

Calculus on Computational Graphs: Backpropagation

Posted on August 31, 2015

Introduction

Backpropagation is the key algorithm that makes training deep models computationally tractable. For modern neural networks, it can make training with gradient descent as much as ten million times faster, relative to a naive implementation. That’s the difference between a model taking a week to train and taking 200,000 years.

Beyond its use in deep learning, backpropagation is a powerful computational tool in many other areas, ranging from weather forecasting to analyzing numerical stability – it just goes by different names. In fact, the algorithm has been reinvented at least dozens of times in different fields (see Griewank (2010) (http://www.math.uiuc.edu/documenta/vol-ismp/52_griewank-andreas-b.pdf)). The general, application independent, name is “reverse-mode differentiation.”

Fundamentally, it’s a technique for calculating derivatives quickly. And it’s an essential trick to have in your bag, not only in deep learning, but in a wide variety of numerical computing situations.

Computational Graphs

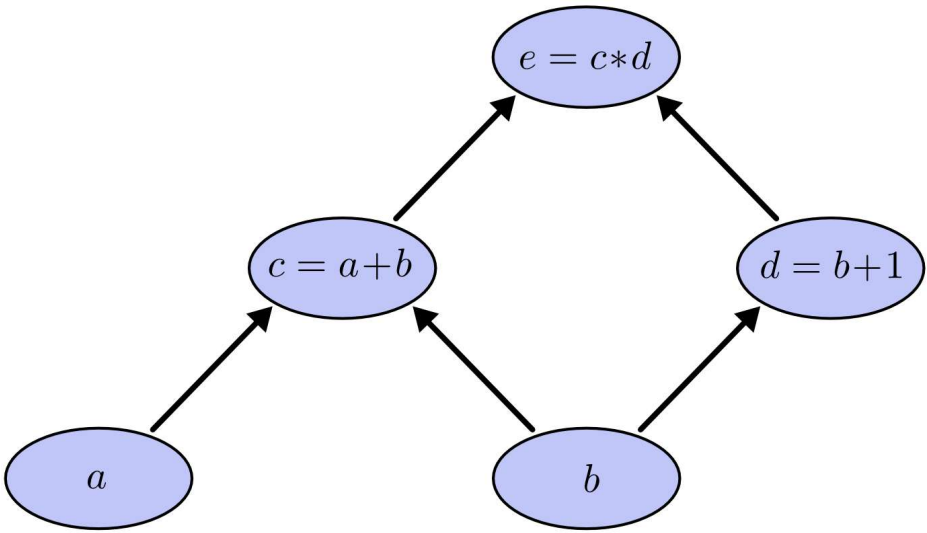
Computational graphs are a nice way to think about mathematical expressions. For example, consider the expression $e = (a + b) * (b + 1)$. There are three operations: two additions and one multiplication. To help us talk about this, let’s introduce two intermediary variables, c and d so that every function’s output has a variable. We now have:

$$c = a + b$$

$$d = b + 1$$

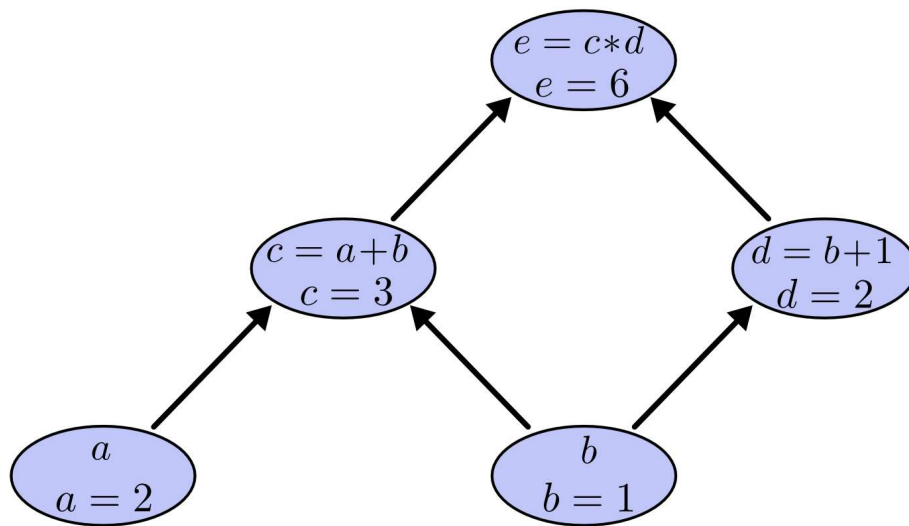
$$e = c * d$$

To create a computational graph, we make each of these operations, along with the input variables, into nodes. When one node’s value is the input to another node, an arrow goes from one to another.



These sorts of graphs come up all the time in computer science, especially in talking about functional programs. They are very closely related to the notions of dependency graphs and call graphs. They’re also the core abstraction behind the popular deep learning framework Theano (<http://deeplearning.net/software/theano/>).

We can evaluate the expression by setting the input variables to certain values and computing nodes up through the graph. For example, let’s set $a = 2$ and $b = 1$:



The expression evaluates to 6.

Derivatives on Computational Graphs

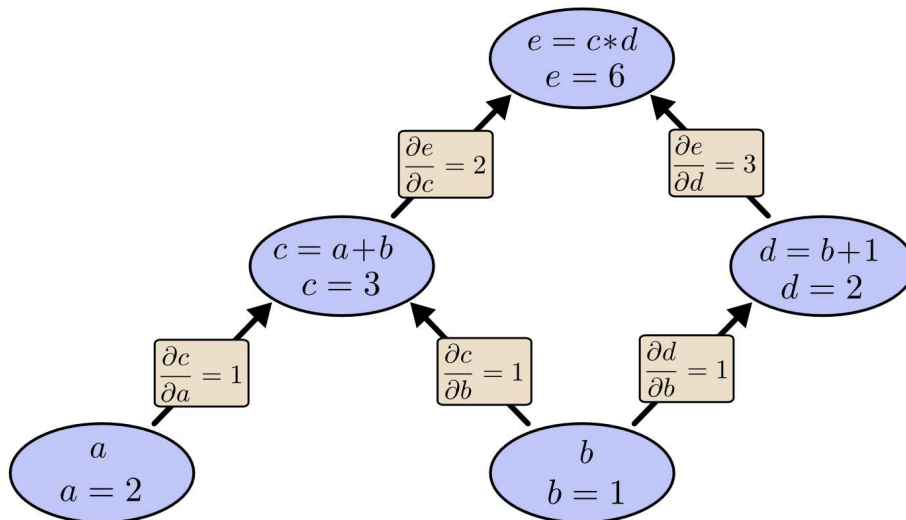
If one wants to understand derivatives in a computational graph, the key is to understand derivatives on the edges. If a directly affects c , then we want to know how it affects c . If a changes a little bit, how does c change? We call this the partial derivative (https://en.wikipedia.org/wiki/Partial_derivative) of c with respect to a .

To evaluate the partial derivatives in this graph, we need the sum rule (https://en.wikipedia.org/wiki/Sum_rule_in_differentiation) and the product rule (https://en.wikipedia.org/wiki/Product_rule):

$$\frac{\partial}{\partial a}(a + b) = \frac{\partial a}{\partial a} + \frac{\partial b}{\partial a} = 1$$

$$\frac{\partial}{\partial u}uv = u \frac{\partial v}{\partial u} + v \frac{\partial u}{\partial u} = v$$

Below, the graph has the derivative on each edge labeled.



What if we want to understand how nodes that aren't directly connected affect each other? Let's consider how e is affected by a . If we change a at a speed of 1, c also changes at a speed of 1. In turn, c changing at a speed of 1 causes e to change at a speed of 2. So e changes at a rate of $1 * 2$ with respect to a .

The general rule is to sum over all possible paths from one node to the other, multiplying the derivatives on each edge of the path together. For example, to get the derivative of e with respect to b we get:

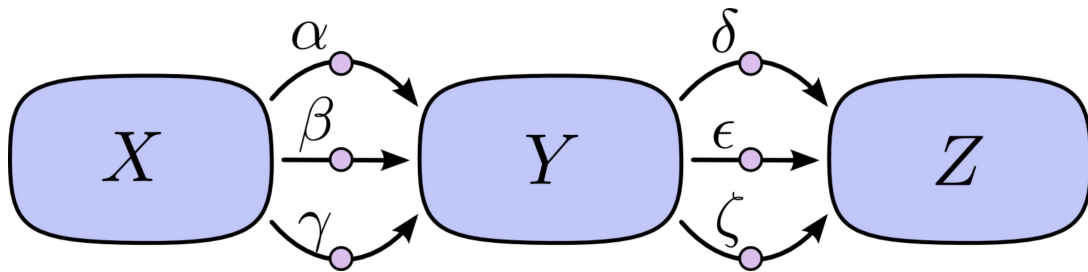
$$\frac{\partial e}{\partial b} = 1 * 2 + 1 * 3$$

This accounts for how b affects e through c and also how it affects it through d .

This general “sum over paths” rule is just a different way of thinking about the multivariate chain rule (https://en.wikipedia.org/wiki/Chain_rule#Higher_dimensions).

Factoring Paths

The problem with just “summing over the paths” is that it's very easy to get a combinatorial explosion in the number of possible paths.



In the above diagram, there are three paths from X to Y , and a further three paths from Y to Z . If we want to get the derivative $\frac{\partial Z}{\partial X}$ by summing over all paths, we need to sum over $3 * 3 = 9$ paths:

$$\frac{\partial Z}{\partial X} = \alpha\delta + \alpha\epsilon + \alpha\zeta + \beta\delta + \beta\epsilon + \beta\zeta + \gamma\delta + \gamma\epsilon + \gamma\zeta$$

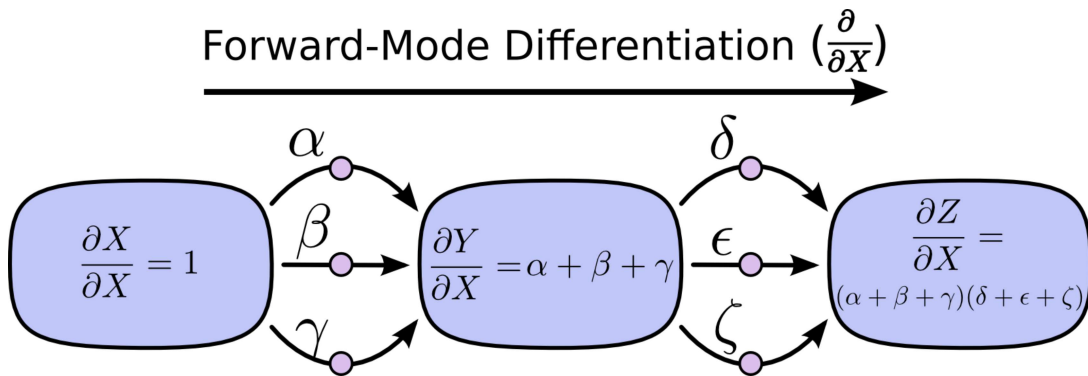
The above only has nine paths, but it would be easy to have the number of paths to grow exponentially as the graph becomes more complicated.

Instead of just naively summing over the paths, it would be much better to factor them:

$$\frac{\partial Z}{\partial X} = (\alpha + \beta + \gamma)(\delta + \epsilon + \zeta)$$

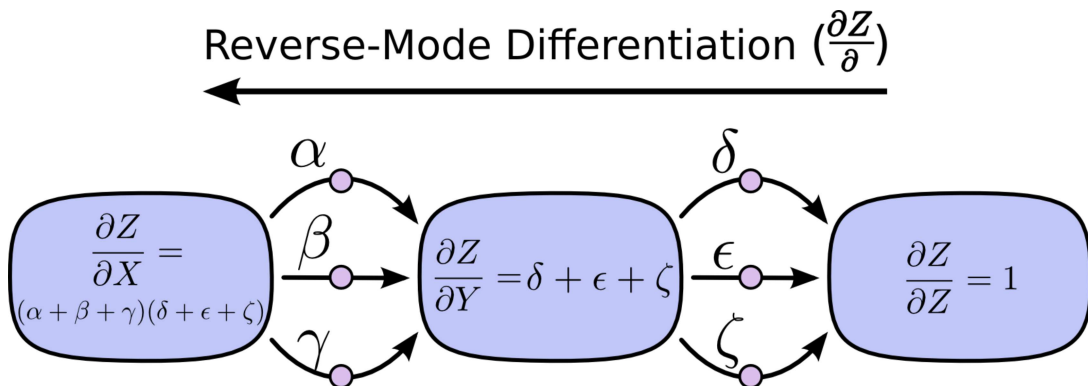
This is where “forward-mode differentiation” and “reverse-mode differentiation” come in. They’re algorithms for efficiently computing the sum by factoring the paths. Instead of summing over all of the paths explicitly, they compute the same sum more efficiently by merging paths back together at every node. In fact, both algorithms touch each edge exactly once!

Forward-mode differentiation starts at an input to the graph and moves towards the end. At every node, it sums all the paths feeding in. Each of those paths represents one way in which the input affects that node. By adding them up, we get the total way in which the node is affected by the input, it’s derivative.



Though you probably didn’t think of it in terms of graphs, forward-mode differentiation is very similar to what you implicitly learned to do if you took an introduction to calculus class.

Reverse-mode differentiation, on the other hand, starts at an output of the graph and moves towards the beginning. At each node, it merges all paths which originated at that node.

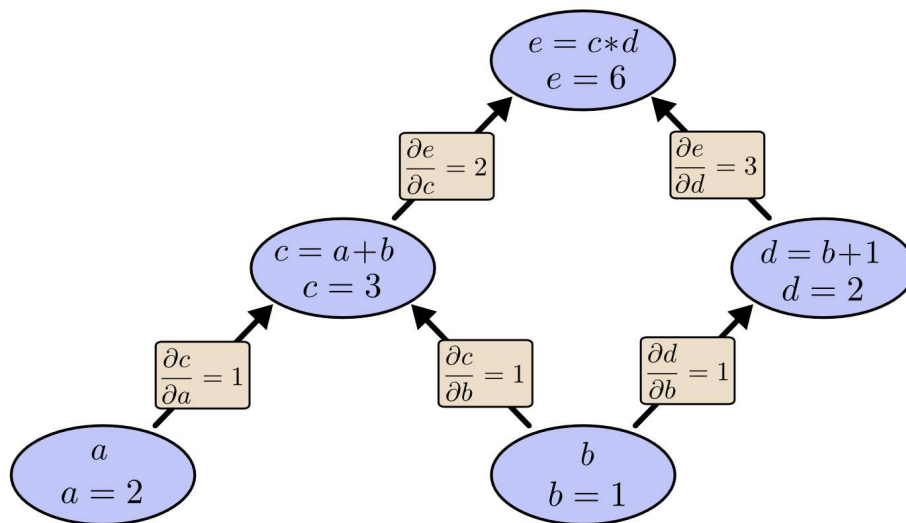


Forward-mode differentiation tracks how one input affects every node. Reverse-mode differentiation tracks how every node affects one output. That is, forward-mode differentiation applies the operator $\frac{\partial}{\partial X}$ to every node, while reverse mode differentiation applies the operator $\frac{\partial Z}{\partial}$ to every node.¹

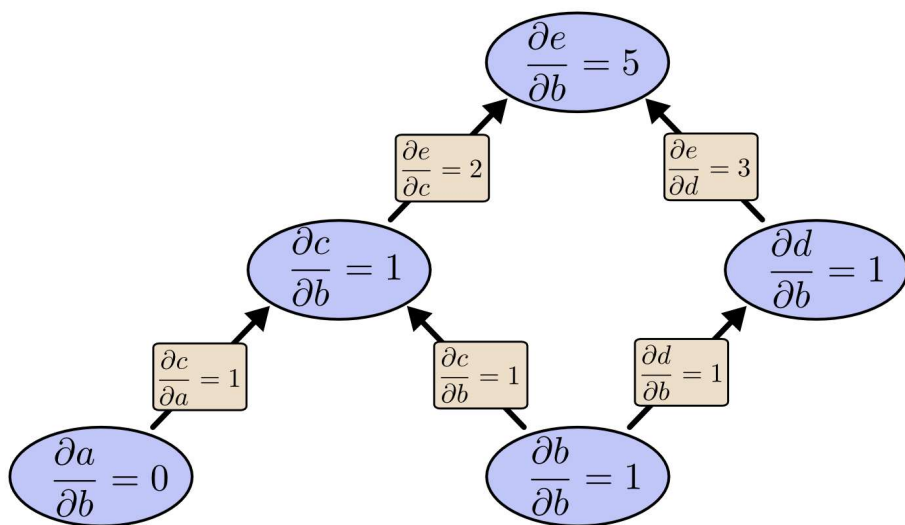
Computational Victories

At this point, you might wonder why anyone would care about reverse-mode differentiation. It looks like a strange way of doing the same thing as the forward-mode. Is there some advantage?

Let’s consider our original example again:

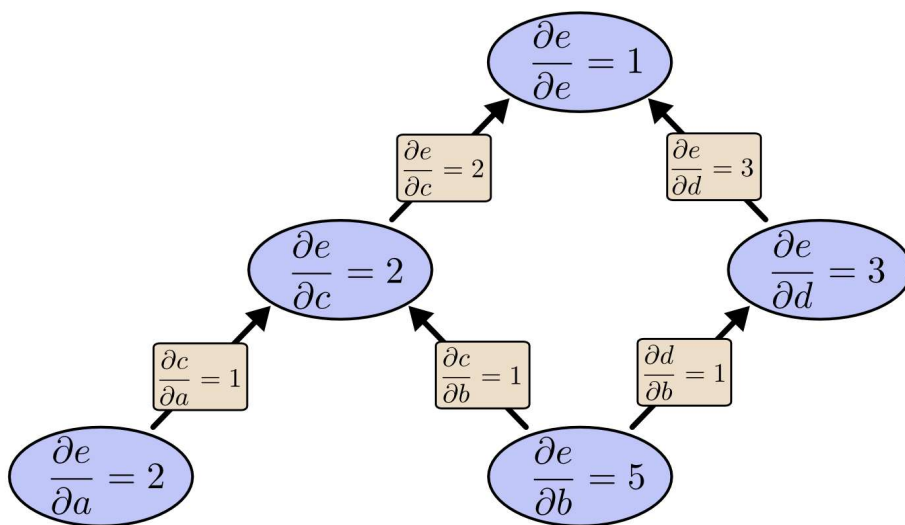


We can use forward-mode differentiation from b up. This gives us the derivative of every node with respect to b .



We've computed $\frac{\partial e}{\partial b}$, the derivative of our output with respect to one of our inputs.

What if we do reverse-mode differentiation from e down? This gives us the derivative of e with respect to every node:



When I say that reverse-mode differentiation gives us the derivative of e with respect to every node, I really do mean *every node*. We get both $\frac{\partial e}{\partial a}$ and $\frac{\partial e}{\partial b}$, the derivatives of e with respect to both inputs. Forward-mode differentiation gave us the derivative of our output with respect to a single input, but reverse-mode differentiation gives us all of them.

For this graph, that's only a factor of two speed up, but imagine a function with a million inputs and one output. Forward-mode differentiation would require us to go through the graph a million times to get the derivatives. Reverse-mode differentiation can get them all in one fell swoop! A speed up of a factor of a million is pretty nice!

