

机器学习系列 (2)

提高深度神经网络性能之--合理初始化

在深度学习的训练中，权重的初始化起着重要的作用，如何选择合适的初始化方式，影响着网络的学习能力，今天我们就来看一下各种初始化方式有什么样的异同吧~

1.零初始化:

- W 和 b 都初始化为零（矩阵、向量，下同），这种方式由于导致前馈过程具有对称性，反向传播也产生对称性，导致权重值一样，相当于降低了隐藏层神经元的个数，从而无法进行有效的学习。

2.随机初始化:

- W 初始化为随机值， b 初始化为0，这种方式尽管 b 都为0，但由于 W 值不同，打破了对称性，可以有效学习，但是若网络深度很深，可能会产生梯度消失或梯度爆炸的问题，需要更好的方式。
- 梯度消失：通常神经网络所用的激活函数是sigmoid函数，这个函数将负无穷到正无穷的输入映射到0和1之间，这个函数求导的结果是 $f'(x)=f(x)(1-f(x))$ ，两个0到1之间的数相乘，得到的结果就会变得很小了。由于神经网络的反向传播是逐层对函数偏导相乘，因此当神经网络层数非常深的时候，最后一层产生的偏差就因为乘了很多的小于1的数而越来越小，非常接近0，从而导致层数比较浅的权重没有得到更新。
- 那么什么是梯度爆炸呢？梯度爆炸就是由于初始化权值过大，前面层会比后面层变化的更快，就会导致权值越来越大，这就是梯度爆炸。
- 以一个三层网络为例：如下结构：

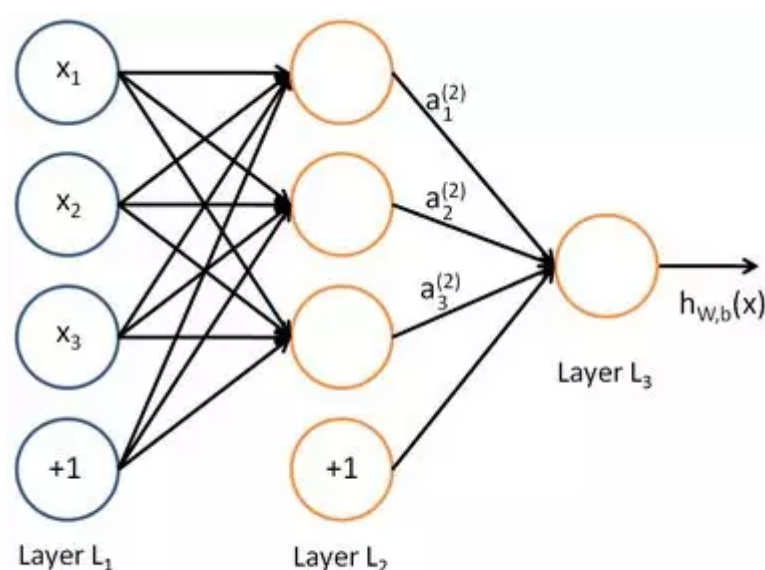


Figure 1:三层前馈举例（图片来源知乎koala tree用户）

- 那么表达式就如下图所示：

$$\begin{aligned} a_1^{(2)} &= f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + W_{13}^{(1)} x_3 + b_1^{(1)}) \\ a_2^{(2)} &= f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{23}^{(1)} x_3 + b_2^{(1)}) \\ a_3^{(2)} &= f(W_{31}^{(1)} x_1 + W_{32}^{(1)} x_2 + W_{33}^{(1)} x_3 + b_3^{(1)}) \\ h_{W,b}(x) &= a_1^{(3)} = f(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)}) \end{aligned}$$

Figure 2:三层前馈计算式（图片来源知乎koala tree用户）

- 从上式可以看出，如果每个权重都一样，那么在多层网络中，从第二层开始，每一层的输入值都是相同的了也就是 $a_1=a_2=a_3=...$ ，那么就相当于一个输入了，那么反向传播算法最后得到的 W 和 b 也都相同，即对称的，就不能拟合任意输入到输出的映射了。

3.He initialization:

- W 初始化为随机值后，再除以 $\sqrt{\frac{2}{\text{dimension of the previous layer}}}$ ， b 初始化为0，这种方式称为"he initialization"，可以很好地控制每一层权重，使得其不会太大，从而有效降低梯度爆炸发生的概率，尤其对relu激活函数有效。

4.Xavier initialization:

- 注："Xavier initialization" $\sqrt{\frac{1}{\text{dimension of the previous layer}}}$

申明

本文原理解释和公式推导均由LSayhi完成，供学习参考，可传播；代码实现的框架由Coursera提供，由LSayhi完成，详细数据和代码可在github中查询，请勿用于Coursera刷分。

<https://github.com/LSayhi/DeepLearning> (<https://github.com/LSayhi/DeepLearning>)

微信公众号：AI有点可ai

Initialization

Welcome to the first assignment of "Improving Deep Neural Networks".

Training your neural network requires specifying an initial value of the weights. A well chosen initialization method will help learning.

If you completed the previous course of this specialization, you probably followed our instructions for weight initialization, and it has worked out so far. But how do you choose the initialization for a new neural network? In this notebook, you will see how different initializations lead to different results.

A well chosen initialization can:

- Speed up the convergence of gradient descent
- Increase the odds of gradient descent converging to a lower training (and generalization) error

To get started, run the following cell to load the packages and the planar dataset you will try to classify.

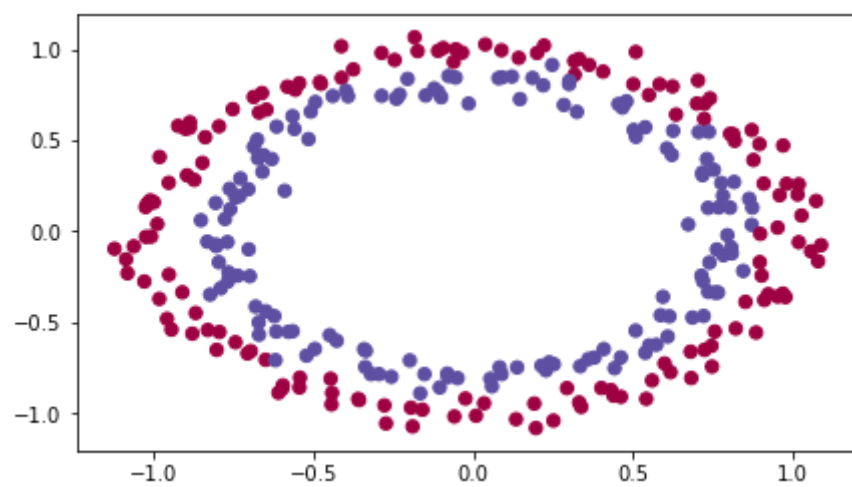
一个优秀的参数初始化方案可以实现:

- 增加训练和泛化误差收敛的概率
- 加速梯度下降的收敛

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import sklearn
import sklearn.datasets
from init_utils import sigmoid, relu, compute_loss, forward_propagation, backward_propagation
from init_utils import update_parameters, predict, load_dataset, plot_decision_boundary, predict_dec

%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'

# load image dataset: blue/red dots in circles
train_X, train_Y, test_X, test_Y = load_dataset()
```



You would like a classifier to separate the blue dots from the red dots.

1 - Neural Network model

You will use a 3-layer neural network (already implemented for you). Here are the initialization methods you will experiment with:

- *Zeros initialization* -- setting `initialization = "zeros"` in the input argument.
- *Random initialization* -- setting `initialization = "random"` in the input argument. This initializes the weights to large random values.
- *He initialization* -- setting `initialization = "he"` in the input argument. This initializes the weights to random values scaled according to a paper by He et al., 2015.

Instructions: Please quickly read over the code below, and run it. In the next part you will implement the three initialization methods that this `model()` calls.

2015年, “he”初始化方式被提出

```

In [2]: def model(X, Y, learning_rate = 0.01, num_iterations = 15000, print_cost = True, initialization = "he"):
        """
        Implements a three-layer neural network: LINEAR->RELU->LINEAR->RELU->LINEAR->SIGMOID.

        Arguments:
        X -- input data, of shape (2, number of examples)
        Y -- true "label" vector (containing 0 for red dots; 1 for blue dots), of shape (1, number of examples)
        learning_rate -- learning rate for gradient descent
        num_iterations -- number of iterations to run gradient descent
        print_cost -- if True, print the cost every 1000 iterations
        initialization -- flag to choose which initialization to use ("zeros", "random" or "he")

        Returns:
        parameters -- parameters learnt by the model
        """

        grads = {}
        costs = [] # to keep track of the loss
        m = X.shape[1] # number of examples
        layers_dims = [X.shape[0], 10, 5, 1]

        # Initialize parameters dictionary.
        if initialization == "zeros":
            parameters = initialize_parameters_zeros(layers_dims)
        elif initialization == "random":
            parameters = initialize_parameters_random(layers_dims)
        elif initialization == "he":
            parameters = initialize_parameters_he(layers_dims)

        # Loop (gradient descent)

        for i in range(0, num_iterations):

            # Forward propagation: LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID.
            a3, cache = forward_propagation(X, parameters)

            # Loss
            cost = compute_loss(a3, Y)

            # Backward propagation.
            grads = backward_propagation(X, Y, cache)

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

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

            # plot the loss
            plt.plot(costs)
            plt.ylabel('cost')
            plt.xlabel('iterations (per hundreds)')
            plt.title("Learning rate =" + str(learning_rate))
            plt.show()

        return parameters

```

2 - Zero initialization

There are two types of parameters to initialize in a neural network:

- the weight matrices ($W^{[1]}, W^{[2]}, W^{[3]}, \dots, W^{[L-1]}, W^{[L]}$)
- the bias vectors ($b^{[1]}, b^{[2]}, b^{[3]}, \dots, b^{[L-1]}, b^{[L]}$)

Exercise: Implement the following function to initialize all parameters to zeros. You'll see later that this does not work well since it fails to "break symmetry", but lets try it anyway and see what happens. Use `np.zeros((...,))` with the correct shapes.

如果对W和b参数初始化为0（零向量、零矩阵），看看效果

```
In [6]: # GRADED FUNCTION: initialize_parameters_zeros

def initialize_parameters_zeros(layers_dims):
    """
    Arguments:
    layer_dims -- python array (list) containing the size of each layer.

    Returns:
    parameters -- python dictionary containing your parameters "W1", "b1", ..., "WL", "bL":
        W1 -- weight matrix of shape (layers_dims[1], layers_dims[0])
        b1 -- bias vector of shape (layers_dims[1], 1)
        ...
        WL -- weight matrix of shape (layers_dims[L], layers_dims[L-1])
        bL -- bias vector of shape (layers_dims[L], 1)
    """

    parameters = {}
    L = len(layers_dims)      # number of layers in the network

    for l in range(1, L):
        ### START CODE HERE ### (≈ 2 lines of code)
        parameters['W' + str(l)] = np.zeros((layers_dims[l], layers_dims[l-1]))
        parameters['b' + str(l)] = np.zeros((layers_dims[l], 1))
        ### END CODE HERE ###
    return parameters
```

```
In [7]: parameters = initialize_parameters_zeros([3, 2, 1])
print("W1 = " + str(parameters["W1"]))
print("b1 = " + str(parameters["b1"]))
print("W2 = " + str(parameters["W2"]))
print("b2 = " + str(parameters["b2"]))
```

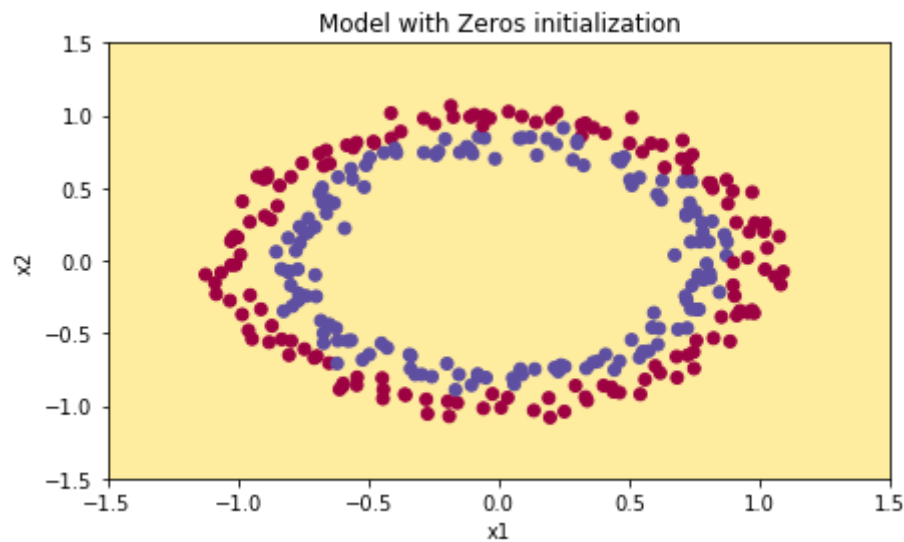
```
W1 = [[ 0.  0.  0.]
 [ 0.  0.  0.]]
b1 = [[ 0.]
 [ 0.]]
W2 = [[ 0.  0.]]
b2 = [[ 0.]]
```

Expected Output:

```
W1 [[ 0. 0. 0.] [ 0. 0. 0.]]
b1      [[ 0.] [ 0.]]
W2      [[ 0. 0.]]
b2      [[ 0.]]
```

Run the following code to train your model on 15,000 iterations using zeros initialization.


```
In [10]: plt.title("Model with Zeros initialization")
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)
```



The model is predicting 0 for every example.

In general, initializing all the weights to zero results in the network failing to break symmetry. This means that every neuron in each layer will learn the same thing, and you might as well be training a neural network with $n^{[l]} = 1$ for every layer, and the network is no more powerful than a linear classifier such as logistic regression.

全部初始化为0（零向量或零矩阵），并没有打破网络的对称性，每一层的每个神经元学习到的知识是一样的，那么神经网络就像是每一层只有一个神经元，无法学习到足够的特征

What you should remember:

- The weights $W^{[l]}$ should be initialized randomly to break symmetry.
- It is however okay to initialize the biases $b^{[l]}$ to zeros. Symmetry is still broken so long as $W^{[l]}$ is initialized randomly.

3 - Random initialization

To break symmetry, let's initialize the weights randomly. Following random initialization, each neuron can then proceed to learn a different function of its inputs. In this exercise, you will see what happens if the weights are initialized randomly, but to very large values.

Exercise: Implement the following function to initialize your weights to large random values (scaled by *10) and your biases to zeros. Use `np.random.randn(...)` * 10 for weights and `np.zeros(...)` for biases. We are using a fixed `np.random.seed(...)` to make sure your "random" weights match ours, so don't worry if running several times your code gives you always the same initial values for the parameters.

所以W不能全部初始化为0，我们将W随机初始化（因为W非对称了，b即使全部初始化为0，网络还是非对称的，因此b可以初始化为0）

```
In [13]: # GRADED FUNCTION: initialize_parameters_random

def initialize_parameters_random(layers_dims):
    """
    Arguments:
    layer_dims -- python array (list) containing the size of each layer.

    Returns:
    parameters -- python dictionary containing your parameters "W1", "b1", ..., "WL", "bL":
        W1 -- weight matrix of shape (layers_dims[1], layers_dims[0])
        b1 -- bias vector of shape (layers_dims[1], 1)
        ...
        WL -- weight matrix of shape (layers_dims[L], layers_dims[L-1])
        bL -- bias vector of shape (layers_dims[L], 1)
    """

    np.random.seed(3) # This seed makes sure your "random" numbers will be the as ours
    parameters = {}
    L = len(layers_dims) # integer representing the number of layers

    for l in range(1, L):
        ### START CODE HERE ### (≈ 2 lines of code)
        parameters['W' + str(l)] = np.random.randn(layers_dims[l], layers_dims[l-1])*10
        parameters['b' + str(l)] = np.zeros((layers_dims[l], 1))
        ### END CODE HERE ###

    return parameters
```

```
In [14]: parameters = initialize_parameters_random([3, 2, 1])
print("W1 = " + str(parameters["W1"]))
print("b1 = " + str(parameters["b1"]))
print("W2 = " + str(parameters["W2"]))
print("b2 = " + str(parameters["b2"]))
```

```
W1 = [[ 17.88628473  4.36509851  0.96497468]
      [-18.63492703 -2.77388203 -3.54758979]]
b1 = [[ 0.]
      [ 0.]]
W2 = [[-0.82741481 -6.27000677]]
b2 = [[ 0.]]
```

Expected Output:

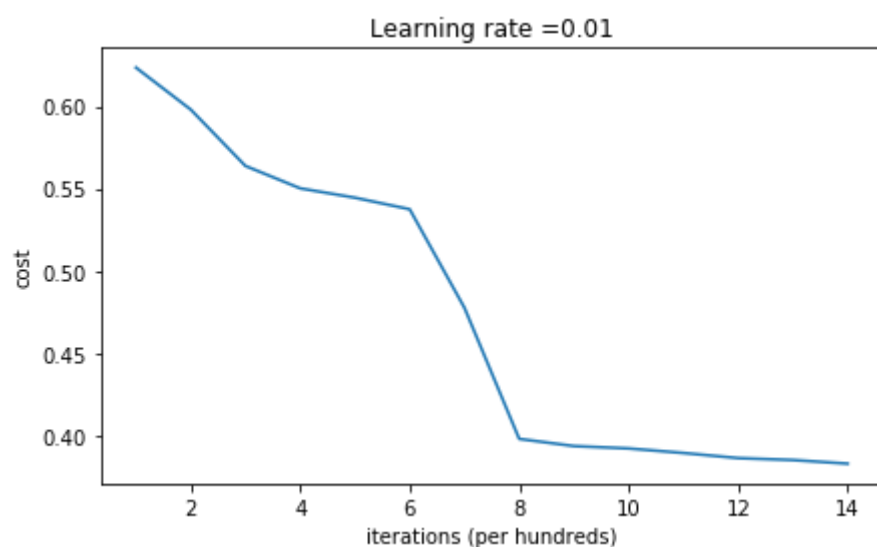
```
W1 [[ 17.88628473 4.36509851 0.96497468] [-18.63492703 -2.77388203 -3.54758979]]
b1 [[ 0.] [ 0.]]
W2 [[-0.82741481 -6.27000677]]
b2 [[ 0.]]
```

Run the following code to train your model on 15,000 iterations using random initialization.

```
In [15]: parameters = model(train_X, train_Y, initialization = "random")
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)
```

```
C:\Users\BD\代码作业\第二课第一周编程作业\assignment1\init_utils.py:145: RuntimeWarning: divide by zero encountered in log
  logprobs = np.multiply(-np.log(a3), Y) + np.multiply(-np.log(1 - a3), 1 - Y)
C:\Users\BD\代码作业\第二课第一周编程作业\assignment1\init_utils.py:145: RuntimeWarning: invalid value encountered in multiply
  logprobs = np.multiply(-np.log(a3), Y) + np.multiply(-np.log(1 - a3), 1 - Y)
```

```
Cost after iteration 0: inf
Cost after iteration 1000: 0.6235719528716395
Cost after iteration 2000: 0.5980821226022246
Cost after iteration 3000: 0.5637996692567824
Cost after iteration 4000: 0.5501754102867465
Cost after iteration 5000: 0.5444767640123352
Cost after iteration 6000: 0.5374657035647926
Cost after iteration 7000: 0.4775406670630984
Cost after iteration 8000: 0.39784053325714386
Cost after iteration 9000: 0.3934817369887478
Cost after iteration 10000: 0.39203280921110983
Cost after iteration 11000: 0.38927347547167324
Cost after iteration 12000: 0.3861625886188003
Cost after iteration 13000: 0.38499044850062425
Cost after iteration 14000: 0.38279756848782404
```



```
On the train set:
Accuracy: 0.83
On the test set:
Accuracy: 0.86
```

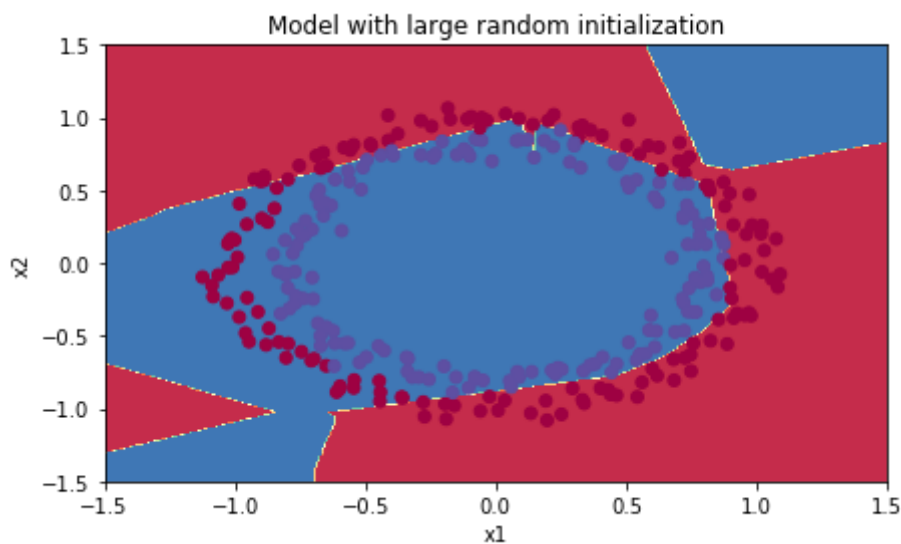
If you see "inf" as the cost after the iteration 0, this is because of numerical roundoff; a more numerically sophisticated implementation would fix this. But this isn't worth worrying about for our purposes.

Anyway, it looks like you have broken symmetry, and this gives better results. than before. The model is no longer outputting all 0s.

```
In [16]: print (predictions_train)
print (predictions_test)
```

```
[[1 0 1 1 0 0 1 1 1 1 1 0 1 0 0 1 0 1 1 0 0 0 1 0 1 1 1 1 1 1 1 0 1 1 0 0 1 1
 1 1 1 1 1 1 0 1 1 1 1 0 1 0 1 1 1 1 0 0 1 1 1 1 0 1 1 0 1 0 1 1 1 1 0 0 0
 0 0 1 0 1 0 1 1 1 0 0 1 1 1 1 1 1 0 0 1 1 1 0 1 1 0 1 0 1 1 0 1 1 0 1 0 1
 1 0 0 1 0 0 1 1 0 1 1 1 0 1 0 0 1 0 1 1 1 1 1 1 1 1 0 1 1 0 0 1 1 0 0 0 1 0
 1 0 1 0 1 1 1 0 0 1 1 1 1 0 1 1 0 1 0 1 1 0 1 0 1 1 1 1 1 0 1 1 1 1 0 1 0 1
 0 1 1 1 1 0 1 1 0 1 1 0 1 1 0 1 0 1 1 1 1 0 1 1 1 0 1 0 1 0 0 1 0 1 1 0 1 1
 0 1 1 0 1 1 1 0 1 1 1 1 0 1 0 0 1 1 0 1 1 1 0 0 0 1 1 0 1 1 1 1 0 1 1 0 1
 1 1 0 0 1 0 0 0 1 0 0 0 1 1 1 1 0 0 0 0 1 1 1 1 0 0 1 1 1 1 1 1 1 1 0 0 0 1
 1 1 1 0]]
[[1 1 1 1 0 1 0 1 1 0 1 1 1 0 0 0 0 1 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0
 1 1 0 0 1 1 1 1 1 0 1 1 1 0 1 0 1 1 0 1 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1
 1 1 1 0 1 0 0 1 0 0 0 1 1 0 1 1 0 0 0 1 1 0 1 1 0 0]]
```

```
In [17]: plt.title("Model with large random initialization")
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)
```



Observations:

- The cost starts very high. This is because with large random-valued weights, the last activation (sigmoid) outputs results that are very close to 0 or 1 for some examples, and when it gets that example wrong it incurs a very high loss for that example. Indeed, when $\log(a^{[3]}) = \log(0)$, the loss goes to infinity.
- Poor initialization can lead to vanishing/exploding gradients, which also slows down the optimization algorithm.
- If you train this network longer you will see better results, but initializing with overly large random numbers slows down the optimization.

In summary:

- Initializing weights to very large random values does not work well.
- Hopefully initializing with small random values does better. The important question is: how small should be these random values be? Lets find out in the next part!
- 从上例中看出，当权重初始化不合适（过大/过小）时，可能会产生梯度爆炸/梯度消失问题，那么如何选择合适的初始化值呢？看下文

4 - He initialization

Finally, try "He Initialization"; this is named for the first author of He et al., 2015. (If you have heard of "Xavier initialization", this is similar except Xavier initialization uses a scaling factor for the weights $W^{[l]}$ of $\sqrt{1./\text{layers_dims}[l-1]}$ where He initialization would use $\sqrt{2./\text{layers_dims}[l-1]}$.)

Exercise: Implement the following function to initialize your parameters with He initialization.

Hint: This function is similar to the previous `initialize_parameters_random(...)`. The only difference is that instead of multiplying `np.random.randn(...)` by 10, you will multiply it by $\sqrt{\frac{2}{\text{dimension of the previous layer}}}$, which is what He initialization recommends for layers with a ReLU activation.

在初始化 W 时，使用“he”方法，即将 W 除以 $\sqrt{\frac{2}{\text{dimension of the previous layer}}}$


```
In [24]: # GRADED FUNCTION: initialize_parameters_he

def initialize_parameters_he(layers_dims):
    """
    Arguments:
    layer_dims -- python array (list) containing the size of each layer.

    Returns:
    parameters -- python dictionary containing your parameters "W1", "b1", ..., "WL", "bL":
        W1 -- weight matrix of shape (layers_dims[1], layers_dims[0])
        b1 -- bias vector of shape (layers_dims[1], 1)
        ...
        WL -- weight matrix of shape (layers_dims[L], layers_dims[L-1])
        bL -- bias vector of shape (layers_dims[L], 1)
    """

    np.random.seed(3)
    parameters = {}
    L = len(layers_dims) - 1 # integer representing the number of layers

    for l in range(1, L + 1):
        ### START CODE HERE ### (≈ 2 lines of code)
        parameters['W' + str(l)] = np.random.randn(layers_dims[l], layers_dims[l-1])*np.sqrt(2/layers_dims[l-1])
        parameters['b' + str(l)] = np.zeros((layers_dims[l], 1))
        ### END CODE HERE ###

    return parameters
```

```
In [25]: parameters = initialize_parameters_he([2, 4, 1])
print("W1 = " + str(parameters["W1"]))
print("b1 = " + str(parameters["b1"]))
print("W2 = " + str(parameters["W2"]))
print("b2 = " + str(parameters["b2"]))
```

```
W1 = [[ 1.78862847  0.43650985]
 [ 0.09649747 -1.8634927 ]
 [-0.2773882  -0.35475898]
 [-0.08274148 -0.62700068]]
b1 = [[ 0.]
 [ 0.]
 [ 0.]
 [ 0.]]
W2 = [[-0.03098412 -0.33744411 -0.92904268  0.62552248]]
b2 = [[ 0.]
```

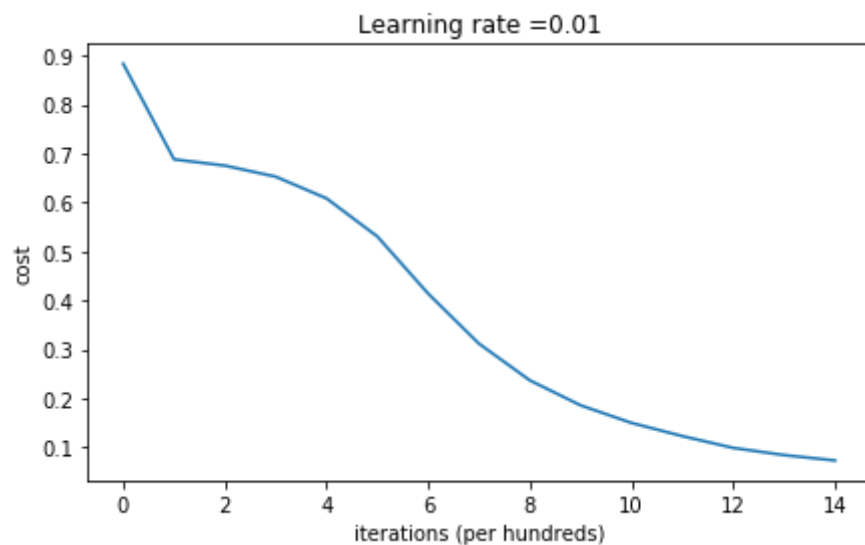
Expected Output:

```
W1 [[ 1.78862847  0.43650985] [ 0.09649747 -1.8634927 ] [-0.2773882 -0.35475898] [-0.08274148 -0.62700068]]
b1 [[ 0.] [ 0.] [ 0.] [ 0.]]
W2 [[-0.03098412 -0.33744411 -0.92904268  0.62552248]]
b2 [[ 0.]
```

Run the following code to train your model on 15,000 iterations using He initialization.

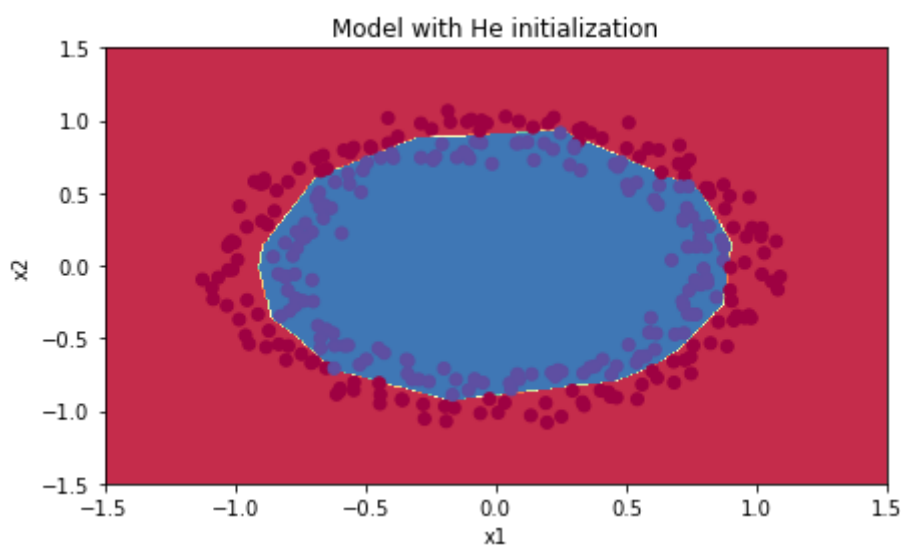
```
In [26]: parameters = model(train_X, train_Y, initialization = "he")
print ("On the train set:")
predictions_train = predict(train_X, train_Y, parameters)
print ("On the test set:")
predictions_test = predict(test_X, test_Y, parameters)
```

```
Cost after iteration 0: 0.8830537463419761
Cost after iteration 1000: 0.6879825919728063
Cost after iteration 2000: 0.6751286264523371
Cost after iteration 3000: 0.6526117768893807
Cost after iteration 4000: 0.6082958970572938
Cost after iteration 5000: 0.5304944491717495
Cost after iteration 6000: 0.4138645817071795
Cost after iteration 7000: 0.31178034648444414
Cost after iteration 8000: 0.23696215330322565
Cost after iteration 9000: 0.18597287209206842
Cost after iteration 10000: 0.15015556280371808
Cost after iteration 11000: 0.12325079292273552
Cost after iteration 12000: 0.09917746546525931
Cost after iteration 13000: 0.08457055954024274
Cost after iteration 14000: 0.07357895962677365
```



```
On the train set:
Accuracy: 0.9933333333333333
On the test set:
Accuracy: 0.96
```

```
In [27]: plt.title("Model with He initialization")
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)
```



Observations:

- The model with He initialization separates the blue and the red dots very well in a small number of iterations.
- “he”初始化方式提高了训练集和测试集的准确度，降低了迭代次数

5 - Conclusions

You have seen three different types of initializations. For the same number of iterations and same hyperparameters the comparison is:

Model	Train accuracy	Problem/Comment
3-layer NN with zeros initialization	50%	fails to break symmetry
3-layer NN with large random initialization	83%	too large weights
3-layer NN with He initialization	99%	recommended method

What you should remember from this notebook:

- Different initializations lead to different results
- Random initialization is used to break symmetry and make sure different hidden units can learn different things
- Don't initialize to values that are too large
- He initialization works well for networks with ReLU activations.