]))
    cache = (A, W, b)

    return Z, cache
```

```
In [10]: A, W, b = linear_forward_test_case()

Z, linear_cache = linear_forward(A, W, b)
print("Z = " + str(Z))

Z = [[ 3.26295337 -1.23429987]]
```

Expected output:

Z	[[3.26295337 -1.23429987]]
-------	-----------------------------

4.2 - Linear-Activation Forward

In this notebook, you will use two activation functions:

- **Sigmoid:** $\sigma(Z) = \sigma(WA + b) = \frac{1}{1+e^{-(WA+b)}}$. We have provided you with the `sigmoid` function. This function returns **two** items: the activation value "a" and a "cache" that contains "Z" (it's what we will feed in to the corresponding backward function). To use it you could just call:

```
A, activation_cache = sigmoid(Z)
```

- **ReLU:** The mathematical formula for ReLU is $A = RELU(Z) = \max(0, Z)$. We have provided you with the `relu` function. This function returns **two** items: the activation value "A" and a "cache" that contains "Z" (it's what we will feed in to the corresponding backward function). To use it you could just call:

```
A, activation_cache = relu(Z)
```

For more convenience, you are going to group two functions (Linear and Activation) into one function (`LINEAR- > ACTIVATION`). Hence, you will implement a function that does the LINEAR forward step followed by an ACTIVATION forward step.

Implement the forward propagation of the `LINEAR->ACTIVATION` layer. Mathematical relation is:

$A^{[l]} = g(Z^{[l]}) = g(W^{[l]} A^{[l-1]} + b^{[l]})$ where the activation "g" can be `sigmoid()` or `relu()`. Use `linear_forward()` and the correct activation function.

```
In [11]: def linear_activation_forward(A_prev, W, b, activation):
    """
        Implement the forward propagation for the LINEAR->ACTIVATION Layer

        Arguments:
            A_prev -- activations from previous Layer (or input data): (size of previous Layer, number of examples)
            W -- weights matrix: numpy array of shape (size of current Layer, size of previous layer)
            b -- bias vector, numpy array of shape (size of the current Layer, 1)
            activation -- the activation to be used in this layer, stored as a text string: "sigmoid" or "relu"

        Returns:
            A -- the output of the activation function, also called the post-activation n value
            cache -- a python dictionary containing "linear_cache" and "activation_cache";
                    stored for computing the backward pass efficiently
    """

    if activation == "sigmoid":
        # Inputs: "A_prev, W, b". Outputs: "A, activation_cache".
        Z, linear_cache = linear_forward(A_prev, W, b)
        A, activation_cache = sigmoid(Z)

    elif activation == "relu":
        # Inputs: "A_prev, W, b". Outputs: "A, activation_cache".
        Z, linear_cache = linear_forward(A_prev, W, b)
        A, activation_cache = relu(Z)

    assert (A.shape == (W.shape[0], A_prev.shape[1]))
    cache = (linear_cache, activation_cache)

    return A, cache
```

```
In [12]: A_prev, W, b = linear_activation_forward_test_case()

A, linear_activation_cache = linear_activation_forward(A_prev, W, b, activation = "sigmoid")
print("With sigmoid: A = " + str(A))

A, linear_activation_cache = linear_activation_forward(A_prev, W, b, activation = "relu")
print("With ReLU: A = " + str(A))

With sigmoid: A = [[ 0.96890023  0.11013289]]
With ReLU: A = [[ 3.43896131  0.          ]]
```

Expected output:

**With sigmoid: A **	[[0.96890023 0.11013289]]
**With ReLU: A **	[[3.43896131 0.]]

Note: In deep learning, the "[LINEAR->ACTIVATION]" computation is counted as a single layer in the neural network, not two layers.

d) L-Layer Model

For even more convenience when implementing the L -layer Neural Net, you will need a function that replicates the previous one (`linear_activation_forward` with RELU) $L - 1$ times, then follows that with one `linear_activation_forward` with SIGMOID.

Implement the forward propagation of the above model.

In the code below, the variable AL will denote $A^{[L]} = \sigma(Z^{[L]}) = \sigma(W^{[L]}A^{[L-1]} + b^{[L]})$. (This is sometimes also called \hat{Y} , i.e., this is \hat{Y} .)

```
In [13]: def L_model_forward(X, parameters):
    """
        Implement forward propagation for the [LINEAR->RELU]*(L-1)->LINEAR->SIGMOID computation

    Arguments:
    X -- data, numpy array of shape (input size, number of examples)
    parameters -- output of initialize_parameters_deep()

    Returns:
    AL -- last post-activation value
    caches -- list of caches containing:
        every cache of linear_activation_forward() (there are L-1 of them, indexed from 0 to L-1)
    """
    caches = []
    A = X
    L = len(parameters) // 2 # number of layers in the neural network

    # Implement [LINEAR -> RELU]*(L-1). Add "cache" to the "caches" list.
    for l in range(1, L):
        A_prev = A
        W = parameters.get("W" + str(l))
        b = parameters.get("b" + str(l))
        A, cache = linear_activation_forward(A_prev, W, b, "relu")
        caches.append(cache)

    # Implement LINEAR -> SIGMOID. Add "cache" to the "caches" list.
    W = parameters.get("W" + str(L))
    b = parameters.get("b" + str(L))
    AL, cache = linear_activation_forward(A, W, b, "sigmoid")
    caches.append(cache)

    assert(AL.shape == (1,X.shape[1]))

    return AL, caches
```

```
In [14]: X, parameters = L_model_forward_test_case_2hidden()
AL, caches = L_model_forward(X, parameters)
print("AL = " + str(AL))
print("Length of caches list = " + str(len(caches)))
```

```
AL = [[ 0.03921668  0.70498921  0.19734387  0.04728177]]
Length of caches list = 3
```

AL	[[0.03921668 0.70498921 0.19734387 0.04728177]]
**Length of caches list **	3

Great! Now you have a full forward propagation that takes the input X and outputs a row vector $A^{[L]}$ containing your predictions. It also records all intermediate values in "caches". Using $A^{[L]}$, you can compute the cost of your predictions.

5 - Cost function

Compute the cross-entropy cost J , using the following formula:

$$-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)})) \quad (7)$$

In [15]: `def compute_cost(AL, Y):`

"""

Implement the cost function defined by equation (7).

Arguments:

AL -- probability vector corresponding to your Label predictions, shape (1, number of examples)

Y -- true "Label" vector (for example: containing 0 if non-cat, 1 if cat), shape (1, number of examples)

Returns:

cost -- cross-entropy cost

"""

`m = Y.shape[1]`

Compute Loss from aL and y.

`cost = (-1/m) * np.sum(np.multiply(Y,np.log(AL)) + np.multiply((1-Y),np.log(1-AL)))`

cost = np.squeeze(cost) # To make sure your cost's shape is what we expect (e.g. this turns [[17]] into 17).

`assert(cost.shape == ())`

`return cost`

In [16]: `Y, AL = compute_cost_test_case()`

`print("cost = " + str(compute_cost(AL, Y)))`

`cost = 0.414931599615`

Expected Output:

cost	0.41493159961539694
----------	---------------------

6 - Backward propagation module

Just like with forward propagation, you will implement helper functions for backpropagation. Remember that back propagation is used to calculate the gradient of the loss function with respect to the parameters.

$$\frac{d\mathcal{L}(a^{[2]}, y)}{dz^{[1]}} = \frac{d\mathcal{L}(a^{[2]}, y)}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{da^{[1]}} \frac{da^{[1]}}{dz^{[1]}} \quad (8)$$

In order to calculate the gradient $dW^{[1]} = \frac{\partial L}{\partial W^{[1]}}$, you use the previous chain rule and you do

$dW^{[1]} = dz^{[1]} \times \frac{\partial z^{[1]}}{\partial W^{[1]}}$. During the backpropagation, at each step you multiply your current gradient by the gradient corresponding to the specific layer to get the gradient you wanted.

Equivalently, in order to calculate the gradient $db^{[1]} = \frac{\partial L}{\partial b^{[1]}}$, you use the previous chain rule and you do

$$db^{[1]} = dz^{[1]} \times \frac{\partial z^{[1]}}{\partial b^{[1]}}.$$

This is why we talk about **backpropagation**. !-->

Now, similar to forward propagation, you are going to build the backward propagation in three steps:

- LINEAR backward
- LINEAR -> ACTIVATION backward where ACTIVATION computes the derivative of either the ReLU or sigmoid activation
- [LINEAR -> RELU] \times (L-1) -> LINEAR -> SIGMOID backward (whole model)

6.1 - Linear backward

For layer l , the linear part is: $Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}$ (followed by an activation).

Suppose you have already calculated the derivative $dZ^{[l]} = \frac{\partial L}{\partial Z^{[l]}}$. You want to get $(dW^{[l]}, db^{[l]}, dA^{[l-1]})$.

The three outputs $(dW^{[l]}, db^{[l]}, dA^{[l]})$ are computed using the input $dZ^{[l]}$. Here are the formulas you need:

$$dW^{[l]} = \frac{\partial \mathcal{L}}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]T} \quad (8)$$

$$db^{[l]} = \frac{\partial \mathcal{L}}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dZ^{[l](i)} \quad (9)$$

$$dA^{[l-1]} = \frac{\partial \mathcal{L}}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]} \quad (10)$$

Use the 3 formulas above to implement `linear_backward()`.

```
In [17]: def linear_backward(dZ, cache):
    """
        Implement the Linear portion of backward propagation for a single layer (layer L)

    Arguments:
    dZ -- Gradient of the cost with respect to the Linear output (of current layer L)
    cache -- tuple of values (A_prev, W, b) coming from the forward propagation in the current layer

    Returns:
    dA_prev -- Gradient of the cost with respect to the activation (of the previous layer L-1), same shape as A_prev
    dW -- Gradient of the cost with respect to W (current layer L), same shape as W
    db -- Gradient of the cost with respect to b (current layer L), same shape as b
    """
    A_prev, W, b = cache
    m = A_prev.shape[1]

    dW = (1/m) * np.dot(dZ, A_prev.T)
    db = (1/m) * np.sum(dZ, axis = 1, keepdims=True)
    dA_prev = np.dot(W.T, dZ)

    assert (dA_prev.shape == A_prev.shape)
    assert (dW.shape == W.shape)
    assert (db.shape == b.shape)

    return dA_prev, dW, db
```

```
In [18]: # Set up some test inputs
dZ, linear_cache = linear_backward_test_case()

dA_prev, dW, db = linear_backward(dZ, linear_cache)
print ("dA_prev = " + str(dA_prev))
print ("dW = " + str(dW))
print ("db = " + str(db))

dA_prev = [[ 0.51822968 -0.19517421]
           [-0.40506361  0.15255393]
           [ 2.37496825 -0.89445391]]
dW = [[-0.10076895  1.40685096  1.64992505]]
db = [[ 0.50629448]]
```

Expected Output:

dA_prev	[[0.51822968 -0.19517421] [-0.40506361 0.15255393] [2.37496825 -0.89445391]]
dW	[[[-0.10076895 1.40685096 1.64992505]]]
db	[[0.50629448]]

6.2 - Linear-Activation backward

Next, you will create a function that merges the two helper functions: `linear_backward` and the backward step for the activation `linear_activation_backward`.

To help you implement `linear_activation_backward`, we provided two backward functions:

- `sigmoid_backward`: Implements the backward propagation for SIGMOID unit. You can call it as follows:

```
dZ = sigmoid_backward(dA, activation_cache)
```

- `relu_backward`: Implements the backward propagation for RELU unit. You can call it as follows:

```
dZ = relu_backward(dA, activation_cache)
```

If $g(\cdot)$ is the activation function, `sigmoid_backward` and `relu_backward` compute

$$dZ^{[l]} = dA^{[l]} * g'(Z^{[l]}) \quad (11)$$

Implement the backpropagation for the *LINEAR->ACTIVATION* layer.

```
In [19]: def linear_activation_backward(dA, cache, activation):
    """
    Implement the backward propagation for the LINEAR->ACTIVATION layer.

    Arguments:
    dA -- post-activation gradient for current layer l
    cache -- tuple of values (linear_cache, activation_cache) we store for computing backward propagation efficiently
    activation -- the activation to be used in this layer, stored as a text string: "sigmoid" or "relu"

    Returns:
    dA_prev -- Gradient of the cost with respect to the activation (of the previous layer l-1), same shape as A_prev
    dW -- Gradient of the cost with respect to W (current layer l), same shape as W
    db -- Gradient of the cost with respect to b (current layer l), same shape as b
    """
    linear_cache, activation_cache = cache

    if activation == "relu":
        dZ = dZ = relu_backward(dA, activation_cache)
        dA_prev, dW, db = linear_backward(dZ, linear_cache)

    elif activation == "sigmoid":
        dZ = sigmoid_backward(dA, activation_cache)
        dA_prev, dW, db = linear_backward(dZ, linear_cache)

    return dA_prev, dW, db
```

```
In [20]: dAL, linear_activation_cache = linear_activation_backward_test_case()

dA_prev, dW, db = linear_activation_backward(dAL, linear_activation_cache, activation = "sigmoid")
print ("sigmoid:")
print ("dA_prev = " + str(dA_prev))
print ("dW = " + str(dW))
print ("db = " + str(db) + "\n")

dA_prev, dW, db = linear_activation_backward(dAL, linear_activation_cache, activation = "relu")
print ("relu:")
print ("dA_prev = " + str(dA_prev))
print ("dW = " + str(dW))
print ("db = " + str(db))

sigmoid:
dA_prev = [[ 0.11017994  0.01105339]
 [ 0.09466817  0.00949723]
 [-0.05743092 -0.00576154]]
dW = [[ 0.10266786  0.09778551 -0.01968084]]
db = [[-0.05729622]]

relu:
dA_prev = [[ 0.44090989  0.          ]
 [ 0.37883606  0.          ]
 [-0.2298228   0.          ]]
dW = [[ 0.44513824  0.37371418 -0.10478989]]
db = [[-0.20837892]]
```

Expected output with sigmoid:

dA_prev	[[0.11017994 0.01105339] [0.09466817 0.00949723] [-0.05743092 -0.00576154]]
dW	[[0.10266786 0.09778551 -0.01968084]]
db	[[-0.05729622]]

Expected output with relu:

dA_prev	[[0.44090989 0.] [0.37883606 0.] [-0.2298228 0.]]
dW	[[0.44513824 0.37371418 -0.10478989]]
db	[[-0.20837892]]

6.3 - L-Model Backward

Now you will implement the backward function for the whole network. Recall that when you implemented the `L_model_forward` function, at each iteration, you stored a cache which contains (X, W, b , and z). In the back propagation module, you will use those variables to compute the gradients. Therefore, in the `L_model_backward` function, you will iterate through all the hidden layers backward, starting from layer L . On each step, you will use the cached values for layer l to backpropagate through layer l . Figure 5 below shows the backward pass.

Initializing backpropagation: To backpropagate through this network, we know that the output is, $A^{[L]} = \sigma(Z^{[L]})$. Your code thus needs to compute $dAL = \frac{\partial \mathcal{L}}{\partial A^{[L]}}$. To do so, use this formula (derived using calculus which you don't need in-depth knowledge of):

```
dAL = - (np.divide(Y, AL) - np.divide(1 - Y, 1 - AL)) # derivative of cost with respect to AL
```

You can then use this post-activation gradient dAL to keep going backward. As seen in Figure 5, you can now feed in dAL into the LINEAR->SIGMOID backward function you implemented (which will use the cached values stored by the `L_model_forward` function). After that, you will have to use a `for` loop to iterate through all the other layers using the LINEAR->RELU backward function. You should store each dA , dW , and db in the `grads` dictionary. To do so, use this formula :

$$grads["dW^{[l]}"] = dW^{[l]} \quad (15)$$

For example, for $l = 3$ this would store $dW^{[l]}$ in `grads["dW3"]`.

Implement backpropagation for the $[LINEAR \rightarrow RELU] \times (L-1) \rightarrow LINEAR \rightarrow SIGMOID$ model.

```
In [21]: def L_model_backward(AL, Y, caches):
    """
        Implement the backward propagation for the [LINEAR->RELU] * (L-1) -> LINEAR
        -> SIGMOID group

        Arguments:
            AL -- probability vector, output of the forward propagation (L_model_forwa
            rd())
            Y -- true "label" vector (containing 0 if non-cat, 1 if cat)
            caches -- list of caches containing:
                every cache of linear_activation_forward() with "relu" (it's c
                aches[l], for l in range(L-1) i.e l = 0...L-2)
                the cache of linear_activation_forward() with "sigmoid" (it's
                caches[L-1])

        Returns:
            grads -- A dictionary with the gradients
                grads["dA" + str(l)] = ...
                grads["dW" + str(l)] = ...
                grads["db" + str(l)] = ...
        """
    grads = {}
    L = len(caches) # the number of layers
    m = AL.shape[1]
    Y = Y.reshape(AL.shape) # after this line, Y is the same shape as AL

    # Initializing the backpropagation
    dAL = - (np.divide(Y, AL) - np.divide(1 - Y, 1 - AL))

    # Lth Layer (SIGMOID -> LINEAR) gradients. Inputs: "dAL, current_cache". O
    utputs: "grads["dAL-1"], grads["dWL"], grads["dbL"]
    current_cache = caches[L-1]
    grads["dA" + str(L-1)], grads["dW" + str(L)], grads["db" + str(L)] = linea
    r_activation_backward(dAL, current_cache, "sigmoid")

    # Loop from L=L-2 to L=0
    for l in reversed(range(L-1)):
        # lth layer: (RELU -> LINEAR) gradients.
        # Inputs: "grads["dA" + str(l + 1)], current_cache". Outputs: "grads
        ["dA" + str(l)], grads["dW" + str(l + 1)], grads["db" + str(l + 1)]
        current_cache = caches[l]
        dA_prev_temp, dw_temp, db_temp = linear_activation_backward(grads.get(
            "dA" + str(l+1)), current_cache, "relu")
        grads["dA" + str(l)] = dA_prev_temp
        grads["dW" + str(l + 1)] = dw_temp
        grads["db" + str(l + 1)] = db_temp

    return grads
```

```
In [22]: AL, Y_assess, caches = L_model_backward_test_case()
grads = L_model_backward(AL, Y_assess, caches)
print_grads(grads)
```

```
dw1 = [[ 0.41010002  0.07807203  0.13798444  0.10502167]
       [ 0.          0.          0.          0.        ]
       [ 0.05283652  0.01005865  0.01777766  0.0135308 ]]
db1 = [[-0.22007063]
       [ 0.        ]
       [-0.02835349]]
dA1 = [[ 0.12913162 -0.44014127]
       [-0.14175655  0.48317296]
       [ 0.01663708 -0.05670698]]
```

Expected Output

dW1	[[0.41010002 0.07807203 0.13798444 0.10502167] [0. 0. 0.] [0.05283652 0.01005865 0.01777766 0.0135308]]
db1	[[-0.22007063] [0.] [-0.02835349]]
dA1	[[0.12913162 -0.44014127] [-0.14175655 0.48317296] [0.01663708 -0.05670698]]

6.4 - Update Parameters

In this section you will update the parameters of the model, using gradient descent:

$$W^{[l]} = W^{[l]} - \alpha dW^{[l]} \quad (16)$$

$$b^{[l]} = b^{[l]} - \alpha db^{[l]} \quad (17)$$

where α is the learning rate. After computing the updated parameters, store them in the parameters dictionary.

Implement `update_parameters()` to update your parameters using gradient descent.

Update parameters using gradient descent on every $W^{[l]}$ and $b^{[l]}$ for $l = 1, 2, \dots, L$

```
In [23]: def update_parameters(parameters, grads, learning_rate):
    """
        Update parameters using gradient descent

    Arguments:
        parameters -- python dictionary containing your parameters
        grads -- python dictionary containing your gradients, output of L_model_backward

    Returns:
        parameters -- python dictionary containing your updated parameters
        parameters["W" + str(l)] = ...
        parameters["b" + str(l)] = ...
    """
    L = len(parameters) // 2 # number of layers in the neural network

    # Update rule for each parameter. Use a for loop.
    for l in range(L):
        parameters["W" + str(l+1)] = parameters.get("W" + str(l+1)) - learning_rate*grads.get("dW" + str(l+1))
        parameters["b" + str(l+1)] = parameters.get("b" + str(l+1)) - learning_rate*grads.get("db" + str(l+1))

    return parameters
```

```
In [24]: parameters, grads = update_parameters_test_case()
parameters = update_parameters(parameters, grads, 0.1)
```

```
print ("W1 = "+ str(parameters["W1"]))
print ("b1 = "+ str(parameters["b1"]))
print ("W2 = "+ str(parameters["W2"]))
print ("b2 = "+ str(parameters["b2"]))
```

```
W1 = [[-0.59562069 -0.09991781 -2.14584584 1.82662008]
      [-1.76569676 -0.80627147 0.51115557 -1.18258802]
      [-1.0535704 -0.86128581 0.68284052 2.20374577]]
b1 = [[-0.04659241]
      [-1.28888275]
      [ 0.53405496]]
W2 = [[-0.55569196 0.0354055 1.32964895]]
b2 = [[-0.84610769]]
```

Expected Output:

W1	<pre>[[-0.59562069 -0.09991781 -2.14584584 1.82662008] [-1.76569676 -0.80627147 0.51115557 -1.18258802] [-1.0535704 -0.86128581 0.68284052 2.20374577]]</pre>
b1	<pre>[[-0.04659241] [-1.28888275] [0.53405496]]</pre>
W2	<pre>[[-0.55569196 0.0354055 1.32964895]]</pre>
b2	<pre>[[-0.84610769]]</pre>