

Regularization

Welcome to the second assignment of this week. Deep Learning models have so much flexibility and capacity that **overfitting can be a serious problem**, if the training dataset is not big enough. Sure it does well on the training set, but the learned network **doesn't generalize to few examples** that it has never seen!

You will learn to: Use regularization in your deep learning models.

Let's get started!

Table of Contents

- 1 - Packages
- 2 - Problem Statement
- 3 - Loading the Dataset
- 4 - Non-Regularized Model
- 5 - L2 Regularization
 - Exercise 1 - compute_cost_with_regularization
 - Exercise 2 - backward_propagation_with_regularization
- 6 - Dropout
 - 6.1 - Forward Propagation with Dropout
 - Exercise 3 - forward_propagation_with_dropout
 - 6.2 - Backward Propagation with Dropout
 - Exercise 4 - backward_propagation_with_dropout
- 7 - Conclusions

1 - Packages

```
In [4]: # import packages
import numpy as np
import matplotlib.pyplot as plt
import sklearn
import sklearn.datasets
import scipy.io

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
from testCases import *
from public_tests import *

%matplotlib inline
plt.rcParams['figure.figsize'] = (7, 4) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

%load_ext autoreload
%autoreload 2
```

2 - Problem Statement

You have just been hired as an expert by the French Football Corporation. They would like you to recommend positions where France's goal keeper should kick the ball so that the French team's players can then hit it with their head.

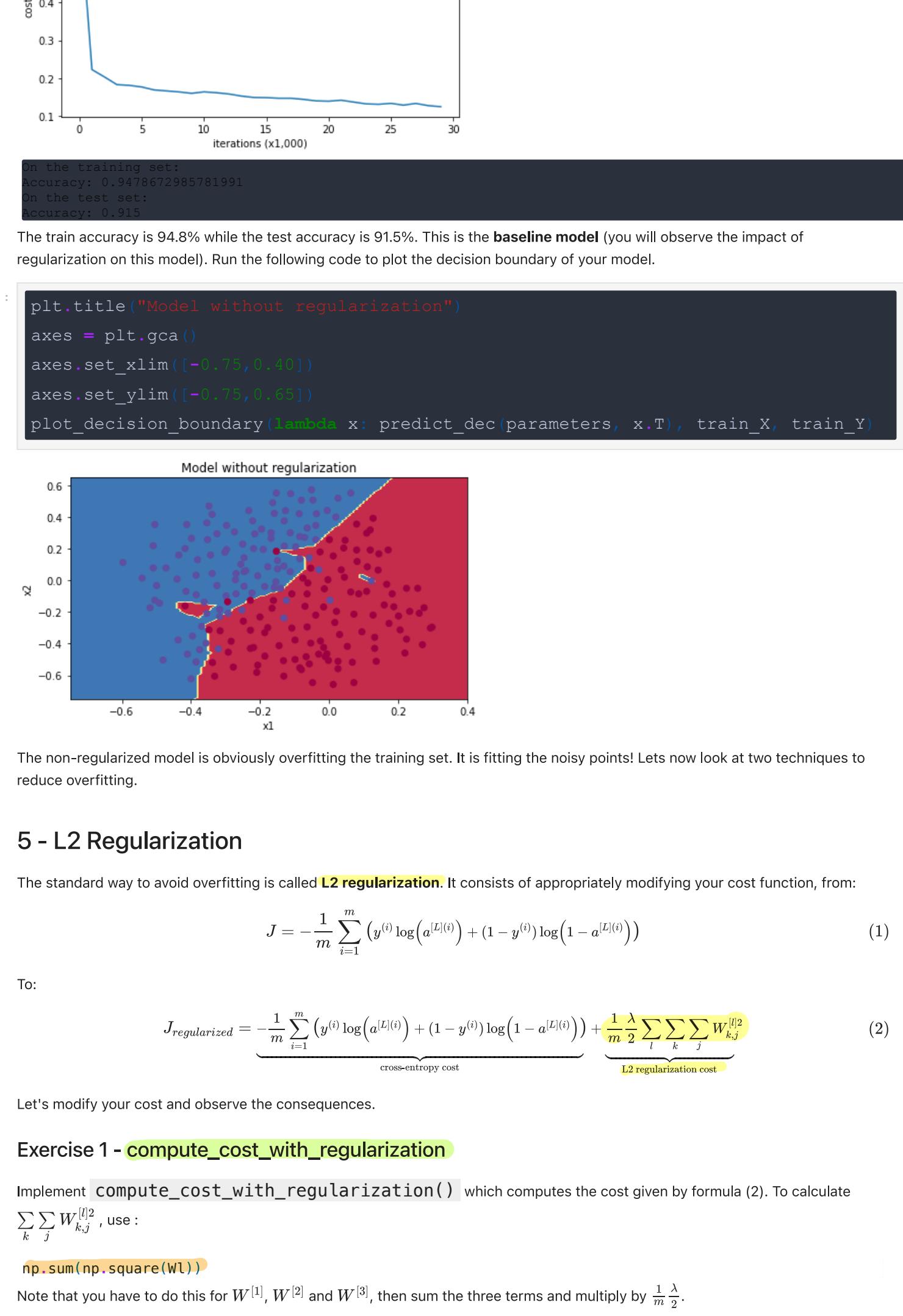
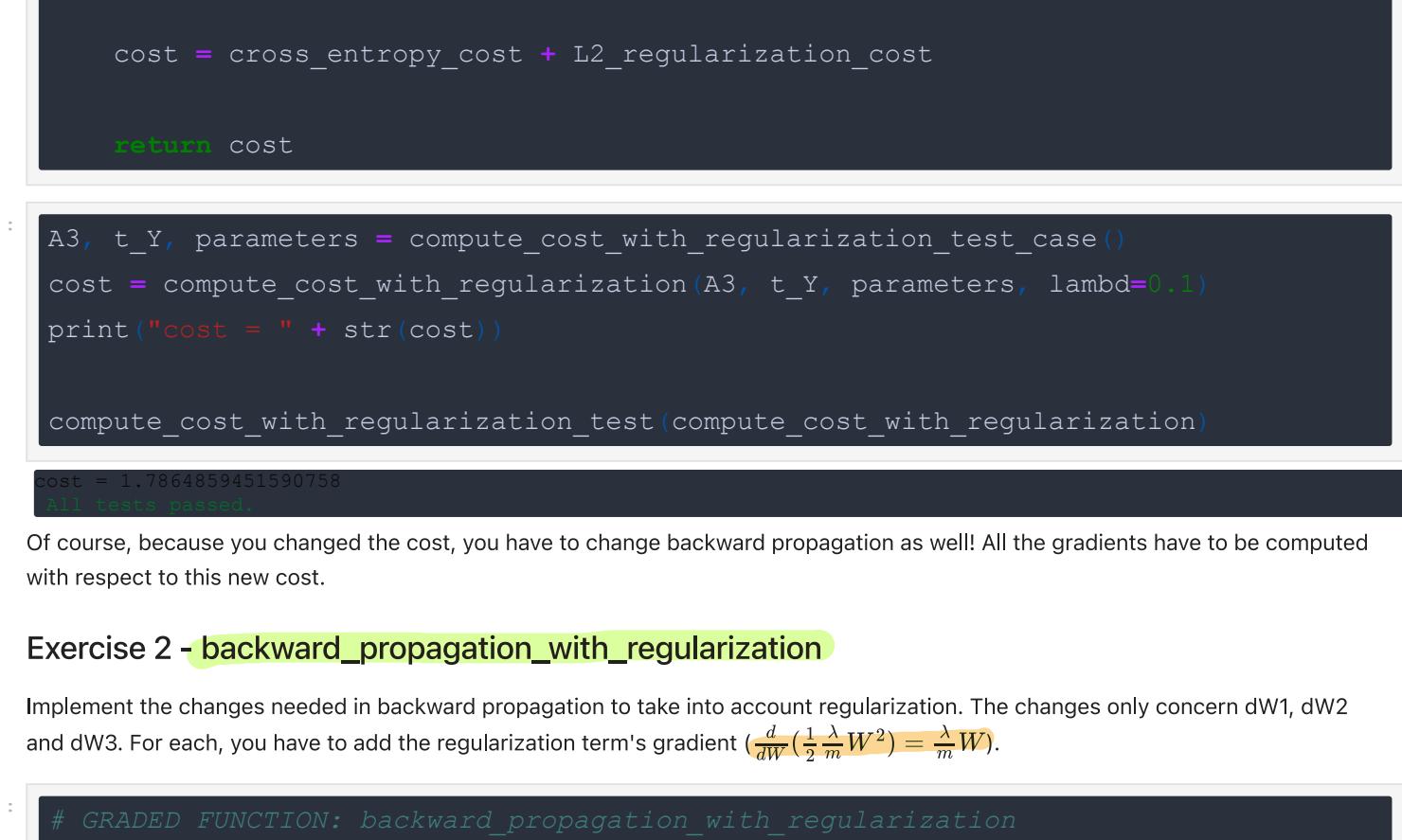


Figure 1: football field. The goal keeper kicks the ball in the air, the players of each team are fighting to hit the ball with their head. They give you the following 2D dataset from France's past 10 games.

3 - Loading the Dataset

```
In [5]: train_X, train_Y, test_X, test_Y = load_2D_dataset()
```



Each dot corresponds to a position on the football field where a football player has hit the ball with his/her head after the French goal keeper has shot the ball from the left side of the football field.

- If the dot is blue, it means the French player managed to hit the ball with his/her head
- If the dot is red, it means the other team's player hit the ball with his/her head

Your goal: Use a deep learning model to find the positions on the field where the goalkeeper should kick the ball.

Analysis of the dataset: This dataset is a little noisy, but it looks like a diagonal line separating the upper left half (blue) from the lower right half (red) would work well.

You will first try a non-regularized model. Then you'll learn how to regularize it and decide which model you will choose to solve the French Football Corporation's problem.

4 - 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 lambda input to a non-zero value. We use "lambda" instead of "lambd" because "lambd" is a reserved keyword in Python.
```

```
• In dropout mode -- by setting the keep_prob to a value less than one
```

```
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 [6]: # GRADED FUNCTION: model
def model(X, Y, learning_rate = 0.05, num_iterations = 10000, print_cost = True,
         lambd = 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
    lambd -- 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 = []
    m = X.shape[1] # number of examples
    layers_dims = [X.shape[0], 20, 1] # 3-layer model

    # 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.
        a1, cache1 = forward_propagation(X, parameters)
        a2, cache2 = forward_propagation(a1, parameters)
        a3, cache3 = forward_propagation(a2, parameters)

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

        # Backward propagation.
        grads["d1"] = backward_propagation(X, Y, cache1)
        grads["d2"] = backward_propagation(a1, cache2, lambd)
        grads["d3"] = backward_propagation(a2, cache3, lambd)

        # 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))
        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 [7]: # GRADED FUNCTION: predict
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: 0.65074105348105
cost after iteration 10000: 0.108510305353574204
cost after iteration 20000: 0.1185164243234922

Learning rate=0.3

accuracy: 0.947867298571991
the best set:

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 [8]: plt.title("Model without regularization")
axes = plt.gca()
axes.set_xlim([-1.0, 1.0])
axes.set_ylim([-1.0, 1.0])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Model without regularization

```
In [17]: # 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)

    # 4 lines of code)           # Steps 1-4 below correspond to the Steps 1-4
described above.
    # D1 =
    # Step 1: initialize matrix
D1 = np.random.rand(..., ...)

    # D1 =
    # Step 2: convert entries of
D1 = 0 or 1 (using keep_prob as the threshold)
    # A1 =
    # Step 3: shut down some
neurons of A1
    # A1 =
    # Step 4: scale the value of
neurons that haven't been shut down
    # YOUR CODE STARTS HERE
    D1 = np.random.rand(A1.shape[0]), A1.shape[1])
    D1 = (D1 < keep_prob).astype int
    A1 = np.multiply(A1, D1)
    A1 = A1 / keep_prob
    # YOUR CODE ENDS HERE
    Z2 = np.dot(W2, A1) + b2
    A2 = relu(Z2)

    # 4 lines of code)
    # D2 =
    # Step 1: initialize matrix
D2 = np.random.rand(..., ...)

    # D2 =
    # Step 2: convert entries of
D2 = 0 or 1 (using keep_prob as the threshold)
    # A2 =
    # Step 3: shut down some
neurons of A2
    # A2 =
    # Step 4: scale the value of
neurons that haven't been shut down
    # YOUR CODE STARTS HERE
    D2 = np.random.rand(A2.shape[0]), A2.shape[1])
    D2 = (D2 < keep_prob).astype int
    A2 = np.multiply(A2, D2)
    A2 = A2 / keep_prob
    # YOUR CODE ENDS HERE
    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 [18]: t_X, parameters = forward_propagation_with_dropout_test_case()

A3, cache = forward_propagation_with_dropout(t_X, parameters, keep_prob=0.5)
print("A3 = "+ str(A3))

forward_propagation_with_dropout_test(forward_propagation_with_dropout)
```

6.2 - Backward Propagation with Dropout

Exercise 4 - backward_propagation_with_dropout

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 $A^{[1]}$. In backpropagation, you will have to **shut down the same neurons**, by reapplying the same mask $D^{[1]}$ to $dA^{[1]}$.

2. During forward propagation, you had divided $A^{[1]}$ by **keep_prob**. In backpropagation, you'll therefore have to **divide $dA^{[1]}$ 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 [19]: # 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)

    # ~ 2 lines of code)
    # dA2 =
    # Step 1: Apply mask D2 to shut down the same neurons
as during the forward propagation
    # dA2 =
    # Step 2: Scale the value of neurons that haven't been
shut down
    # YOUR CODE STARTS HERE
    dA2 = np.multiply(dA2, D2)
    dA2 = dA2 / keep_prob
    # YOUR CODE ENDS 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)
    # ~ 2 lines of code)
    # dA1 =
    # Step 1: Apply mask D1 to shut down the same neurons
as during the forward propagation
    # dA1 =
    # Step 2: Scale the value of neurons that haven't been
shut down
    # YOUR CODE STARTS HERE
    dA1 = np.multiply(dA1, D1)
    dA1 = dA1 / keep_prob
    # YOUR CODE ENDS 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 [20]: t_X, t_Y, cache = backward_propagation_with_dropout_test_case()

gradients = backward_propagation_with_dropout(t_X, t_Y, cache, keep_prob=0.5)
print("dA1 = "+ str(gradients["dA1"]))
print("dA2 = "+ str(gradients["dA2"]))

backward_propagation_with_dropout_test(backward_propagation_with_dropout)
```

It'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 [21]: parameters = model(train_X, train_Y, keep_prob = 0.86, learning_rate = 0.1)

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)

print("Accuracy on the train set: %f" % (accuracy_lr(train_X, train_Y, predictions_train)))
print("Accuracy on the test set: %f" % (accuracy_lr(test_X, test_Y, predictions_test)))
```

Dropout works great! The test accuracy has increased again (to 95%). Your model is not overfitting the training set and does a great job on the test set. The French football team will be forever grateful to you!

Run the code below to plot the decision boundary.

```
In [22]: plt.title("Model with dropout")
axes = plt.gca()
axes.set_xlim([-1.5, 1.5])
axes.set_ylim([-1.5, 1.5])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Note:

- A common mistake when using dropout is to use it both in training and testing. You should use **dropout (randomly eliminate nodes) only in training**.
- Deep learning frameworks like **TensorFlow**, **PaddlePaddle**, **Keras** or **Caffe** come with a dropout layer implementation. Don't stress - you will soon learn some of these frameworks.

What you should remember about dropout:

- Dropout is a regularization technique.
- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.
- During training time, divide each dropout layer by **keep_prob** to keep the same expected value for the activations. For example, if **keep_prob** is 0.5, then we will on average shut down half the nodes, so the output will be scaled by 0.5 since only the remaining half are contributing to the solution. Dividing by 0.5 is equivalent to multiplying by 2. Hence, the output now has the same expected value. You can check that this works even when **keep_prob** is other values than 0.5.

7 - Conclusions

Here are the results of our three models:

model	train accuracy	test accuracy
3-layer NN without regularization	95%	91.5%
3-layer NN with L2-regularization	94%	93%
3-layer NN with dropout	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.

Congratulations for finishing this assignment! And also for revolutionizing French football. :-)

What we want you to remember from this notebook:

- Regularization will help you reduce overfitting.
- Regularization will drive your weights to lower values.
- L2 regularization and Dropout are two very effective regularization techniques.

Now run the code below to plot the decision boundary.

```
In [23]: plt.title("Model with dropout")
axes = plt.gca()
axes.set_xlim([-1.5, 1.5])
axes.set_ylim([-1.5, 1.5])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Note:

- The common mistake when using dropout is to use it both in training and testing. You should use **dropout (randomly eliminate nodes) only in training**.
- Deep learning frameworks like **TensorFlow**, **PaddlePaddle**, **Keras** or **Caffe** come with a dropout layer implementation. Don't stress - you will soon learn some of these frameworks.

What you should remember about dropout:

- Dropout is a regularization technique.
- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.
- During training time, divide each dropout layer by **keep_prob** to keep the same expected value for the activations. For example, if **keep_prob** is 0.5, then we will on average shut down half the nodes, so the output will be scaled by 0.5 since only the remaining half are contributing to the solution. Dividing by 0.5 is equivalent to multiplying by 2. Hence, the output now has the same expected value. You can check that this works even when **keep_prob** is other values than 0.5.

8 - Bibliography

Deep learning is a very active research area. If you want to learn more about it, here are some resources:

• Yoshua Bengio's homepage: <http://www.cs.toronto.edu/~bengioy/>

• Andrew Ng's Coursera course on machine learning: <https://www.coursera.org/learn/machine-learning>

• Yoshua Bengio's book: <http://www.deeplearningbook.org/>

• Ian Goodfellow, Yoshua Bengio and Aaron Courville's book: <http://www.deeplearningbook.org/>

• Yoshua Bengio's book: <http://www.cs.toronto.edu/~bengioy/DeepLearningBook.html>

• Yoshua Bengio's book: [http://](http://www.cs.toronto.edu/~bengioy/DeepLearningBook.html)