

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-ismmp/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

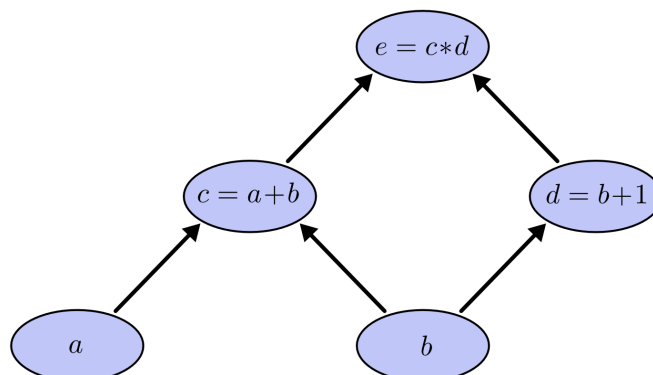
Computational graphs are a nice way to think about mathematical expressions. For example, consider the expression $e = (a + b) * (b + 1)$. There are three operations: two additions and one multiplication. To help us talk about this, let's introduce two intermediary variables, c and d so that every function's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

$$c = a + b$$

$$d = b + 1$$

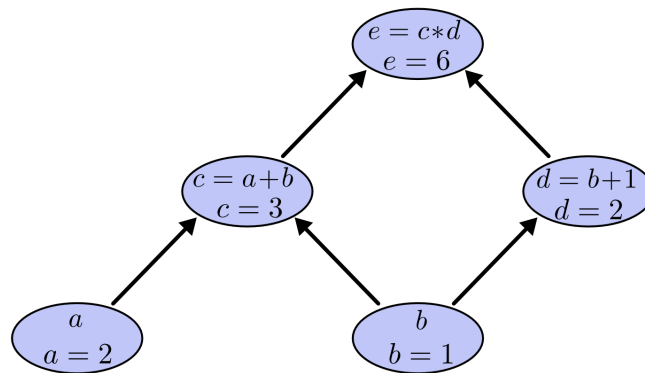
$$e = c * d$$

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

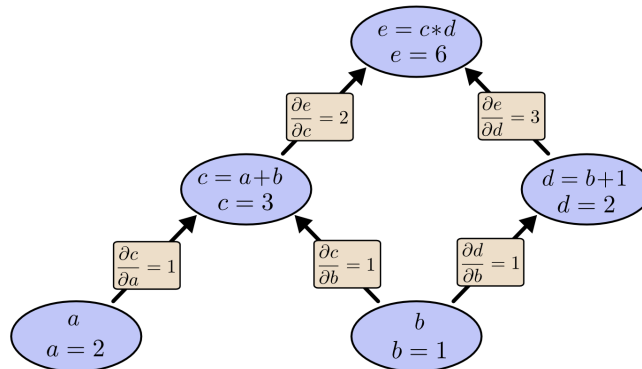
If one wants to understand derivatives in a computational graph, the key is to understand derivatives on the edges. If a directly affects c , then we want to know how it affects c . If a changes a little bit, how does c change? We call this the partial derivative (https://en.wikipedia.org/wiki/Partial_derivative)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

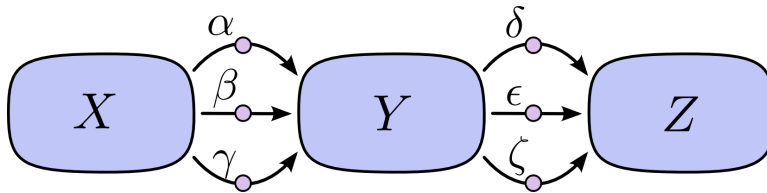
$$\frac{\partial e}{\partial b} = 1 * 2 + 1 * 3$$

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:

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

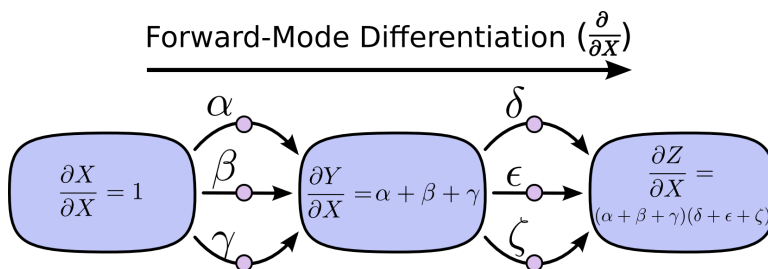
The above only has nine paths, but it would be easy to have the number of paths to grow exponentially as the graph becomes more complicated.

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

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

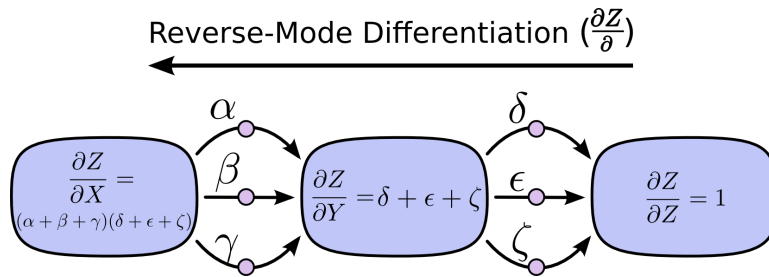
This is where “forward-mode differentiation” and “reverse-mode differentiation” come in. They’re algorithms for efficiently computing the sum by factoring the paths. Instead of summing over all of the paths explicitly, they compute the same sum more efficiently by merging paths back together at every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

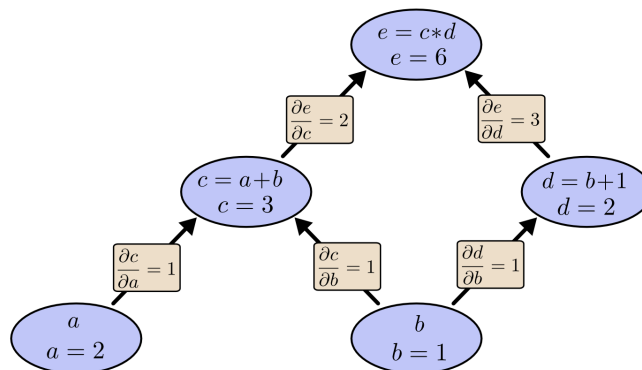


Forward-mode differentiation tracks how one input affects every node. Reverse-mode differentiation tracks how every node affects one output. That is, forward-mode differentiation applies the operator $\frac{\partial}{\partial x}$ to every node, while reverse mode differentiation applies the operator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

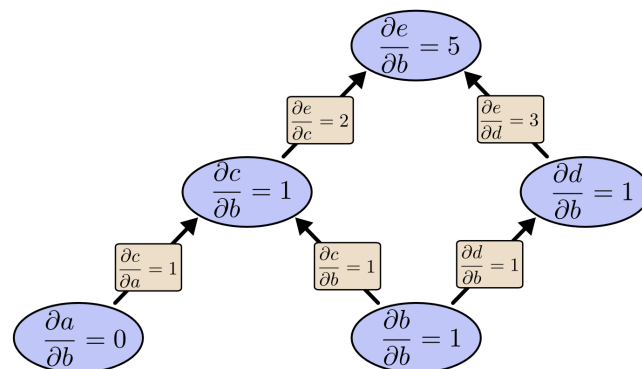
Computational Victories

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:

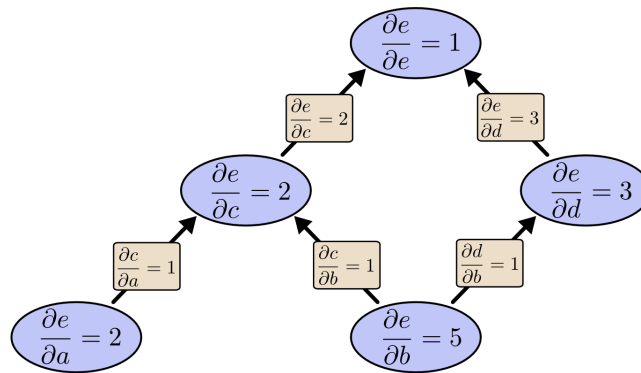


We can use forward-mode differentiation from b up. This gives us the derivative of every node with respect to b .



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!

When training neural networks, we think of the cost (a value describing how bad a neural network performs) as a function of the parameters (numbers describing how the network behaves). We want to calculate the derivatives of the cost with respect to all the parameters, for use in gradient descent (https://en.wikipedia.org/wiki/Gradient_descent). Now, there's often millions, or even tens of millions of parameters in a neural network. So, reverse-mode differentiation, called backpropagation in the context of neural networks, gives us a massive speed up!

(Are there any cases where forward-mode differentiation makes more sense? Yes, there are! Where the reverse-mode gives the derivatives of one output with respect to all inputs, the forward-mode gives us the derivatives of all outputs with respect to one input. If one has a function with lots of outputs, forward-mode differentiation can be much, much, much faster.)

Isn't This Trivial?

When I first understood what backpropagation was, my reaction was: "Oh, that's just the chain rule! How did it take us so long to figure out?" I'm not the only one who's had that reaction. It's true that if you ask "is there a smart way to calculate derivatives in feedforward neural networks?" the answer isn't that difficult.

But I think it was much more difficult than it might seem. You see, at the time backpropagation was invented, people weren't very focused on the feedforward neural networks that we study. It also wasn't obvious that derivatives were the right way to train them. Those are only obvious once you realize you can quickly calculate derivatives. There was a circular dependency.

Worse, it would be very easy to write off any piece of the circular dependency as impossible on casual thought. Training neural networks with derivatives? Surely you'd just get stuck in local minima. And obviously it would be expensive to compute all those derivatives. It's only because we know this approach works that we don't immediately start listing reasons it's likely not to.

That's the benefit of hindsight. Once you've framed the question, the hardest work is already done.

Conclusion

Derivatives are cheaper than you think. That's the main lesson to take away from this post. In fact, they're unintuitively cheap, and us silly humans have had to repeatedly rediscover this fact. That's an important thing to understand in deep learning. It's also a really useful thing to know in other fields, and only more so if it isn't common

knowledge.

Are there other lessons? I think there are.

Backpropagation is also a useful lens for understanding how derivatives flow through a model. This can be extremely helpful in reasoning about why some models are difficult to optimize. The classic example of this is the problem of vanishing gradients in recurrent neural networks.

Finally, I claim there is a broad algorithmic lesson to take away from these techniques. Backpropagation and forward-mode differentiation use a powerful pair of tricks (linearization and dynamic programming) to compute derivatives more efficiently than one might think possible. If you really understand these techniques, you can use them to efficiently calculate several other interesting expressions involving derivatives. We'll explore this in a later blog post.

This post gives a very abstract treatment of backpropagation. I strongly recommend reading Michael Nielsen's chapter on it (<http://neuralnetworksanddeeplearning.com/chap2.html>) for an excellent discussion, more concretely focused on neural networks.

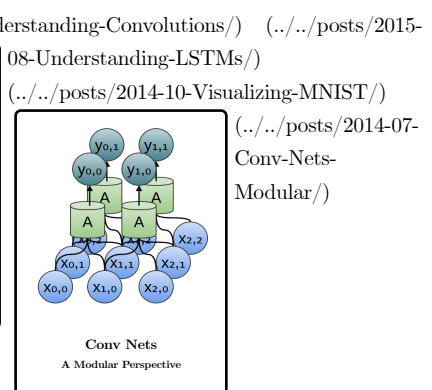
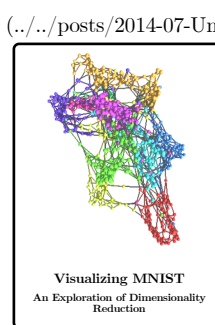
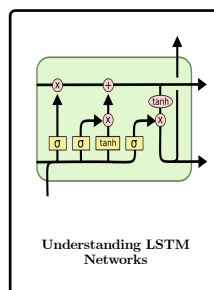
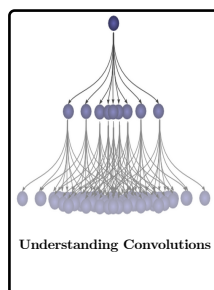
Acknowledgments

Thank you to Greg Corrado (<http://research.google.com/pubs/GregCorrado.html>), Jon Shlens (<https://shlens.wordpress.com/>), Samy Bengio (<http://bengio.abracadoudou.com/>) and Anelia Angelova (<http://www.vision.caltech.edu/anelia/>) for taking the time to proofread this post.

Thanks also to Dario Amodei (<https://www.linkedin.com/pub/dario-amodei/4/493/393>), Michael Nielsen (<http://michaelnielsen.org/>) and Yoshua Bengio (http://www.iro.umontreal.ca/~bengioy/yoshua_en/index.html) for discussion of approaches to explaining backpropagation. Also thanks to all those who tolerated me practicing explaining backpropagation in talks and seminar series!

-
1. This might feel a bit like dynamic programming (https://en.wikipedia.org/wiki/Dynamic_programming). That's because it is!↔
-

More Posts



52 Comments (/posts/2015-08-Backprop/#disqus_thread)

52 Comments colah's blog

Login

Recommend 39

Tweet

Share

Sort by Best



Join the discussion...

LOG IN WITH

OR SIGN UP WITH DISQUS ?

Name

**Ryan Jansekok** · 2 years ago

These posts are absolutely brilliant - and have helped me so much. Thank you for your work!

52 ^ | v · Reply · Share ›

**zachahuy** · 2 years ago

Thank you! This is the best explanation of computational graph and when to use back or forward propagation.

30 ^ | v · Reply · Share ›

**Damodharan J** · 2 years ago

Hi chrisolah! I happened to stumble upon your blog from tensorflow! You happened to answer the very question that was bothering me! About the circular dependency. Thanks a lot for this blog!

22 ^ | v · Reply · Share ›

**yzzzd** · 5 months ago

Thank you so much. This is really helpful.

4 ^ | v · Reply · Share ›

**Yaroslav Bulatov** · 3 years ago

There's another cool algebraic view: for $f(g(h(\dots)))$ the derivative is F^*G^*H where $*$ is matmul and F, G, H are Jacobian matrices. If you have many inputs and one output, f is $\mathbb{R}^n \rightarrow \mathbb{R}^1$, then your last matrix is skinny and tall, then Matrix Chain Multiplication solution tells you to do $(F G)H$, which is reverse mode AD. But if you have many outputs and one input, your H is wide and short, so most efficient is to do $F(G H)$ which is forward mode AD. But also there are cases where neither forward nor reverse mode AD are the most efficient, and those are the "other" solutions of the MCM problem

5 ^ | v · Reply · Share ›

**computereasy** · 3 years ago

Maybe this is irrelevant. But may I ask how you draw the nodes&edges in such a elegant way?

2 ^ | v · Reply · Share ›

**chrisolah** Mod → **computereasy** · 3 years ago

The diagrams in this post (and most of my diagrams) were illustrated in Inkscape.

3 ^ | v · Reply · Share ›

**Omar** → **computereasy** · 3 years ago

I think he is drawing those nodes and edges using LaTeX and TikZ

^ | v · Reply · Share ›

**chrisolah** Mod → **Omar** · 3 years ago

Good guess, but I actually used Inkscape. :)

^ | v · Reply · Share ›

**Emm Kay** · 4 years ago · edited

Hi!

In



, I see why the two nodes aren't really switched or anything. $de/dd = c$ and $de/dc = d$ but Are we adding those derivatives as we go back down the tree?

The way I see it, the derivative for $e=c*d$ with respect to c should be just c , and d for respect to d ? Why doesn't $de/da = 3$?

1 ^ | v • Reply • Share ›



Axio → Emm Kay • 3 years ago

(Using "\$" for the derivative)

$\$e/\$c = \$(c*d)/\$c = d = 2$. Remember, for example that $d(x^7)/dx$ is 7. In our case, we have "d", not 7. You had incorrectly calculated the derivative. Similarly $\$e/\$d = \$(c*d)/\$d = c = 3$.

Then, $\$e/\$a = \$(d*c)/\$a = d*\$c/\$a + c*\$d/\$a = d*\$(a+b)/\$a + c*\$(b+1)/\$a = d*1 + c*0 = 2$.

I hope this helps.

^ | v • Reply • Share ›



Chen Chen → Axio • 3 years ago

Hi Axio,

Why it is not $\$e/\$a = \$e/\$c * \$c/\$a = d*\$c/\$a = 2*1 = 2$? Should it be full derivative instead of partial derivative?

^ | v • Reply • Share ›



Garima Jain • 7 months ago

The best explanation! Thank you so much!

^ | v • Reply • Share ›



Izumi • 7 months ago

By the way, without considering the storage consumption, we can achieve the same time complexity as back propagation using forward-mode. As the blog says, this is dynamic programming (in tree). Let's see it as a tree, so that we can get the dp function: $dp[u] = \sum_{v \in \text{children}(u)} dp[v]$, v is the set of u 's children.

^ | v • Reply • Share ›



Qi Zhang • 8 months ago

This is the best explanation for forward/back propagation I have every seen. No second.

^ | v • Reply • Share ›



Muhammad Usman • a year ago

Best Explanation of Backpropagation

^ | v • Reply • Share ›



Toan Truong • a year ago

As always, your explanation is very intuitive. Thank you very much.

^ | v • Reply • Share ›



Danny Julian • a year ago

Thank you SO MUCH! I cant thank you enough! You have cleared my head a lot!
Thanks!!!!

^ | v • Reply • Share ›



Atul • 2 years ago

Brilliant explanation! Thanks for taking time out to write this.

^ | v • Reply • Share ›

^ | v · [Reply](#) · [Share](#) ›



Enoch Sit · 2 years ago · edited

"Once you've framed the question, the hardest work is already done." like this quote very much
Thank you very much! Please do one blog on dynamic programming !!!

^ | v · [Reply](#) · [Share](#) ›



lei zhao · 2 years ago

Damn! It's so good! Can i quote it and translate it into Chinese to benefit more people?

^ | v · [Reply](#) · [Share](#) ›



Enoch Sit → lei zhao · 2 years ago · edited

Plz do =]

^ | v · [Reply](#) · [Share](#) ›



thecity2 · 2 years ago

Is it more accurate (or actually correct) to say that backprop can be run on trees but not on graphs, in general?
Does backprop still work on graphs that contain loops?

^ | v · [Reply](#) · [Share](#) ›



Shane Lee · 2 years ago

Thanks a lot! I was having troubles understanding backprop and this post makes everything clear!

^ | v · [Reply](#) · [Share](#) ›



Kevin Lu · 3 years ago

I thought this would be interesting to post, it is a worked example with real numbers for the backpropagation algorithm: <https://mattmazur.com/2015/...>

^ | v · [Reply](#) · [Share](#) ›



Shashank Gupta · 3 years ago

I don't have enough words in my vocabulary to thank you for this blog post. This notion of computation graph gave me a nice framework to reason about any complicated Deep Network in general. Thanks again! :)

^ | v · [Reply](#) · [Share](#) ›



Joseph Catrambone · 3 years ago

Far be it for me to look a gift horse in the mouth, but is there any chance you could do a worked example with matrices instead of single elements? I ask because a resource of that sort is quite hard to find online. There's a presentation from IBM which says that it's an open problem (or ill posed), but I have to assume it's solved if it works for backprop.

^ | v · [Reply](#) · [Share](#) ›



Kevin Lu → Joseph Catrambone · 2 years ago

here is an example! <https://cookedsashimi.wordp...>

1 ^ | v · [Reply](#) · [Share](#) ›



Daniel Seita → Joseph Catrambone · 3 years ago

This would also be nice. I'm thinking of trying that at some point. It definitely works; it's how we have to do it for the CS 231n assignments from Stanford. :)

^ | v · [Reply](#) · [Share](#) ›



Pedro Marcal · 3 years ago

Basically you have to solve a set of quasi-linear set of equations many times with the residuals,

$$R_j = A_{ji} \cdot W_i - B_j$$

The back-propagation process is the only one that can be used because the equations are non-positive definite. In most physics problems such equations are positive definite so that a Gauss-Seidel procedure can be used where only the diagonal terms A_{ii} are used as divisors.

^ | v · [Reply](#) · [Share](#) ›



Afreen · 3 years ago

Thank you for this intuitive explanation of backpropagation.

^ | v · [Reply](#) · [Share](#) ›

**tobe433** · 3 years ago

Great. That helps a lot!

^ | v · Reply · Share ›

**Alex Telfar** · 3 years ago

Holy crap, autograd is backprop... and vice versa. I have been using both and didnt even realise. Awesome

^ | v · Reply · Share ›

**Chanchana Sornsoontorn** → Alex Telfar · 3 years ago · edited

Backprop is not Adagrad.

Backprop is a way to get gradients. What you do with the gradients is up to you. One simple approach is SGD, and the rest like adagrad, adam, adadelta, RMSprop are just a modification to SGD.

Backprop return gradients

Adagrad use gradients to optimize weights.

^ | v · Reply · Share ›

**Joseph Catrambone** → Chanchana Sornsoontorn · 3 years ago

Not AdaGrad, AUTOgrad. Autograd _is_ generalized backprop.

Adagrad and autograd are different things. The parent commenter is correct.

3 ^ | v · Reply · Share ›

**Chanchana Sornsoontorn** → Joseph Catrambone · 3 years ago

OK. Thanks it's my fault.

^ | v · Reply · Share ›

**Axio** · 3 years ago

That was very insightful. Thanks for taking the time to explain in so much details.

^ | v · Reply · Share ›

**Ivan Kryukov** · 4 years ago

Great article, thank you so much!

In the conclusion, you mention "surely, you'd just get stuck at local minima", indicating that this is not the case.

Are there any general principles that guarantee that you would not?

Can you perhaps provide a few pointers on this?

Thanks!

^ | v · Reply · Share ›

**chrisolah** Mod → Ivan Kryukov · 4 years ago

When we train neural networks, we don't find the global minima. Actually, we also don't even find a local minima: recent research suggests we get stuck on saddle points (eg. see some Ian Goodfellow and collaborators recent papers).

But, surprisingly, things still work really well! What's going on? Well, we still find good enough points that neural networks work well, even if they aren't minima.

2 ^ | v · Reply · Share ›

**Ivan Kryukov** → chrisolah · 4 years ago

Thank you.

I will look into Ian's papers.

^ | v · Reply · Share ›

**Amr El-Desoky Mousa** → Ivan Kryukov · 3 years ago

I would like to also point that there some number of techniques that aim at preventing the network from getting stuck at local minima or saddle points, like learning rate tuning and pre-training.

^ | v · Reply · Share ›

**Dp** · 4 years agoInsightful! In fact this post shares the viewpoint with <http://neuralnetworksanddee...>, Chap2, section The Big Picture.

^ | v · Reply · Share ›



mrityunjay pande · 4 years ago

If we find all paths from target node to all source nodes, that will also give all paths between target and intermediate nodes. for a feed forward, adding the edges is ok, for speedup but for recurrent, is all path a better option for back propagation.

^ | v · Reply · Share ›



Jack Scully · 4 years ago

Great post, love this blog. Keep it up!

^ | v · Reply · Share ›



OB · 4 years ago

So, if we have small number of input units and huge number of output units, forward mode works more efficient than the back propagation (reverse mode), and it's almost always not the case (# of input units are huge, and # of output units relatively few)?

Thanks

^ | v · Reply · Share ›



chrisolah Mod → **OB** · 4 years ago

If you want the derivatives of all the outputs with respect to all the inputs, yes.

In machine learning, we usually only care about one output (the cost) with respect to lots of inputs (the parameters) so backprop is great.

2 ^ | v · Reply · Share ›



OB → **chrisolah** · 4 years ago

Thank you, Chris.

^ | v · Reply · Share ›



Kunal Bohra · 4 years ago

Loved it.

^ | v · Reply · Share ›



Alejandro Martínez · 4 years ago

Excellent.

Thank you for effort in simplified .



^ | v · Reply · Share ›



Masaki Nakada · 4 years ago · edited

Loved this article! Thank you so much Christopher!

^ | v · Reply · Share ›

[Load more comments](#)

ALSO ON COLAH'S BLOG

Visual Information Theory

1 comment · 4 years ago



Darwin Kim — I think there's a much simpler way to prove that the optimal code makes the cost equal to the word frequency. However, many of the clever tricks and

Collaboration & Credit Principles

2 comments · 3 months ago



odedbd — This is great; I think most of these recommendations can apply, with relevant modification, to other forms of collaboration, beyond scientific

Visualizing Representations: Deep Learning and Human Beings

1 comment · 3 years ago






Andrew Brereton — Have you ever read Blindsight by Peter Watts? You might enjoy it. Great blog btw.

Calculus on Computational Graphs: Backpropagation

1 comment · 4 years ago



atomicthumbs — Dang it, I KNEW I shouldn't have chosen Business Math!

 [Subscribe](#)  [Add Disqus to your site](#)[Add Disqus](#)[Add](#)  [Disqus' Privacy Policy](#)[Privacy Policy](#)[Privacy Policy](#)[Privacy Policy](#)