

# COSE474-Deep Learning

## 2.5. Automatic Differentiation

Fortunately all modern deep learning frameworks take this work off our plates by offering automatic differentiation (often shortened to autograd). As we pass data through each successive function, the framework builds a computational graph that tracks how each value depends on others. To calculate derivatives, automatic differentiation works backwards through this graph applying the chain rule. The computational algorithm for applying the chain rule in this fashion is called backpropagation.

While autograd libraries have become a hot concern over the past decade, they have a long history. In fact the earliest references to autograd date back over half of a century :cite:Wengert.1964. The core ideas behind modern backpropagation date to a PhD thesis from 1980 :cite:Speelpenning.1980 and were further developed in the late 1980s :cite:Griewank.1989. While backpropagation has become the default method for computing gradients, it is not the only option. For instance, the Julia programming language employs forward propagation :cite:Revels.Lubin.Papamarkou.2016. Before exploring methods, let's first master the autograd package.

```
import torch
```

### 2.5.1. A Simple Function

Let's assume that we are interested in (**differentiating the function  $y = 2\mathbf{x}^\top \mathbf{x}$  with respect to the column vector  $\mathbf{x}$ .**) To start, we assign  $\mathbf{x}$  an initial value.

```
x = torch.arange(4.0)
x
→ tensor([0., 1., 2., 3.]
```

**[Before we calculate the gradient of  $y$  with respect to  $\mathbf{x}$ , we need a place to store it.]** In general, we avoid allocating new memory every time we take a derivative because deep learning requires successively computing derivatives with respect to the same parameters a great many times, and we might risk running out of memory. Note that the gradient of a scalar-valued function with respect to a vector  $\mathbf{x}$  is vector-valued with the same shape as  $\mathbf{x}$ .

```
# Can also create x = torch.arange(4.0, requires_grad=True)
x.requires_grad_(True)
x.grad # The gradient is None by default
```

We now calculate our function of  $\mathbf{x}$  and assign the result to  $y$ .

```
y = 2 * torch.dot(x, x)
y
→ tensor(28., grad_fn=<MulBackward0>)
```

We can now take the gradient of  $y$  with respect to  $\mathbf{x}$  by calling its backward method. Next, we can access the gradient via  $\mathbf{x}$ 's grad attribute.

```
y.backward()
x.grad
→ tensor([ 0.,  4.,  8., 12.]
```

**(We already know that the gradient of the function  $y = 2\mathbf{x}^\top \mathbf{x}$  with respect to  $\mathbf{x}$  should be  $4\mathbf{x}$ .)** We can now verify that the automatic gradient computation and the expected result are identical.

```
x.grad == 4 * x
→ tensor([ True,  True,  True,  True])
```

Now let's calculate another function of  $x$  and take its gradient. Note that PyTorch does not automatically reset the gradient buffer when we record a new gradient. Instead, the new gradient is added to the already-stored gradient. This behavior comes in handy when we want to optimize the sum of multiple objective functions. To reset the gradient buffer, we can call `x.grad.zero_()` as follows:

```
x.grad.zero_() # Reset the gradient
y = x.sum()
y.backward()
x.grad
```

→ tensor([1., 1., 1., 1.])

Start coding or [generate](#) with AI.

## ✓ 2.5.2. Backward for Non-Scalar Variables

When  $y$  is a vector, the most natural representation of the derivative of  $y$  with respect to a vector  $x$  is a matrix called the Jacobian that contains the partial derivatives of each component of  $y$  with respect to each component of  $x$ . Likewise, for higher-order  $y$  and  $x$ , the result of differentiation could be an even higher-order tensor.

While Jacobians do show up in some advanced machine learning techniques, more commonly we want to sum up the gradients of each component of  $y$  with respect to the full vector  $x$ , yielding a vector of the same shape as  $x$ . For example, we often have a vector representing the value of our loss function calculated separately for each example among a batch of training examples. Here, we just want to sum up the gradients computed individually for each example.

Because deep learning frameworks vary in how they interpret gradients of non-scalar tensors, PyTorch takes some steps to avoid confusion. Invoking `backward` on a non-scalar elicits an error unless we tell PyTorch how to reduce the object to a scalar. More formally, we need to provide some vector  $\mathbf{v}$  such that `backward` will compute  $\mathbf{v}^\top \partial_x \mathbf{y}$  rather than  $\partial_x \mathbf{y}$ . This next part may be confusing, but for reasons that will become clear later, this argument (representing  $\mathbf{v}$ ) is named `gradient`. For a more detailed description, see Yang Zhang's [Medium post](#).

```
x.grad.zero_()
y = x * x
y.backward(gradient=torch.ones(len(y))) # Faster: y.sum().backward()
x.grad
```

→ tensor([0., 2., 4., 6.])

## ✓ 2.5.3. Detaching Computation

Sometimes, we wish to move some calculations outside of the recorded computational graph. For example, say that we use the input to create some auxiliary intermediate terms for which we do not want to compute a gradient. In this case, we need to detach the respective computational graph from the final result. The following toy example makes this clearer: suppose we have  $z = x * y$  and  $y = x * x$  but we want to focus on the direct influence of  $x$  on  $z$  rather than the influence conveyed via  $y$ . In this case, we can create a new variable  $u$  that takes the same value as  $y$  but whose provenance (how it was created) has been wiped out. Thus  $u$  has no ancestors in the graph and gradients do not flow through  $u$  to  $x$ . For example, taking the gradient of  $z = x * u$  will yield the result  $u$ , (not  $3 * x * x$  as you might have expected since  $z = x * x * x$ ).

```
x.grad.zero_()
y = x * x
u = y.detach()
z = u * x
```

```
z.sum().backward()
x.grad == u
```

→ tensor([True, True, True, True])

Note that while this procedure detaches  $y$ 's ancestors from the graph leading to  $z$ , the computational graph leading to  $y$  persists and thus we can calculate the gradient of  $y$  with respect to  $x$ .

```
x.grad.zero_()
y.sum().backward()
```

```
x.grad == 2 * x
```

```
→ tensor([True, True, True, True])
```

## ✓ 2.5.4. Gradients and Python Control Flow

So far we reviewed cases where the path from input to output was well defined via a function such as  $z = x * x * x$ . Programming offers us a lot more freedom in how we compute results. For instance, we can make them depend on auxiliary variables or condition choices on intermediate results. One benefit of using automatic differentiation is that even if building the computational graph of a function required passing through a maze of Python control flow (e.g., conditionals, loops, and arbitrary function calls), we can still calculate the gradient of the resulting variable. To illustrate this, consider the following code snippet where the number of iterations of the while loop and the evaluation of the if statement both depend on the value of the input  $a$ .

```
def f(a):
    b = a * 2
    while b.norm() < 1000:
        b = b * 2
    if b.sum() > 0:
        c = b
    else:
        c = 100 * b
    return c
```

Below, we call this function, passing in a random value, as input. Since the input is a random variable, we do not know what form the computational graph will take. However, whenever we execute  $f(a)$  on a specific input, we realize a specific computational graph and can subsequently run backward.

```
a = torch.randn(size=(), requires_grad=True)
d = f(a)
d.backward()
```

Even though our function  $f$  is, for demonstration purposes, a bit contrived, its dependence on the input is quite simple: it is a linear function of  $a$  with piecewise defined scale. As such,  $f(a) / a$  is a vector of constant entries and, moreover,  $f(a) / a$  needs to match the gradient of  $f(a)$  with respect to  $a$ .

```
a.grad == d / a
```

```
→ tensor(True)
```

Dynamic control flow is very common in deep learning. For instance, when processing text, the computational graph depends on the length of the input. In these cases, automatic differentiation becomes vital for statistical modeling since it is impossible to compute the gradient a priori.

## ✓ 2.5.5. Discussion

You have now gotten a taste of the power of automatic differentiation. The development of libraries for calculating derivatives both automatically and efficiently has been a massive productivity booster for deep learning practitioners, liberating them so they can focus on less menial. Moreover, autograd lets us design massive models for which pen and paper gradient computations would be prohibitively time consuming. Interestingly, while we use autograd to optimize models (in a statistical sense) the optimization of autograd libraries themselves (in a computational sense) is a rich subject of vital interest to framework designers. Here, tools from compilers and graph manipulation are leveraged to compute results in the most expedient and memory-efficient manner.

For now, try to remember these basics: (i) attach gradients to those variables with respect to which we desire derivatives; (ii) record the computation of the target value; (iii) execute the backpropagation function; and (iv) access the resulting gradient.

Start coding or [generate](#) with AI.

