

# Gradient checking and advanced optimization

## From Ufldl

Backpropagation is a notoriously difficult algorithm to debug and get right, especially since many subtly buggy implementations of it—for example, one that has an off-by-one error in the indices and that thus only trains some of the layers of weights, or an implementation that omits the bias term—will manage to learn something that can look surprisingly reasonable (while performing less well than a correct implementation). Thus, even with a buggy implementation, it may not at all be apparent that anything is amiss. In this section, we describe a method for numerically checking the derivatives computed by your code to make sure that your implementation is correct. Carrying out the derivative checking procedure described here will significantly increase your confidence in the correctness of your code.

Suppose we want to minimize  $J(\theta)$  as a function of  $\theta$ . For this example, suppose  $J : \mathbb{R} \mapsto \mathbb{R}$ , so that  $\theta \in \mathbb{R}$ . In this 1-dimensional case, one iteration of gradient descent is given by

$$\theta := \theta - \alpha \frac{d}{d\theta} J(\theta).$$

Suppose also that we have implemented some function  $g(\theta)$  that purportedly computes  $\frac{d}{d\theta} J(\theta)$ , so that we implement gradient descent using the update  $\theta := \theta - \alpha g(\theta)$ . How can we check if our implementation of  $g$  is correct?

Recall the mathematical definition of the derivative as

$$\frac{d}{d\theta} J(\theta) = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}.$$

Thus, at any specific value of  $\theta$ , we can numerically approximate the derivative as follows:

$$\frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}$$

In practice, we set EPSILON to a small constant, say around  $10^{-4}$ . (There's a large range of values of EPSILON that should work well, but we don't set EPSILON to be "extremely" small, say  $10^{-20}$ , as that would lead to numerical roundoff errors.)

Thus, given a function  $g(\theta)$  that is supposedly computing  $\frac{d}{d\theta} J(\theta)$ , we can now numerically verify its correctness by checking that

$$g(\theta) \approx \frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}.$$

The degree to which these two values should approximate each other will depend on the details of  $J$ . But assuming  $\text{EPSILON} = 10^{-4}$ , you'll usually find that the left- and right-hand sides of the above will agree to at least 4 significant digits (and often many more).

Now, consider the case where  $\theta \in \mathbb{R}^n$  is a vector rather than a single real number (so that we have  $n$  parameters that we want to learn), and  $J : \mathbb{R}^n \mapsto \mathbb{R}$ . In our neural network example we used " $J(W, b)$ ," but one can imagine "unrolling" the parameters  $W, b$  into a long vector  $\theta$ . We now generalize our derivative checking procedure to the case where  $\theta$  may be a vector.

Suppose we have a function  $g_i(\theta)$  that purportedly computes  $\frac{\partial}{\partial \theta_i} J(\theta)$ ; we'd like to check if  $g_i$  is outputting correct derivative values. Let  $\theta^{(i+)} = \theta + \text{EPSILON} \times \vec{e}_i$ , where

$$\vec{e}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

is the  $i$ -th basis vector (a vector of the same dimension as  $\theta$ , with a "1" in the  $i$ -th position and "0"s everywhere else). So,  $\theta^{(i+)}$  is the same as  $\theta$ , except its  $i$ -th element has been incremented by EPSILON. Similarly, let  $\theta^{(i-)} = \theta - \text{EPSILON} \times \vec{e}_i$  be the corresponding vector with the  $i$ -th element decreased by EPSILON. We can now numerically verify  $g_i(\theta)$ 's correctness by checking, for each  $i$ , that:

$$g_i(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2 \times \text{EPSILON}}.$$

When implementing backpropagation to train a neural network, in a correct implementation we will have that

$$\begin{aligned} \nabla_{W^{(l)}} J(W, b) &= \left( \frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \\ \nabla_{b^{(l)}} J(W, b) &= \frac{1}{m} \Delta b^{(l)}. \end{aligned}$$

This result shows that the final block of pseudo-code in Backpropagation Algorithm is indeed implementing gradient descent. To make sure your implementation of gradient descent is correct, it is usually very helpful to use the method described above to numerically compute the derivatives of  $J(W, b)$ , and thereby verify that your computations of  $\left( \frac{1}{m} \Delta W^{(l)} \right) + \lambda W$  and  $\frac{1}{m} \Delta b^{(l)}$  are indeed giving the derivatives you want.

Finally, so far our discussion has centered on using gradient descent to minimize  $J(\theta)$ . If you have implemented a function that computes  $J(\theta)$  and  $\nabla_{\theta} J(\theta)$ , it turns out there are more sophisticated algorithms than gradient descent for trying to minimize  $J(\theta)$ . For example, one can envision an algorithm that uses gradient descent, but automatically tunes the learning rate  $\alpha$  so as to try to use a step-size that causes  $\theta$  to approach a local optimum as quickly as possible. There are other algorithms that are even more sophisticated than this; for example, there are algorithms that try to find an approximation to the Hessian matrix, so that it can take more rapid steps towards a local optimum (similar to Newton's method). A full discussion of these algorithms is beyond the scope of these notes, but one example is the **L-BFGS** algorithm. (Another example is the **conjugate gradient** algorithm.) You will use one of these algorithms in the programming exercise. The main thing you need to provide to these advanced optimization algorithms is that for any  $\theta$ , you have to be able to compute  $J(\theta)$  and  $\nabla_{\theta} J(\theta)$ . These optimization algorithms will then do their own internal tuning of the learning rate/step-size  $\alpha$  (and compute its own approximation to the Hessian, etc.) to automatically search for a value of  $\theta$  that minimizes  $J(\theta)$ . Algorithms such as L-BFGS and conjugate gradient can often be much faster than gradient descent.

Language : 中文

Retrieved from

["http://deeplearning.stanford.edu/wiki/index.php/Gradient\\_checking\\_and\\_advanced\\_optimization"](http://deeplearning.stanford.edu/wiki/index.php/Gradient_checking_and_advanced_optimization)

---

- This page was last modified on 7 April 2013, at 12:40.