



deeplearning.ai

Setting up your
optimization problem

Gradient Checking

Gradient check for a neural network

This is Done so we can ↓ error from gradients calculated

Take $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a big vector θ .

*This is a matrix
⇒ Need to reshape to vector*

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}, \dots, W^{[L]}, b^{[L]}) = J(\theta)$$

Is Already a vector

Concatenate all of these

Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

Big Question

Is $d\theta$ the gradient of $J(\theta)$?

Gradient checking (Grad check) $J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$

for each i : \rightarrow i^{th} component of θ [the level of w]

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

} large Δ calculation

Also $d\theta[i] = \frac{\partial J}{\partial \theta_i}$ (By traditional calculus)

what we hope is

$$d\theta[i] = d\theta_{\text{approx}}[i]$$

How to check if they are approx equal?

Take Euclidean dist & normalize by size of vectors

$$\text{Blah} = \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

} Euclidean dist
} size of vectors

What is the role of normalizing by length?

If Euclidean dist is large/small, the denominator scales it accordingly to small/large (Respectively)

If $\epsilon = 10^{-7}$ Above in the calculation

& suppose $\text{Blah} = 10^{-7}$, then no error

If $\text{Blah} = 10^{-5}$, double check components of θ , none of them should be too large, if large, then possible bug

If $\text{Blah} = 10^{-3}$, there is a bug will have to track down which one made it large



deeplearning.ai

Setting up your
optimization problem

Gradient Checking
implementation notes

Gradient checking implementation notes

- Don't use in training – only to debug

calc of $d\Theta_{\text{approx}}$ [i] is very slow

- If algorithm fails grad check, look at components to try to identify bug.

- Remember regularization.

If you Add $\frac{\lambda}{2m} \leq \|W\|_2$ Remember to include this when finding $\frac{\partial J}{\partial \Theta}$

- Doesn't work with dropout.

→ we're eliminating nodes randomly
⇒ can't find consistent diff as there is no consistent J
 $J = \text{Some Random Subset of } J$

- Run at random initialization; perhaps again after some training.

look at values of the i^{th} $d\Theta_{\text{approx}}$ & $d\Theta$, say for layer 5, we see large diff
then look at decomposition of Θ
 $\Theta = \text{vectorized}(w) + \text{vectorized}(b)$ [In this case, for layer 5]
So if the large diff is being caused by the "b" component
look into how "b" is calculated

$[0.1, 0.1, 0.01, 0.8, 0.9]$
w part b part
⇒ b is contributing to the large diff