

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x

```

Running the above script results in the following error:

```

RuntimeError: grad_fn is None for input tensor at index 0

```

This is because gradients can't be computed with respect to scalar values by definition. You can't really differentiate a scalar with respect to another scalar. This is why scalars are often called **Scalars**, the dimension of which is beyond the scope of this article.

There are two ways to overcome this:

If you just make a small change in the above code setting `z` to be the sum of all the errors, our problem is solved.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))

```

Once that's done, you can access the gradient by calling the `.grad` attribute of `z`.

Second way is for some reason you have to directly set `z` instead of `sum`, function, you can pass a `grad_fn` of size of shape of the tensor you are trying to call backward with.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)

```

Note how `my_fn` used to take incoming gradients as it's input. Doing the above makes the `z` gradient think that incoming gradient is just `Tensor` of `size(1)` and hence it can't be differentiated.

In this case, we can have gradients for every `z`, and we can update them using optimization algorithms of our choice.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)
    z.backward()

```

And here it is.

## How are PyTorch's graphs different from TensorFlow graphs

PyTorch creates something called **Dynamic Computation Graph**, which means the graph is created on the fly.

Until the `backward()` of a Variable is called, there exists no node for the tensor `0.0`'s `grad_fn` in the graph.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)
    z.backward()

```

The graph is created as a result of `backward` of many `Tensor`s being called. It is a linked list of operations and their gradients, and introduced values for computing gradients later. Whenever `backward` is called, as the gradients are computed, these before the node (which is `z`) are converted into `grad_fn` of `z`. These `grad_fn`s are then passed to `backward` through it since it has been holding values to compute gradients on its own.

Now we will discuss the issue of `grad_fn`, the `non_leaf` buffers from the previous run will be shared, while the `non_leaf` nodes buffer will be created again.

If we do `backward` more than once on a graph with non leaf nodes, you'll be met with the following error:

```

RuntimeError: grad_fn is None for input tensor at index 0

```

This is because the `non_leaf` buffers get destroyed the first time `backward` is called and hence, there is no path to update the `grad_fn` when `backward` is invoked the second time. We can avoid this non-leaf buffer destroying behavior by setting `retain_graph=True` as argument to the `backward` function.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)
    z.backward(retain_graph=True)

```

If we do the above, we are able to backpropagate again through the same graph and the gradients will be accumulated, i.e., the next backpropagation, the gradients will be added to those already stored in the previous back pass.

This is in contrast to the **Static Computation Graphs**, such as **TensorFlow**, where the graph is defined before running the program. That is the graph is "frozen" by feeding values to the predefined graph.

The dynamic graph paradigm allows you to make changes to your network architecture during training, as a graph is created only when a piece of code is run.

This means a graph may be redefined during the lifetime for a program since products have a define it beforehand.

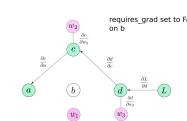
This, however, is not possible with static graphs when graphs are created before running the program, and rarely modified.

TensorFlow makes debugging very easy since it's easier to locate the source of error.

## Some Tricks of Trade

`requires_grad`

This is to attribute the `grad_fn`. By default, it's `None`. It prevents back propagation to those layers, and stops any further updating parameters while training. You can set it at the `requires_grad` to `False`, and those layers are won't participate in the computation graph.



Then, no gradient would be propagated to them, or to those layers which depend upon those layers for gradient flow `requires_grad`. When set to `True`, it propagates gradients to all the layers. If you've created a operation like `grad_fn` set to `True`, as well as the `grad_fn`.

`torch.no_grad()`

When we are computing gradients, we need to cache input values, and intermediate values to compute gradients. This is done by `grad_fn` of a particular node. The `grad_fn` is  $\lambda = \text{out} \rightarrow \text{out} \times \text{grad}$  and  $\text{out} \times \text{grad}$  is responsible for gradient computation along the backward pass. This affects the memory footprint of the network.

While we are performing inference, we don't compute gradients, and thus, we don't need to cache input values, and intermediate values during inference as it will lead to undue consumption of memory.

PyTorch offers a context manager, called `no_grad`, just for this purpose.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)
    z.backward()

```

No graph is defined for operations enclosed under this context manager.

## Conclusion

Understanding how `backward` and computation graph works can make life easier. While we are computing gradients, we need to cache input values, and intermediate values to compute gradients. This is done by `grad_fn` of a particular node. The `grad_fn` is  $\lambda = \text{out} \rightarrow \text{out} \times \text{grad}$  and  $\text{out} \times \text{grad}$  is responsible for gradient computation along the backward pass. This affects the memory footprint of the network.

While we are performing inference, we don't compute gradients, and thus, we don't need to cache input values, and intermediate values during inference as it will lead to undue consumption of memory.

PyTorch offers a context manager, called `no_grad`, just for this purpose.

```

graph LR
    x((x)) --> y((y))
    y --> z((z))
    z --> w((w))
    w --> x
    z == sum((sum))
    z.grad_fn = my_fn(z)
    z.backward()

```

No graph is defined for operations enclosed under this context manager.

## Further Reading

1. [Chain Rule](#)
2. [Backpropagation](#)

[Add speed and simplicity to your Machine Learning workflow today](#)

[GET STARTED](#) [CONTACT SALES](#)

[View Case Study](#)

[View Testimonials](#)

[Contact Us](#)

[Request a Demo](#)

[Download Brochure](#)

[Read Whitepaper](#)

[Watch Video](#)

[Read Case Study](#)

[Read Whitepaper](#)

[Read Case Study](#)

[Read Whitepaper](#)</p