

# Task01- Regularization

It is highly pivotal for deep learning models to be trained in a fashion that prevents issues like over-fitting, especially when the dataset is small. When a model over-fits, it fails to generalize well to new samples (unseen during training).

To avoid over-fitting, you will implement two techniques:

- L2-Regularization
- Dropout

```
In [31]: !git clone https://github.com/SanVik2000/EE5179-Final.git
```

```
Cloning into 'EE5179-Final'...
remote: Enumerating objects: 88, done.
remote: Counting objects: 100% (29/29), done.
remote: Compressing objects: 100% (27/27), done.
remote: Total 88 (delta 11), reused 0 (delta 0), pack-reused 59
Unpacking objects: 100% (88/88), done.
```

```
In [32]: !cp /content/EE5179-Final/Tutorial-3/reg_utils.py /content
!cp /content/EE5179-Final/Tutorial-3/testCases.py /content
```

```
In [33]: # import packages
import numpy as np
import matplotlib.pyplot as plt
from reg_utils import sigmoid, relu, plot_decision_boundary, initialize_parameters, load_2D_dataset, predict_dec
from reg_utils import compute_cost, predict, forward_propagation, backward_propagation, update_parameters
import sklearn
import sklearn.datasets
import scipy.io
from testCases import *

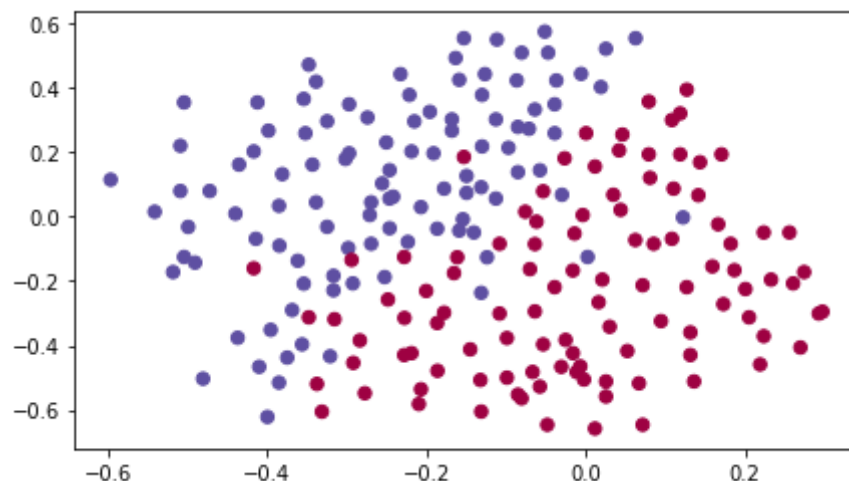
%matplotlib inline
plt.rcParams['figure.figsize'] = (7.0, 4.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'
```

```
In [36]: def load_2D_dataset():
data = scipy.io.loadmat('/content/EE5179-Final/Tutorial-3/data.mat')
train_X = data['X'].T
train_Y = data['y'].T
test_X = data['Xval'].T
test_Y = data['yval'].T

plt.scatter(train_X[0, :], train_X[1, :], c=train_Y, s=40, cmap=plt.cm.Spectral);

return train_X, train_Y, test_X, test_Y

train_X, train_Y, test_X, test_Y = load_2D_dataset()
```



Each dot corresponds to a data point, where

- blue dot represents class-1
- red dot represents class 2

**Your goal:** Use a deep learning model to classify the data points among two classes.

## 1 - Non-regularized model

You will use the following neural network (already implemented for you below). This model can be used:

- in *regularization mode* -- by setting the `lamdb` input to a non-zero value.
- in *dropout mode* -- by setting the `keep_prob` to a value less than one and greater than zero.

You will first try the model without any regularization. Then, you will implement:

- *L2 regularization* -- functions: `"compute_cost_with_regularization()"` and `"backward_propagation_with_regularization()"`
- *Dropout* -- functions: `"forward_propagation_with_dropout()"` and `"backward_propagation_with_dropout()"`

In each part, you will run this model with the correct inputs so that it calls the functions you've implemented. Take a look at the code below to familiarize yourself with the model.

```
In [37]: def model(X, Y, learning_rate = 0.3, num_iterations = 30000, print_cost = True, lamdb = 0, keep_prob = 1):
        """
        Implements a three-layer neural network: LINEAR->RELU->LINEAR->RELU->LINEAR->SIGMOID.

        Arguments:
        X -- input data, of shape (input size, number of examples)
        Y -- true "label" vector (1 for blue dot / 0 for red dot), of shape (output size, number of examples)
        learning_rate -- learning rate of the optimization
        num_iterations -- number of iterations of the optimization loop
        print_cost -- If True, print the cost every 10000 iterations
        lamdb -- regularization hyperparameter, scalar
        keep_prob - probability of keeping a neuron active during drop-out, scalar.

        Returns:
        parameters -- parameters learned by the model. They can then be used to predict.
        """

        grads = {}
        costs = []                # to keep track of the cost
        m = X.shape[1]           # number of examples
        layers_dims = [X.shape[0], 20, 3, 1]

        # Initialize parameters dictionary.
        parameters = initialize_parameters(layers_dims)

        # Loop (gradient descent)

        for i in range(0, num_iterations):
```

```

# Forward propagation: LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID.
if keep_prob == 1:
    a3, cache = forward_propagation(X, parameters)
elif keep_prob < 1:
    a3, cache = forward_propagation_with_dropout(X, parameters, keep_prob)

# Cost function
if lambd == 0:
    cost = compute_cost(a3, Y)
else:
    cost = compute_cost_with_regularization(a3, Y, parameters, lambd)

# Backward propagation.
assert(lambd==0 or keep_prob==1)    # it is possible to use both L2 regularization and dropout,
                                    # but this assignment will only explore one at a time

if lambd == 0 and keep_prob == 1:
    grads = backward_propagation(X, Y, cache)
elif lambd != 0:
    grads = backward_propagation_with_regularization(X, Y, cache, lambd)
elif keep_prob < 1:
    grads = backward_propagation_with_dropout(X, Y, cache, keep_prob)

# Update parameters.
parameters = update_parameters(parameters, grads, learning_rate)

# Print the loss every 10000 iterations
if print_cost and i % 10000 == 0:
    print("Cost after iteration {}: {}".format(i, cost))
if print_cost and i % 1000 == 0:
    costs.append(cost)

# plot the cost
plt.plot(costs)
plt.ylabel('cost')
plt.xlabel('iterations (x1,000)')
plt.title("Learning rate =" + str(learning_rate))
plt.show()

return parameters

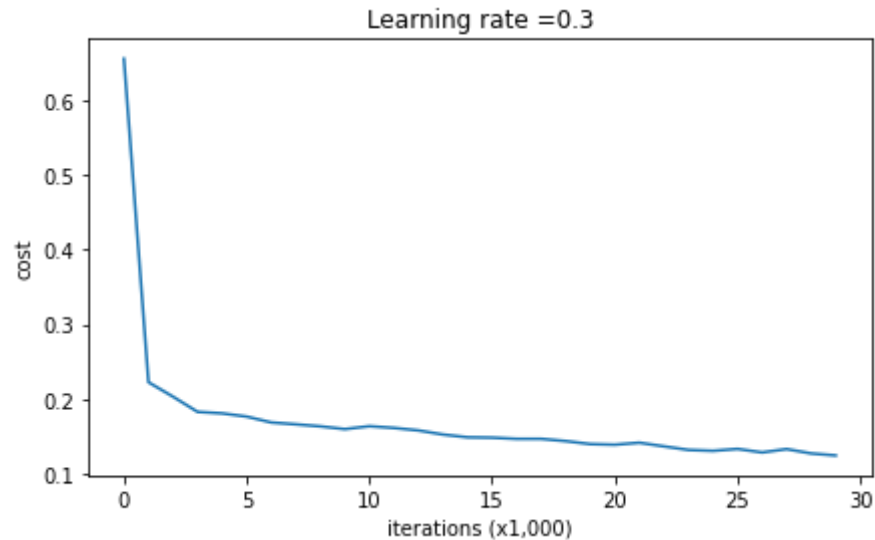
```

Let's train the model without any regularization, and observe the accuracy on the train/test sets.

In [38]: `parameters = model(train_X, train_Y)`

```
print ("On the training set:")
predictions_train = predict(train_X, train_Y, parameters)
print ("On the test set:")
predictions_test = predict(test_X, test_Y, parameters)
```

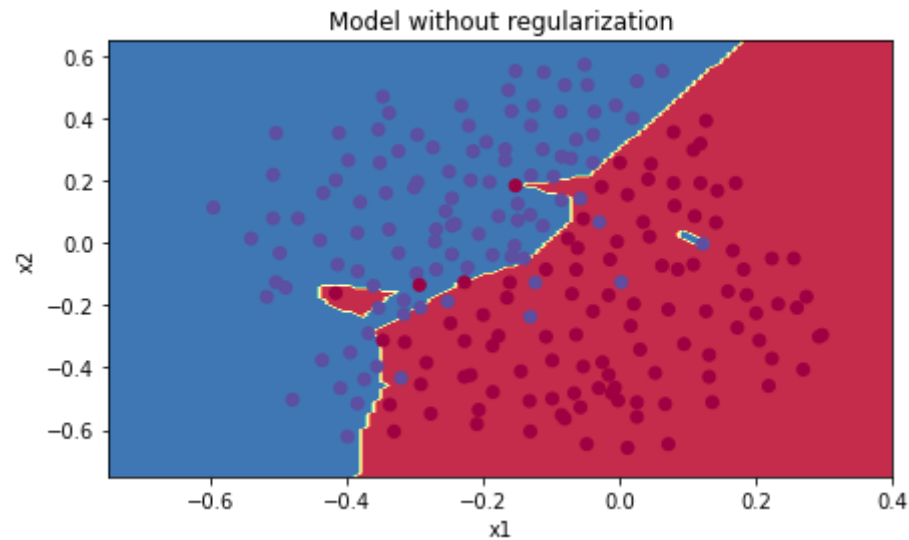
Cost after iteration 0: 0.6557412523481002  
 Cost after iteration 10000: 0.1632998752572417  
 Cost after iteration 20000: 0.13851642423284755



On the training set:  
 Accuracy: 0.9478672985781991  
 On the test set:  
 Accuracy: 0.915

The train accuracy is 94.8% while the test accuracy is 91.5%. This is the **baseline model** (you will observe the impact of regularization on this model). Run the following code to plot the decision boundary of your model.

```
In [39]: plt.title("Model without regularization")
axes = plt.gca()
axes.set_xlim([-0.75,0.40])
axes.set_ylim([-0.75,0.65])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```



The non-regularized model is obviously overfitting the training set. It is fitting the noisy points! Lets now look at two techniques to reduce overfitting.

## 2 - L2 Regularization

The standard way to avoid overfitting is called **L2 regularization**. The new modified cost function now becomes:

$$J_{regularized} = \underbrace{-\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}) \right)}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{k,j}^{[l]2}}_{\text{L2 regularization cost}} \quad (2)$$

**Exercise:** Implement `compute_cost_with_regularization()` which computes the cost given by formula (2). To calculate  $\sum_k \sum_j W_{k,j}^{[l]2}$ , use :

```
np.sum(np.square(Wl))
```

Note that you have to do this for  $W^{[1]}$ ,  $W^{[2]}$  and  $W^{[3]}$ , then sum the three terms and multiply by  $\frac{1}{m} \frac{\lambda}{2}$ .

```
In [40]: # GRADED FUNCTION: compute_cost_with_regularization
```

```
def compute_cost_with_regularization(A3, Y, parameters, lambd):
    """
    Implement the cost function with L2 regularization. See formula (2) above.

    Arguments:
    A3 -- post-activation, output of forward propagation, of shape (output size, number of examples)
    Y -- "true" labels vector, of shape (output size, number of examples)
    parameters -- python dictionary containing parameters of the model

    Returns:
    cost - value of the regularized loss function (formula (2))
    """
    m = Y.shape[1]
    W1 = parameters["W1"]
    W2 = parameters["W2"]
    W3 = parameters["W3"]

    cross_entropy_cost = compute_cost(A3, Y) # This gives you the cross-entropy part of the cost

    ### START CODE HERE ### (approx. 1 line)
    L2_regularization_cost = lambd/(2*m)*(np.sum(np.square(W1))+np.sum(np.square(W2))+np.sum(np.square(W3)))
    ### END CODER HERE ###

    cost = cross_entropy_cost + L2_regularization_cost

    return cost
```

```
In [41]: A3, Y_assess, parameters = compute_cost_with_regularization_test_case()

print("cost = " + str(compute_cost_with_regularization(A3, Y_assess, parameters, lambd = 0.1)))
```

```
cost = 1.7864859451590758
```

**Expected Output:**

```
**cost** 1.78648594516
```

Since the cost function is changed now, we have to also include the effects of the gradients of this additional term in our gradient descent algorithm.

**Exercise:** Implement the changes needed in backward propagation to take into account regularization. The changes only concern  $dW1$ ,  $dW2$  and  $dW3$ . For each, you have to add the regularization term's gradient ( $\frac{d}{dW}(\frac{1}{2} \frac{\lambda}{m} W^2) = \frac{\lambda}{m} W$ ).

```
In [42]: # GRADED FUNCTION: backward_propagation_with_regularization

def backward_propagation_with_regularization(X, Y, cache, lambd):
    """
    Implements the backward propagation of our baseline model to which we added an L2 regularization.

    Arguments:
    X -- input dataset, of shape (input size, number of examples)
    Y -- "true" labels vector, of shape (output size, number of examples)
    cache -- cache output from forward_propagation()
    lambd -- regularization hyperparameter, scalar

    Returns:
    gradients -- A dictionary with the gradients with respect to each parameter, activation and pre-activation variable
    """

    m = X.shape[1]
    (Z1, A1, W1, b1, Z2, A2, W2, b2, Z3, A3, W3, b3) = cache

    dZ3 = A3 - Y

    ### START CODE HERE ### (approx. 1 line)
    dW3 = 1./m * np.dot(dZ3, A2.T) + (lambd/m)*W3
    ### END CODE HERE ###
    db3 = 1./m * np.sum(dZ3, axis=1, keepdims = True)

    dA2 = np.dot(W3.T, dZ3)
    dZ2 = np.multiply(dA2, np.int64(A2 > 0))
    ### START CODE HERE ### (approx. 1 line)
    dW2 = 1./m * np.dot(dZ2, A1.T) + (lambd/m)*W2
    ### END CODE HERE ###
    db2 = 1./m * np.sum(dZ2, axis=1, keepdims = True)

    dA1 = np.dot(W2.T, dZ2)
    dZ1 = np.multiply(dA1, np.int64(A1 > 0))
    ### START CODE HERE ### (approx. 1 line)
    dW1 = 1./m * np.dot(dZ1, X.T) + (lambd/m)*W1
    ### END CODE HERE ###
    db1 = 1./m * np.sum(dZ1, axis=1, keepdims = True)
```



```

gradients = {"dZ3": dZ3, "dW3": dW3, "db3": db3, "dA2": dA2,
             "dZ2": dZ2, "dW2": dW2, "db2": db2, "dA1": dA1,
             "dZ1": dZ1, "dW1": dW1, "db1": db1}

return gradients

```

```

In [43]: X_assess, Y_assess, cache = backward_propagation_with_regularization_test_case()

grads = backward_propagation_with_regularization(X_assess, Y_assess, cache, lambd = 0.7)
print ("dW1 = \n"+ str(grads["dW1"]))
print ("dW2 = \n"+ str(grads["dW2"]))
print ("dW3 = \n"+ str(grads["dW3"]))

dW1 =
[[-0.25604646  0.12298827 -0.28297129]
 [-0.17706303  0.34536094 -0.4410571 ]]
dW2 =
[[ 0.79276486  0.85133918]
 [-0.0957219  -0.01720463]
 [-0.13100772 -0.03750433]]
dW3 =
[[-1.77691347 -0.11832879 -0.09397446]]

```

#### Expected Output:

```

dW1 =
[[-0.25604646  0.12298827 -0.28297129]
 [-0.17706303  0.34536094 -0.4410571 ]]
dW2 =
[[ 0.79276486  0.85133918]
 [-0.0957219  -0.01720463]
 [-0.13100772 -0.03750433]]
dW3 =
[[-1.77691347 -0.11832879 -0.09397446]]

```

Let's now run the model with L2 regularization ( $\lambda = 0.7$ ). The `model()` function will call:

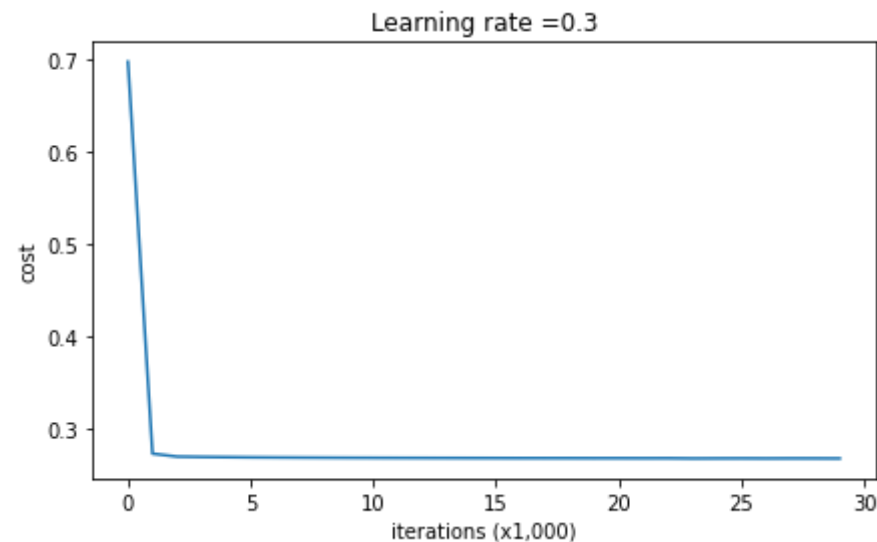
- `compute_cost_with_regularization` instead of `compute_cost`
- `backward_propagation_with_regularization` instead of `backward_propagation`

```
In [44]: parameters = model(train_X, train_Y, lambd = 0.7)
print ("On the train set:")
predictions_train = predict(train_X, train_Y, parameters)
print ("On the test set:")
predictions_test = predict(test_X, test_Y, parameters)
```

Cost after iteration 0: 0.6974484493131264

Cost after iteration 10000: 0.2684918873282239

Cost after iteration 20000: 0.2680916337127301



On the train set:

Accuracy: 0.9383886255924171

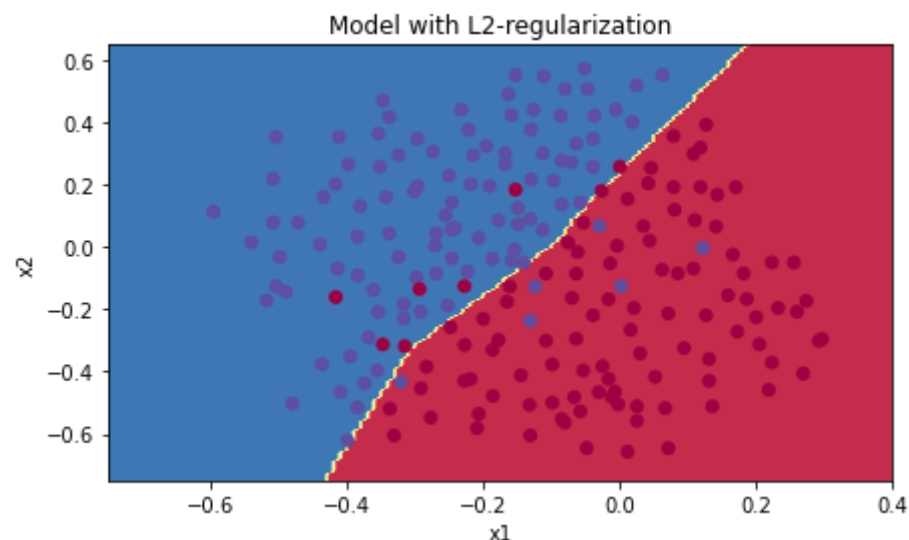
On the test set:

Accuracy: 0.93

The test set accuracy has now increased to 93%.

You are not overfitting the training data anymore. Let's plot the decision boundary.

```
In [45]: plt.title("Model with L2-regularization")
axes = plt.gca()
axes.set_xlim([-0.75,0.40])
axes.set_ylim([-0.75,0.65])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```



### Observations:

- The value of  $\lambda$  is a hyperparameter that you can tune using a validation set.
- L2 regularization makes your decision boundary smoother. If  $\lambda$  is too large, it is also possible to "oversmooth", resulting in a model with high bias.

### What is L2-regularization actually doing?:

L2-regularization relies on the assumption that a model with small weights is simpler than a model with large weights. Thus, by penalizing the square values of the weights in the cost function you drive all the weights to smaller values. It becomes too costly for the cost to have large weights! This leads to a smoother model in which the output changes more slowly as the input changes.

## 3 - Dropout

**Dropout** is a widely used regularization technique that is specific to deep learning. **It randomly shuts down some neurons in each iteration.**

When you shut some neurons down, you actually modify your model. The idea behind drop-out is that at each iteration, you train a different model that uses only a subset of your neurons. With dropout, your neurons thus become less sensitive to the activation of one other specific neuron, because that other neuron might be shut down at any time.

### 3.1 - Forward propagation with dropout

**Exercise:** Implement the forward propagation with dropout. You are using a 3 layer neural network, and will add dropout to the first and second hidden layers. We will not apply dropout to the input layer or output layer.

**Hint:** Let's say that `keep_prob = 0.8`, which means that we want to keep about 80% of the neurons and drop out about 20% of them. We want to generate a vector that has 1's and 0's, where about 80% of them are 1 and about 20% are 0. This python statement:

```
X = (X < keep_prob).astype(int)
```

is conceptually the same as this if-else statement (for the simple case of a one-dimensional array) :

```
for i,v in enumerate(x):
    if v < keep_prob:
        x[i] = 1
    else: # v >= keep_prob
        x[i] = 0
```

Note that the `X = (X < keep_prob).astype(int)` works with multi-dimensional arrays, and the resulting output preserves the dimensions of the input array.

Also note that without using `.astype(int)`, the result is an array of booleans `True` and `False`, which Python automatically converts to 1 and 0 if we multiply it with numbers. (However, it's better practice to convert data into the data type that we intend, so try using `.astype(int)`.)

**Steps to be followed:**

1. Set  $A^{[1]}$  to  $A^{[1]} * D^{[1]}$ . (You are shutting down some neurons). You can think of  $D^{[1]}$  as a mask, so that when it is multiplied with another matrix, it shuts down some of the values.
2. Divide  $A^{[1]}$  by `keep_prob`. By doing this you are assuring that the result of the cost will still have the same expected value as without drop-out. (This technique is also called inverted dropout.)

```
In [46]: # GRADED FUNCTION: forward_propagation_with_dropout
```

```
def forward_propagation_with_dropout(X, parameters, keep_prob = 0.5):
    """
    Implements the forward propagation: LINEAR -> RELU + DROPOUT -> LINEAR -> RELU + DROPOUT -> LINEAR -> SIGMOID.
```

**Arguments:**

X -- input dataset, of shape (2, number of examples)  
 parameters -- python dictionary containing your parameters "W1", "b1", "W2", "b2", "W3", "b3":  
     W1 -- weight matrix of shape (20, 2)  
     b1 -- bias vector of shape (20, 1)  
     W2 -- weight matrix of shape (3, 20)  
     b2 -- bias vector of shape (3, 1)  
     W3 -- weight matrix of shape (1, 3)  
     b3 -- bias vector of shape (1, 1)  
 keep\_prob - probability of keeping a neuron active during drop-out, scalar

**Returns:**

A3 -- last activation value, output of the forward propagation, of shape (1,1)  
 cache -- tuple, information stored for computing the backward propagation  
 """

```
np.random.seed(1)
```

```
# retrieve parameters
```

```
W1 = parameters["W1"]
```

```
b1 = parameters["b1"]
```

```
W2 = parameters["W2"]
```

```
b2 = parameters["b2"]
```

```
W3 = parameters["W3"]
```

```
b3 = parameters["b3"]
```

```
# LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID
```

```
Z1 = np.dot(W1, X) + b1
```

```
A1 = relu(Z1)
```

```
### START CODE HERE ### (approx. 4 lines)
```

```
D1 = np.random.rand(A1.shape[0], A1.shape[1])
```

```
D1 = (D1 < keep_prob).astype(int)
```

```
A1 = A1*D1
```

```
A1 = A1/keep_prob
```

```
### END CODE HERE ###
```

```
Z2 = np.dot(W2, A1) + b2
```

```
A2 = relu(Z2)
```

```
### START CODE HERE ### (approx. 4 lines)
```

```
D2 = np.random.rand(A2.shape[0], A2.shape[1])
```

```
D2 = (D2 < keep_prob)
```

```
A2 = A2*D2
```

```
A2 = A2/keep_prob
```

```
### END CODE HERE ###
```

```
# Steps 1-4 below correspond to the Steps 1-4 described above.
```

```
# Step 1: initialize matrix D1 = np.random.rand(..., ...)
```

```
# Step 2: convert entries of D1 to 0 or 1 (using keep_prob as threshold)
```

```
# Step 3: shut down some neurons of A1
```

```
# Step 4: scale the value of neurons that haven't been shut down
```

```
# Step 1: initialize matrix D2 = np.random.rand(..., ...)
```

```
# Step 2: convert entries of D2 to 0 or 1 (using keep_prob as threshold)
```

```
# Step 3: shut down some neurons of A2
```

```
# Step 4: scale the value of neurons that haven't been shut down
```

```

Z3 = np.dot(W3, A2) + b3
A3 = sigmoid(Z3)

cache = (Z1, D1, A1, W1, b1, Z2, D2, A2, W2, b2, Z3, A3, W3, b3)

return A3, cache

```

```

In [47]: X_assess, parameters = forward_propagation_with_dropout_test_case()

A3, cache = forward_propagation_with_dropout(X_assess, parameters, keep_prob = 0.7)
print ("A3 = " + str(A3))

A3 = [[0.36974721 0.00305176 0.04565099 0.49683389 0.36974721]]

```

**Expected Output:**

```

**A3** [[ 0.36974721 0.00305176 0.04565099 0.49683389 0.36974721]]

```

## 3.2 - Backward propagation with dropout

**Exercise:** Implement the backward propagation with dropout. As before, you are training a 3 layer network. Add dropout to the first and second hidden layers, using the masks  $D^{[1]}$  and  $D^{[2]}$  stored in the cache.

**Instruction:** Backpropagation with dropout is actually quite easy. You will have to carry out 2 Steps:

1. You had previously shut down some neurons during forward propagation, by applying a mask  $D^{[1]}$  to  $A1$ . In backpropagation, you will have to shut down the same neurons, by reapplying the same mask  $D^{[1]}$  to  $dA1$ .
2. During forward propagation, you had divided  $A1$  by `keep_prob`. In backpropagation, you'll therefore have to divide  $dA1$  by `keep_prob` again (the calculus interpretation is that if  $A^{[1]}$  is scaled by `keep_prob`, then its derivative  $dA^{[1]}$  is also scaled by the same `keep_prob`).

```

In [48]: # GRADED FUNCTION: backward_propagation_with_dropout

def backward_propagation_with_dropout(X, Y, cache, keep_prob):
    """
    Implements the backward propagation of our baseline model to which we added dropout.

    Arguments:
    X -- input dataset, of shape (2, number of examples)
    """

```

```

Y -- "true" labels vector, of shape (output size, number of examples)
cache -- cache output from forward_propagation_with_dropout()
keep_prob - probability of keeping a neuron active during drop-out, scalar

Returns:
gradients -- A dictionary with the gradients with respect to each parameter, activation and pre-activation variables
"""

m = X.shape[1]
(Z1, D1, A1, W1, b1, Z2, D2, A2, W2, b2, Z3, A3, W3, b3) = cache

dZ3 = A3 - Y
dW3 = 1./m * np.dot(dZ3, A2.T)
db3 = 1./m * np.sum(dZ3, axis=1, keepdims = True)
dA2 = np.dot(W3.T, dZ3)
### START CODE HERE ### (~ 2 lines of code)
dA2 = dA2*D2          # Step 1: Apply mask D2 to shut down the same neurons as during the forward propagation
dA2 = dA2/keep_prob    # Step 2: Scale the value of neurons that haven't been shut down
### END CODE HERE ###
dZ2 = np.multiply(dA2, np.int64(A2 > 0))
dW2 = 1./m * np.dot(dZ2, A1.T)
db2 = 1./m * np.sum(dZ2, axis=1, keepdims = True)

dA1 = np.dot(W2.T, dZ2)
### START CODE HERE ### (~ 2 lines of code)
dA1 = dA1*D1          # Step 1: Apply mask D1 to shut down the same neurons as during the forward propagation
dA1 = dA1/keep_prob    # Step 2: Scale the value of neurons that haven't been shut down
### END CODE HERE ###
dZ1 = np.multiply(dA1, np.int64(A1 > 0))
dW1 = 1./m * np.dot(dZ1, X.T)
db1 = 1./m * np.sum(dZ1, axis=1, keepdims = True)

gradients = {"dZ3": dZ3, "dW3": dW3, "db3": db3, "dA2": dA2,
            "dZ2": dZ2, "dW2": dW2, "db2": db2, "dA1": dA1,
            "dZ1": dZ1, "dW1": dW1, "db1": db1}

return gradients

```

```

In [49]: X_assess, Y_assess, cache = backward_propagation_with_dropout_test_case()

gradients = backward_propagation_with_dropout(X_assess, Y_assess, cache, keep_prob = 0.8)

```

```
print ("dA1 = \n" + str(gradients["dA1"]))
print ("dA2 = \n" + str(gradients["dA2"]))
```

```
dA1 =
[[ 0.36544439  0.          -0.00188233  0.          -0.17408748]
 [ 0.65515713  0.          -0.00337459  0.          -0.          ]]
dA2 =
[[ 0.58180856  0.          -0.00299679  0.          -0.27715731]
 [ 0.          0.53159854 -0.          0.53159854 -0.34089673]
 [ 0.          0.          -0.00292733  0.          -0.          ]]
```

### Expected Output:

```
dA1 =
[[ 0.36544439  0.          -0.00188233  0.          -0.17408748]
 [ 0.65515713  0.          -0.00337459  0.          -0.          ]]
dA2 =
[[ 0.58180856  0.          -0.00299679  0.          -0.27715731]
 [ 0.          0.53159854 -0.          0.53159854 -0.34089673]
 [ 0.          0.          -0.00292733  0.          -0.          ]]
```

Let's now run the model with dropout ( `keep_prob = 0.86` ). It means at every iteration you shut down each neurons of layer 1 and 2 with 14% probability. The function `model()` will now call:

- `forward_propagation_with_dropout` instead of `forward_propagation`.
- `backward_propagation_with_dropout` instead of `backward_propagation`.

```
In [50]: parameters = model(train_X, train_Y, keep_prob = 0.86, learning_rate = 0.3)
```

```
print ("On the train set:")
predictions_train = predict(train_X, train_Y, parameters)
print ("On the test set:")
predictions_test = predict(test_X, test_Y, parameters)
```

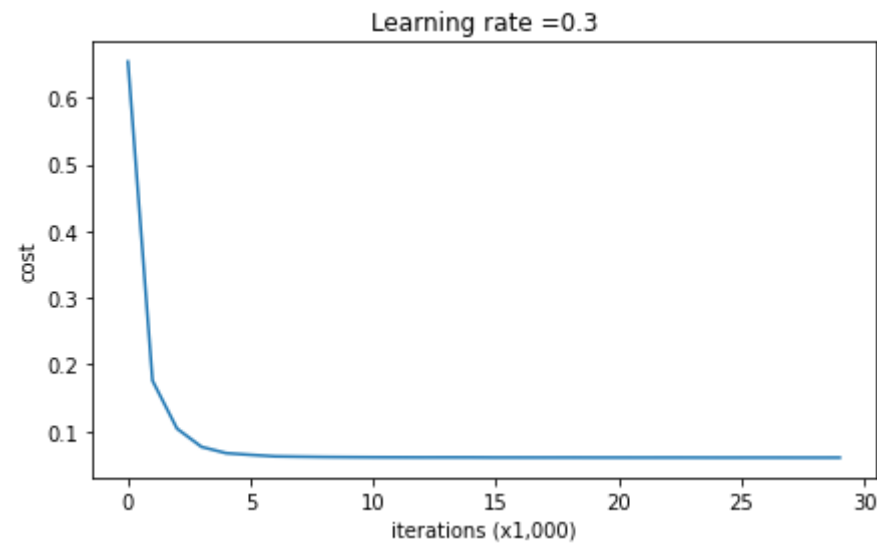
Cost after iteration 0: 0.6543912405149825

```
/content/reg_utils.py:236: RuntimeWarning: divide by zero encountered in log
  logprobs = np.multiply(-np.log(a3),Y) + np.multiply(-np.log(1 - a3), 1 - Y)
/content/reg_utils.py:236: RuntimeWarning: invalid value encountered in multiply
  logprobs = np.multiply(-np.log(a3),Y) + np.multiply(-np.log(1 - a3), 1 - Y)
```

Cost after iteration 10000: 0.0610169865749056

Cost after iteration 20000: 0.060582435798513114





On the train set:

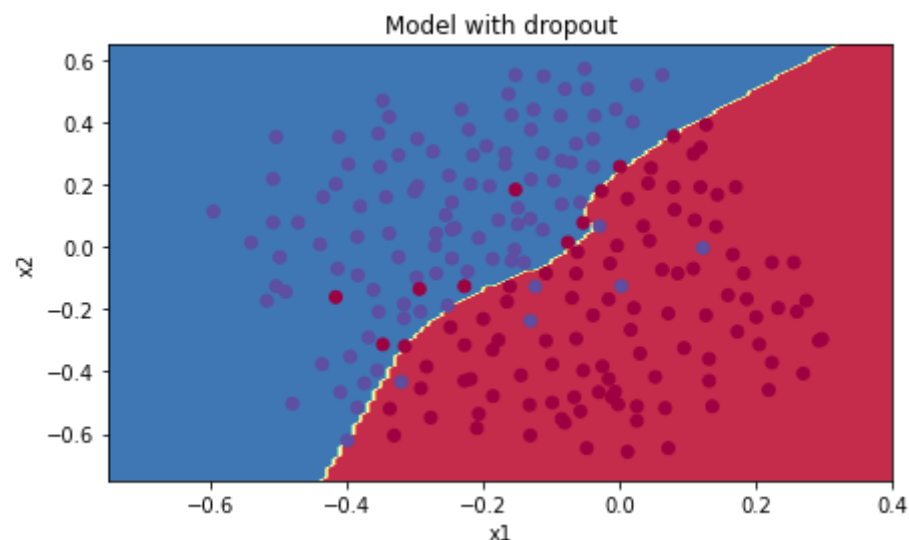
Accuracy: 0.9289099526066351

On the test set:

Accuracy: 0.95

The test accuracy has increased again (to 95%)! Your model is not overfitting the training set and performs well on the test set.

```
In [51]: plt.title("Model with dropout")
axes = plt.gca()
axes.set_xlim([-0.75,0.40])
axes.set_ylim([-0.75,0.65])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

**Note:**

- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.

## 4 - Conclusions

Here are the results of our three models:

<b>**Model**</b>	<b>**Train Accuracy**</b>	<b>**Test Accuracy**</b>
3-layer NN without regularization	95	91.5
3-layer NN with L2-regularization	94	93
3-layer NN with Dropout regularization	93	95

Note that regularization hurts training set performance! This is because it limits the ability of the network to overfit to the training set. But since it ultimately gives better test accuracy, it is helping your system.