

Week 1 - prog ex: Initialization, Regularization & GC

笔记本: DL 2 - Deep NN Hyperparameter Tuning, Regularization & Optimization

创建时间: 2021/1/9 10:07

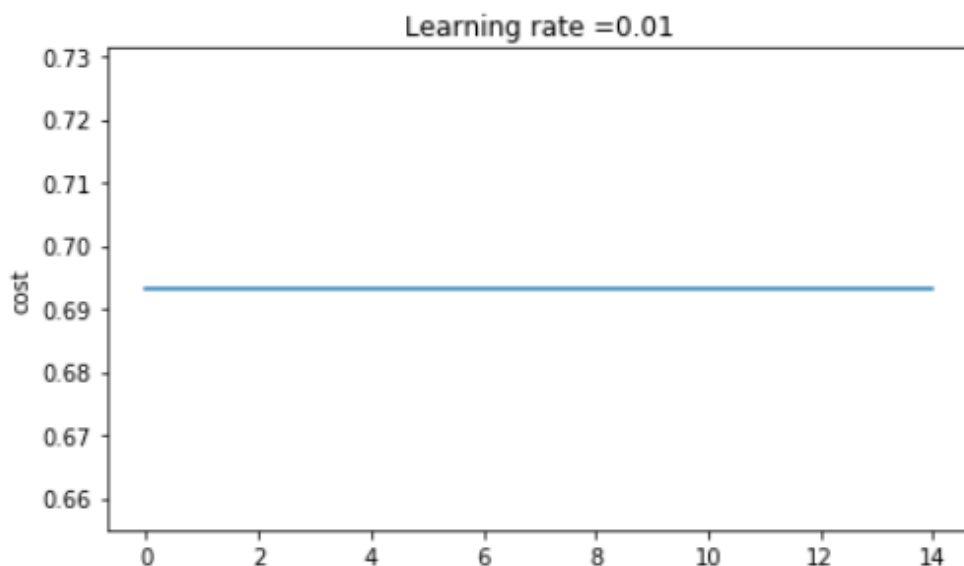
更新时间: 2021/1/9 10:58

Welcome to the first assignment of the hyper parameters tuning specialization. It is very important that you regularize your model properly because it could dramatically improve your results.

By completing this assignment you will:

- Understand that different initialization methods and their impact on your model performance
- Implement *zero initialization* and see it fails to "break symmetry",
- Recognize that *random initialization* "breaks symmetry" and yields more efficient models,
- Understand that you could use both *random initialization* and *scaling* to get even better training performance on your model.

```
Cost after iteration 0: 0.6931471805599453
Cost after iteration 1000: 0.6931471805599453
Cost after iteration 2000: 0.6931471805599453
Cost after iteration 3000: 0.6931471805599453
Cost after iteration 4000: 0.6931471805599453
Cost after iteration 5000: 0.6931471805599453
Cost after iteration 6000: 0.6931471805599453
Cost after iteration 7000: 0.6931471805599453
Cost after iteration 8000: 0.6931471805599453
Cost after iteration 9000: 0.6931471805599453
Cost after iteration 10000: 0.6931471805599455
Cost after iteration 11000: 0.6931471805599453
Cost after iteration 12000: 0.6931471805599453
Cost after iteration 13000: 0.6931471805599453
Cost after iteration 14000: 0.6931471805599453
```



The model is predicting 0 for every example.

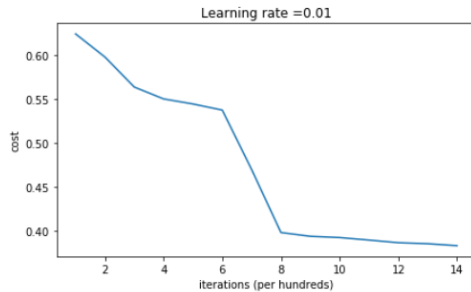
In general, initializing all the weights to zero results in the network failing to break symmetry. This means that every neuron in each layer will learn the same thing, and you might as well be training a neural network with $n^{[l]} = 1$ for every layer, and the network is no more powerful than a linear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

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.

```
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 ###
```



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.

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 $\log(a^{31}) = \log(0)$, 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!

He Initialization

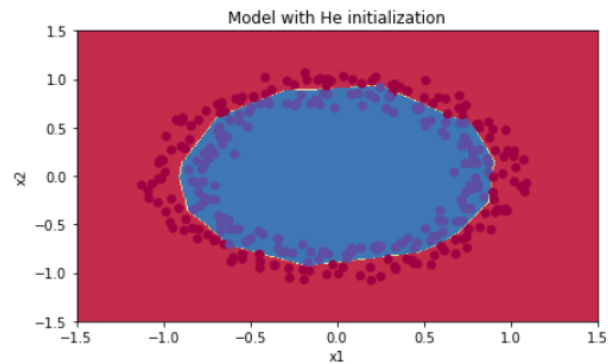
4 - He initialization

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 $\sqrt{1./layers_dims[l-1]}$ where He initialization would use $\sqrt{2./layers_dims[l-1]}$.)

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.

```
In [26]: 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.

5 - Conclusions

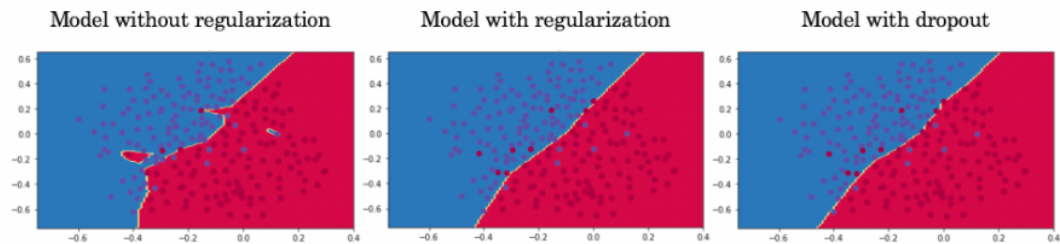
You have seen three different types of initializations. For the same number of iterations and same hyperparameters the comparison is:

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.

Welcome to the first assignment of the hyper parameters tuning specialization. It is very important that you regularize your model properly because it could dramatically improve your results.



By completing this assignment you will:

- Understand that different regularization methods that could help your model.
- Implement dropout and see it work on data.
- Recognize that a model without regularization gives you a better accuracy on the training set but not necessarily on the test set.
- Understand that you could use **both dropout and regularization** on your model.

Problem Statement: You have just been hired as an AI 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.

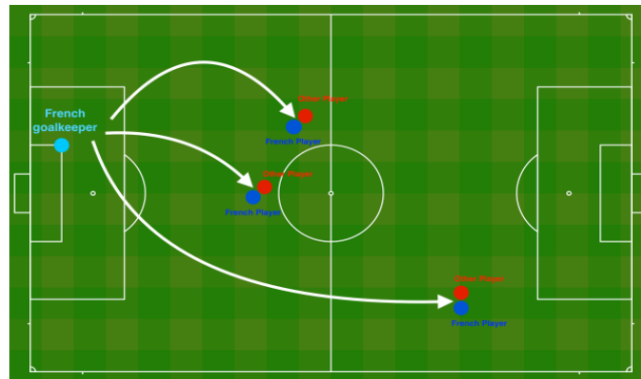
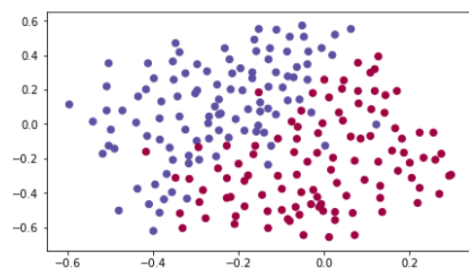


Figure 1: Football field

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.



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 their head

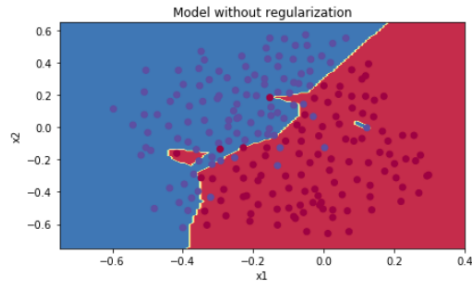
On the training set:

Accuracy: 0.947867298578

On the test set:

Accuracy: 0.915

```
In [5]: 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.

L2

2 - L2 Regularization

The standard way to avoid overfitting is called **L2 regularization**. It consists of appropriately modifying your cost function, from:

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

To:

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

Let's modify your cost and observe the consequences.

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

```
np.sum(np.square(W1))
```

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

Of course, because you changed the cost, you have to change backward propagation as well! All the gradients have to be computed with respect to this new cost.

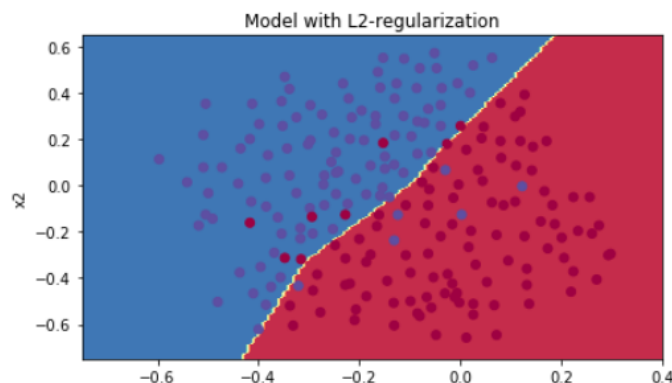
Exercise: Implement the changes needed in backward propagation to take into account regularization. The changes only concern $dW1$, $dW2$ and $dW3$. For each, you have to add the regularization term's gradient ($\frac{d}{dW}(\frac{1}{2} \frac{\lambda}{m} W^2) = \frac{\lambda}{m} W$).

```
On the train set:
Accuracy: 0.938388625592
On the test set:
Accuracy: 0.93
```

Congrats, the test set accuracy increased to 93%. You have saved the French football team!

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

```
In [13]: 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 λ is a hyperparameter that you can tune using a dev set.
- L2 regularization makes your decision boundary smoother. If λ 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.

What you should remember – the implications of L2-regularization on:

- The cost computation:
 - A regularization term is added to the cost
- The backpropagation function:
 - There are extra terms in the gradients with respect to weight matrices
- Weights end up smaller ("weight decay"):
 - Weights are pushed to smaller values.

Dropout

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.

Instructions: You would like to shut down some neurons in the first and second layers. To do that, you are going to carry out 4 Steps:

1. In lecture, we discussed creating a variable $d^{[1]}$ with the same shape as $A^{[1]}$ using `np.random.rand()` to randomly get numbers between 0 and 1. Here, you will use a vectorized implementation, so create a random matrix $D^{[1]} = [d^{1} d^{[1](2)} \dots d^{[1](m)}]$ of the same dimension as $A^{[1]}$.
2. Set each entry of $D^{[1]}$ to be 1 with probability (`keep_prob`), and 0 otherwise.

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)
```

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`).

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.

4 - 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.

Welcome to this week's third programming assignment! You will be implementing **gradient checking** to make sure that your backpropagation implementation is correct. By completing this assignment you will:

- Implement gradient checking from scratch.
- Understand how to use the *difference* formula to check your backpropagation implementation.
- Recognize that your backpropagation algorithm should give you similar results as the ones you got by computing the difference formula.
- Learn how to identify which parameter's gradient was computed incorrectly.

1) How does gradient checking work? ¶

Backpropagation computes the gradients $\frac{\partial J}{\partial \theta}$, where θ denotes the parameters of the model. J is computed using forward propagation and your loss function.

Because forward propagation is relatively easy to implement, you're confident you got that right, and so you're almost 100% sure that you're computing the cost J correctly. Thus, you can use your code for computing J to verify the code for computing $\frac{\partial J}{\partial \theta}$.

Let's look back at the definition of a derivative (or gradient):

$$\frac{\partial J}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad (1)$$

If you're not familiar with the " $\lim_{\epsilon \rightarrow 0}$ " notation, it's just a way of saying "when ϵ is really really small."

We know the following:

- $\frac{\partial J}{\partial \theta}$ is what you want to make sure you're computing correctly.
- You can compute $J(\theta + \epsilon)$ and $J(\theta - \epsilon)$ (in the case that θ is a real number), since you're confident your implementation for J is correct.

Lets use equation (1) and a small value for ϵ to convince your CEO that your code for computing $\frac{\partial J}{\partial \theta}$ is correct!

2) 1-dimensional gradient checking

Consider a 1D linear function $J(\theta) = \theta x$. The model contains only a single real-valued parameter θ , and takes x as input.

You will implement code to compute $J(\cdot)$ and its derivative $\frac{\partial J}{\partial \theta}$. You will then use gradient checking to make sure your derivative computation for J is correct.

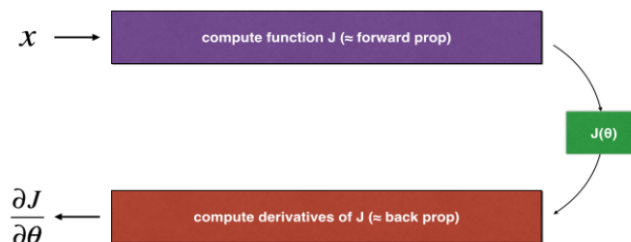


Figure 1: 1D linear model

The diagram above shows the key computation steps: First start with x , then evaluate the function $J(x)$ ("forward propagation"). Then compute the derivative $\frac{\partial J}{\partial \theta}$ ("backward propagation").

Exercise: implement "forward propagation" and "backward propagation" for this simple function. I.e., compute both $J(\cdot)$ ("forward propagation") and its derivative with respect to θ ("backward propagation"), in two separate functions.

Instructions:

- First compute "gradapprox" using the formula above (1) and a small value of ϵ . Here are the Steps to follow:

- $\theta^+ = \theta + \epsilon$
- $\theta^- = \theta - \epsilon$
- $J^+ = J(\theta^+)$
- $J^- = J(\theta^-)$
- $gradapprox = \frac{J^+ - J^-}{2\epsilon}$

- Then compute the gradient using backward propagation, and store the result in a variable "grad"
- Finally, compute the relative difference between "gradapprox" and the "grad" using the following formula:

$$difference = \frac{\|grad - gradapprox\|_2}{\|grad\|_2 + \|gradapprox\|_2} \quad (2)$$

You will need 3 Steps to compute this formula:

- 1'. compute the numerator using `np.linalg.norm(...)`
- 2'. compute the denominator. You will need to call `np.linalg.norm(...)` twice.
- 3'. divide them.
- If this difference is small (say less than 10^{-7}), you can be quite confident that you have computed your gradient correctly. Otherwise, there may be a mistake in the gradient computation.

```
# Compute gradapprox using left side of formula (1). epsilon is small enough, you don't need to worry about the limit.
### START CODE HERE ### (approx. 5 lines)
thetaplus = theta + epsilon          # Step 1
thetaminus = theta - epsilon         # Step 2
J_plus = thetaplus * x               # Step 3
J_minus = thetaminus * x            # Step 4
gradapprox = (J_plus - J_minus) / (2 * epsilon) # Step 5
### END CODE HERE ###

# Check if gradapprox is close enough to the output of backward_propagation()
### START CODE HERE ### (approx. 1 line)
grad = backward_propagation(x, theta)
### END CODE HERE ###

### START CODE HERE ### (approx. 1 line)
numerator = np.linalg.norm(grad - gradapprox, keepdims = True) # Step 1'
denominator = np.linalg.norm(grad, keepdims = True) + np.linalg.norm(gradapprox, keepdims = True) # Step 2'
difference = numerator / denominator # Step 3'
### END CODE HERE ###

if difference < 1e-7:
    print("The gradient is correct!")
else:
    print("The gradient is wrong!")
```

3) N-dimensional gradient checking

The following figure describes the forward and backward propagation of your fraud detection model.

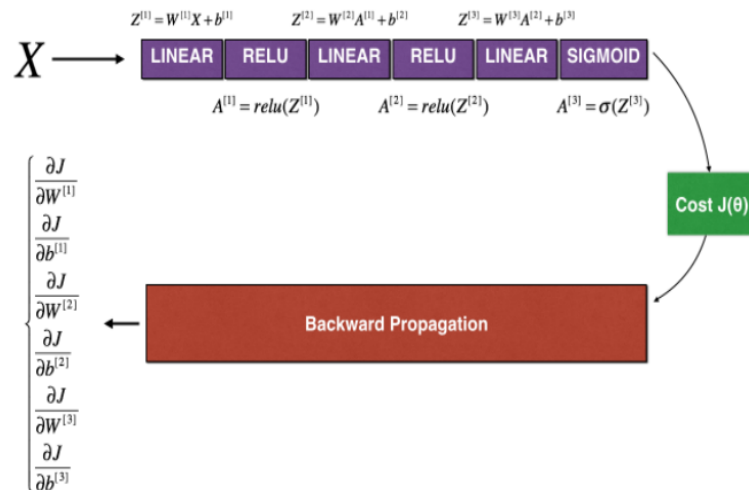


Figure 2 : deep neural network
 LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID

How does gradient checking work?

As in 1) and 2), you want to compare "gradapprox" to the gradient computed by backpropagation. The formula is still:

$$\frac{\partial J}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad (1)$$

However, θ is not a scalar anymore. It is a dictionary called "parameters". We implemented a function "dictionary_to_vector()" for you. It converts the "parameters" dictionary into a vector called "values", obtained by reshaping all parameters (W_1 , b_1 , W_2 , b_2 , W_3 , b_3) into vectors and concatenating them.

The inverse function is "vector_to_dictionary" which outputs back the "parameters" dictionary.

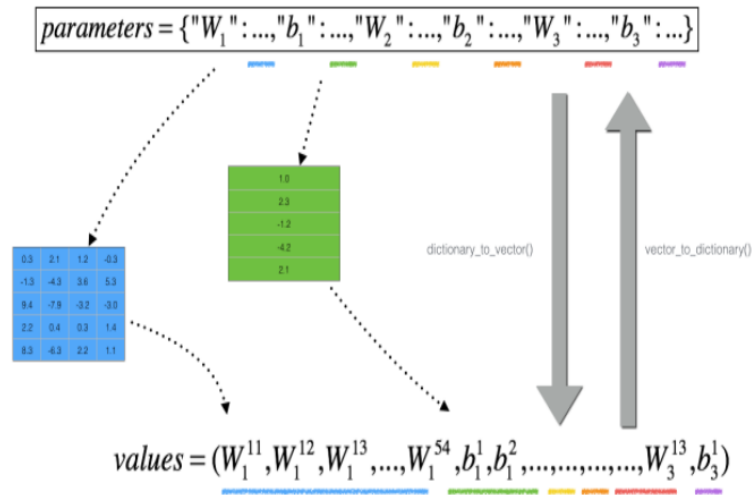


Figure 2: dictionary_to_vector() and vector_to_dictionary()

You will need these functions in gradient_check_n()

We have also converted the "gradients" dictionary into a vector "grad" using gradients_to_vector(). You don't need to worry about that.

Exercise: Implement gradient_check_n().

Exercise: Implement gradient_check_n().

Instructions: Here is pseudo-code that will help you implement the gradient check.

For each i in num_parameters:

- To compute $J_{plus}[i]$:
 - Set θ^+ to $\text{np.copy}(\text{parameters_values})$
 - Set θ^+ to $\theta^+ + \epsilon$
 - Calculate J_i^+ using $\text{forward_propagation}_n(x, y, \text{vector_to_dictionary}(\theta^+))$.
- To compute $J_{minus}[i]$: do the same thing with θ^-
- Compute $\text{gradapprox}[i] = \frac{J_i^+ - J_i^-}{2\epsilon}$

Thus, you get a vector gradapprox, where gradapprox[i] is an approximation of the gradient with respect to parameter_values[i]. You can now compare this gradapprox vector to the gradients vector from backpropagation. Just like for the 1D case (Steps 1', 2', 3'), compute:

$$\text{difference} = \frac{\|\text{grad} - \text{gradapprox}\|_2}{\|\text{grad}\|_2 + \|\text{gradapprox}\|_2} \quad (3)$$

Expected output:

There is a mistake in the backward propagation! difference = 0.285093156781

It seems that there were errors in the backward_propagation_n code we gave you! Good that you've implemented the gradient check. Go back to backward_propagation and try to find/correct the errors (Hint: check dW2 and db1). Rerun the gradient check when you think you've fixed it. Remember you'll need to re-execute the cell defining backward_propagation_n() if you modify the code.

Can you get gradient check to declare your derivative computation correct? Even though this part of the assignment isn't graded, we strongly urge you to try to find the bug and re-run gradient check until you're convinced backprop is now correctly implemented.

Note

- Gradient Checking is slow! Approximating the gradient with $\frac{\partial J}{\partial \theta} \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$ is computationally costly. For this reason, we don't run gradient checking at every iteration during training. Just a few times to check if the gradient is correct.
- Gradient Checking, at least as we've presented it, doesn't work with dropout. You would usually run the gradient check algorithm without dropout to make sure your backprop is correct, then add dropout.

Congrats, you can be confident that your deep learning model for fraud detection is working correctly! You can even use this to convince your CEO. :)

What you should remember from this notebook:

- Gradient checking verifies closeness between the gradients from backpropagation and the numerical approximation of the gradient (computed using forward propagation).
- Gradient checking is slow, so we don't run it in every iteration of training. You would usually run it only to make sure your code is correct, then turn it off and use backprop for the actual learning process.

