



torch.func

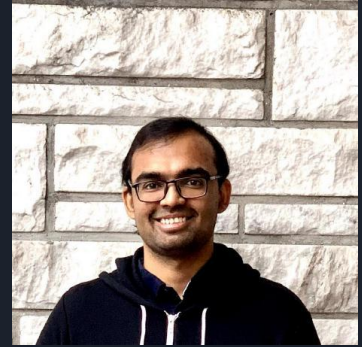
Functional Transforms in PyTorch

Toronto Machine Learning
Summit 2024

@shagunsodhani



About Me



1. Tech Lead and Staff Research Engineer @ Meta AI
2. Training foundation models to develop adaptive neural interfaces



Agenda

1. Quick overview of PyTorch
2. Why torch.func
3. Overview of torch.func API
4. Where can you start using torch.func today
5. Gotachs to look out for
6. Questions and Answers



PyTorch

1. Open-source Machine Learning framework
2. Provides Numpy-like arrays with GPU acceleration
3. Enables training deep neural networks
4. Well-known for its ease of use



Why torch.func

Following use cases are tricky to do in PyTorch

1. computing per-sample-gradients
2. running model ensembles on a single machine
3. efficiently batching together tasks in the inner-loop of MAML-like algorithms
4. efficiently computing Jacobians and Hessians (with or without batching)

These can be supported by introducing composable function transforms



Why torch.func

Composable Function Transforms



Why torch.func

Composable **Function Transforms**

1. Higher-order function that accepts functions as input and returns function as output.



Why torch.func

Composable **Function Transforms**

1. Higher-order function that accepts functions as input and returns function as output.
2. Examples include auto-differentiation transforms, `grad(f)` that returns a function that computes the gradient of `f` or vectorization/batching transform, `vmap(f)` that returns a function that computes `f` over batches of inputs.



Why torch.func

Composable Function Transforms

These function transforms can compose with each other arbitrarily. For example, composing `vmap(grad(f))` computes per-sample-gradients!



Why torch.func

1. In general, having "pure" functions makes it easier to compose them.
 - a. Functions that always produce the same output for the same input and have no side effects
2. In PyTorch, commonly used constructs like modules are stateful.
3. `torch.func` makes it easier to embrace the functional programming style, which in-turn simplifies some workflows easier.



API | `vmap`

Given a function `func` that runs on a single example, we can lift it to a function that can take batches of examples with `vmap(func)`

`vmap(func)` adds a dimension to all tensor operations in `func`

It can be invoked as `vmap(func)(*inputs)`

API | vmap

```
1 x: torch.Tensor = torch.randn(100)
2 y: torch.Tensor = torch.randn(100)
3 x_dot_y = torch.dot(x, y)
4 print(f"{x_dot_y}")
```

```
x_dot_y=tensor(-3.0892)
```

API | vmap

```
1 x: torch.Tensor = torch.randn(100)
2 y: torch.Tensor = torch.randn(100)
3 x_dot_y = torch.dot(x, y)
4 print(f"{x_dot_y}")
```

```
x_dot_y=tensor(-3.0892)
```

```
1 x: torch.Tensor = torch.randn(10, 100)
2 y: torch.Tensor = torch.randn(10, 100)
3 x_dot_y = torch.dot(x, y)
```



RuntimeError:1D tensors expected, but got 2D and 2D tensors



API | `vmap`

```
1 x: torch.Tensor = torch.randn(10, 100)
2 y: torch.Tensor = torch.randn(10, 100)
3 x_dot_y = torch.ones(10)
4 for i in range(x.shape[0]):
5     x_dot_y[i] = torch.dot(x[i], y[i])
6 print(f"{x_dot_y}")
```

```
x_dot_y=tensor([  8.2551,  13.3984, -11.2821, ...,
                -7.5680, -3.9727, -0.6121])
```


API | vmap

```
1  batched_dot_product = torch.func.vmap(torch.dot)
2  x_dot_y_using_vmap = batched_dot_product(x, y)
3  print(f"{x_dot_y_using_vmap=}")
```

```
x_dot_y_using_vmap=tensor([  8.2551,  13.3984, -11.282
                          -7.5680,  -3.9727,  -0.6121])
```

```
1  assert torch.allclose(x_dot_y, x_dot_y_using_vmap)
```




API | `grad`

`grad` operator helps computing gradients of `func`.

This operator can be nested to compute higher-order gradients.



API | grad

```
1  sin_x = lambda x: torch.sin(x)
2  grad_sin_x = torch.func.grad(sin_x)
3  x = torch.randn(1)
4  assert torch.allclose(grad_sin_x(x), x.cos())
```



API | grad

```
1 grad_grad_sin_x = torch.func.grad(grad_sin_x)
2 assert torch.allclose(grad_grad_sin_x(x), -x.sin())
```



API | `grad` + `vmap`

When composed with `vmap`, `grad` can be used to compute per-sample-gradients:

API | grad + vmap

```
1  from torch.func import grad, vmap
2
3  batch_size, feature_size = 3, 5
4
5  def model(weights: torch.Tensor, feature_vec: torch.Tensor) -> torch.Tensor:
6      return feature_vec.dot(weights).relu()
7
8  def compute_loss(weights: torch.Tensor, example: torch.Tensor, target: torch.Tensor) -> torch.Tensor:
9      y = model(weights, example)
10     return ((y - target) ** 2).mean() # MSE Loss'
11
12  weights = torch.randn(feature_size, requires_grad=True)
13  examples = torch.randn(batch_size, feature_size)
14  targets = torch.randn(batch_size)
```

API | grad + vmap

```
1 inputs = (weights, examples, targets)
2 grad_of_loss = grad(compute_loss)
3 grad_of_loss_per_sample = vmap(grad_of_loss, in_dims=(None, 0, 0))
4
5 grad_weight_per_example = grad_of_loss_per_sample(*inputs)
6 print(grad_weight_per_example)
```

```
tensor([[ -0.0000,  0.0000, -0.0000, -0.0000, -0.0000],
        [-0.0352,  0.0429,  0.0334, -0.0174,  0.0097],
        [ 0.0000, -0.0000, -0.0000,  0.0000, -0.0000]], grad_fn=<MulBackward0>)
```



API | `functional_call`

Performs a functional call on the module by replacing the module parameters and buffers with the provided ones.



API | `functional_call`

```
1 x = torch.randn(4, 3)
2 t = torch.randn(4, 3)
3
4 model = nn.Linear(3, 3)
5
6 params = dict(model.named_parameters())
7 y = functional_call(model, params, x)
8
9 assert torch.allclose(y, model(x))
```




API | `functional_call`

```
1  def compute_loss(  
2  |     params: dict[str, torch.Tensor], x: torch.Tensor, t: torch.Tensor  
3  ) -> torch.Tensor:  
4  |     y = functional_call(model, params, x)  
5  |     return nn.functional.mse_loss(y, t)  
6  
7  
8  grad_of_loss = grad(compute_loss)  
9  grad_weights = grad_of_loss(dict(model.named_parameters()), x, t)
```



API | `functional_call`

```
1 num_models = 5
2 batch_size = 64
3 in_features, out_features = 3, 3
4 models = [torch.nn.Linear(in_features, out_features) for i in range(num_models)]
5 data = torch.randn(batch_size, 3)
```



API | `stack_module_state`

```
1  def forward_call(params, buffers, data):
2      |   return torch.func.functional_call(models[0], (params, buffers), data)
3
4  vmap_forward_call = vmap(forward_call, (0, 0, None))
5
6  params, buffers = torch.func.stack_module_state(models)
7
8  output = vmap_forward_call(params, buffers, data)
9
10 assert output.shape == (num_models, batch_size, out_features)
```



API

Other examples include

- `vjp` (vector jacobian product)
- `jvp` (jacobian vector product)
- `hessian`
- ...

Gotchas

Using PyTorch `torch.no_grad` together **with** `grad`.

Case 1: Using `torch.no_grad` inside a function:

```
>>> def f(x):
>>>     with torch.no_grad():
>>>         c = x ** 2
>>>     return x - c
```

In this case, `grad(f)(x)` will respect the inner `torch.no_grad`.

Case 2: Using `grad` inside `torch.no_grad` context manager:

```
>>> with torch.no_grad():
>>>     grad(f)(x)
```

In this case, `grad` will respect the inner `torch.no_grad`, but **not** the outer one. This **is** because `grad` **is** a “function transform”: its result should **not** depend on the result of a context manager outside of `f`.



Gotchas

1. Functions with side-effects / global effects can be problematic
2. `vmap` does not work with some inplace operations
3. `vmap` does not work with some data dependent conditionals
4. [Batchnorm](#) requires special handling
5. For more gotachs, checkout [this](#)



Thank you!

<https://shagunsodhani.com/talks/>

Toronto Machine Learning
Summit 2024

@shagunsodhani