

Autodiff & PyTorch for the gifted amateur

Wissen für Morgen



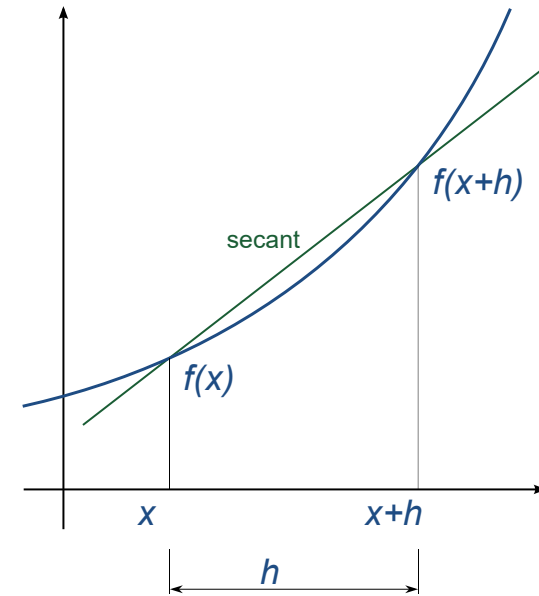
Autodiff – *Why should I care?*

- Reverse mode automatic differentiation (***autodiff***):
 - One of most important algorithms of the last century
 - Invented 1970 by Finish master student (Seppo Linnainmaa)
- Libraries (**autograd**, PyTorch, TF, ...) allow us to write down forward path
- ... and do the rest (**back propagation**) for us
- Still:
 - We should understand what we use all the time
 - Important to understand our tools (including its limitations!)
 - Helps to debug



Automatic vs. numeric / symbolic differentiation

- **Numeric** differentiation
 - Based on finite differences
 - Numerically unstable
 - Small errors accumulate
 - **#evaluations scales with #variables**
- **Symbolic** differentiation
 - General symbolic expression of differentiated function **for arbitrary values**
 - High computational & memory costs
- Automatic differentiation (***autodiff***)
 - Exact derivative of function for a specific value
 - **Requires only a single function evaluation**



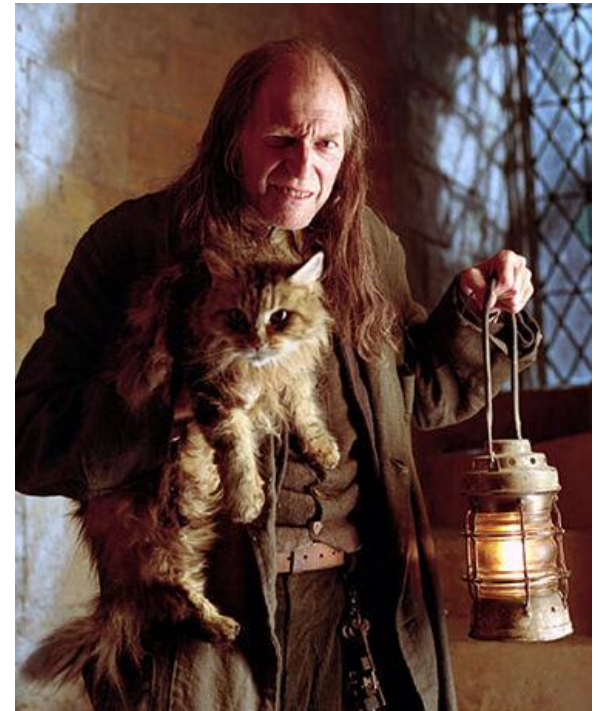
Before we start ...

During this talk

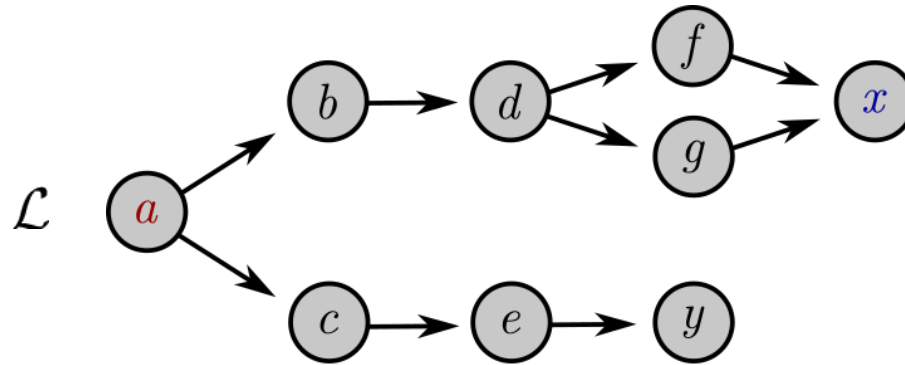
- Autodiff for ML
 - autodiff \leftrightarrow back-propagation
 - Loss function, \mathcal{L} / F : multivariate function (*weights* / x) returning a scalar, y
 - Why? We need the gradient!

$$\vec{x}_{n+1} = \vec{x}_n - \gamma \underbrace{\vec{\nabla} \mathcal{L}(\vec{x}_n)}_{\text{grad } y}$$

- Ask questions!



Back propagation is a fancy name for chain rule



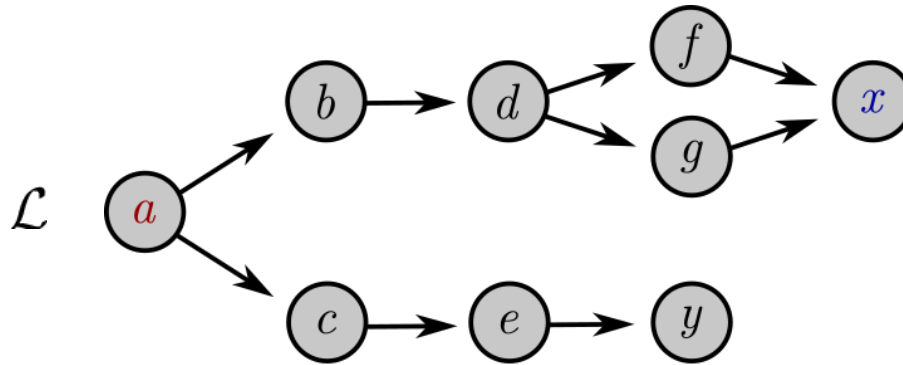
$$\begin{aligned}
 a &= A(b, c) \\
 b &= B(d) \\
 d &= D(f, g) \\
 f &= F(x) \\
 &\vdots
 \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial x} \equiv \frac{\partial a}{\partial x} = \frac{\partial a}{\partial b} \frac{\partial b}{\partial x} + \frac{\partial a}{\partial c} \frac{\partial c}{\partial x}$$

Goal: accumulate gradients for x in $x.\text{grad}$



Chain rule FTW



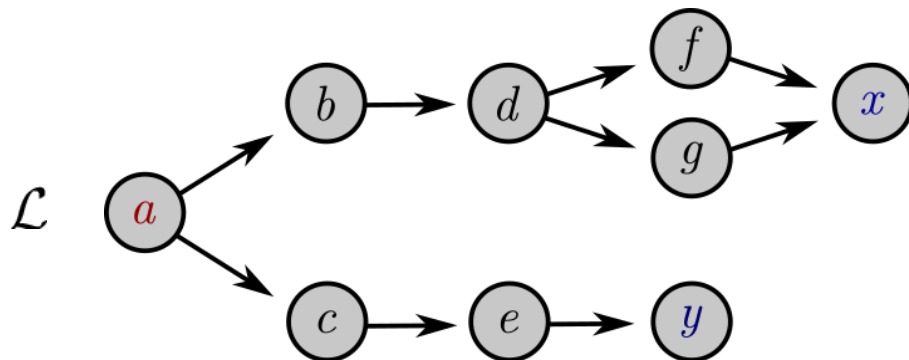
$$\begin{aligned}
 a &= A(b, c) \\
 b &= B(d) \\
 d &= D(f, g) \\
 f &= F(x) \\
 &\vdots
 \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial x} \equiv \frac{\partial a}{\partial x} = \frac{\partial a}{\partial b} \frac{\partial b}{\partial x} + \underbrace{\frac{\partial a}{\partial c} \frac{\partial c}{\partial x}}_{=0}$$

Goal: accumulate gradients for x in $x.\text{grad}$



Chain rule FTW



$$\begin{aligned}
 a &= A(b, c) \\
 b &= B(d) \\
 d &= D(f, g) \\
 f &= F(x) \\
 &\vdots
 \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial x} \equiv \frac{\partial a}{\partial x} = \frac{\partial a}{\partial b} \frac{\partial b}{\partial d} \frac{\partial d}{\partial x} + \frac{\partial a}{\partial c} \frac{\partial c}{\partial e} \frac{\partial e}{\partial x}$$

$$\frac{\partial \mathcal{L}}{\partial y} \equiv \frac{\partial a}{\partial y} = \frac{\partial a}{\partial b} \frac{\partial b}{\partial d} \frac{\partial d}{\partial y} + \frac{\partial a}{\partial c} \frac{\partial c}{\partial e} \frac{\partial e}{\partial y}$$

Goal: accumulate gradients for x in $x.\text{grad}$
 ... and gradients for y in $y.\text{grad}$



Let's implement autodiff!

Let's implement autodiff!

- ... for a rank 0 tensor (scalar)
- ... supporting only addition and multiplication

Strategies:

1. Compiling

- Read & generate source code ahead-of-time (e.g., Python to X)
- Hard to implement
- Libraries: TensorFlow

2. Taping / Recording

- Monitor function execution at runtime
- Easier to implement
- Libraries: PyTorch, JAX



Let's implement autodiff!

```
1  class Scalar:
2      def __init__(self, val):
3          self.val = float(val)
4          self.grad = 0.0
5          self.backward = lambda: None
6
7      def __mul__(self, other):
8          res = Scalar(self.val * other.val)
9
10         def backward():...
11
12         res.backward = backward
13         return res
14
15
16
17
18
19     def __add__(self, other):
20         res = Scalar(self.val + other.val)
21
22         def backward():...
23
24         res.backward = backward
25         return res
26
27
28
29
```



Let's implement autodiff!

```
1 class Scalar:
2     def __init__(self, val):
3         self.val = float(val)
4         self.grad = 0.0
5         self.backward = lambda: None
6
7     def __mul__(self, other):
8         res = Scalar(self.val * other.val)
9
10        def backward():
11            self.grad += other.val * res.grad
12            other.grad += self.val * res.grad
13            self.backward()
14            other.backward()
15
16        res.backward = backward
17        return res
```

$$f(a, b) = ab$$
$$\hookrightarrow df = \frac{\partial f}{\partial a} da + \frac{\partial f}{\partial b} db = b da + a db$$



Let's implement autodiff!

```
1 class Scalar:
2     def __init__(self, val):
3         self.val = float(val)
4         self.grad = 0.0
5         self.backward = lambda: None
6
7     def __mul__(self, other):...
18
19     def __add__(self, other):
20         res = Scalar(self.val + other.val)
21
22         def backward():
23             self.grad += res.grad
24             other.grad += res.grad
25             self.backward()
26             other.backward()
27
28         res.backward = backward
29         return res
```

$$f(a, b) = a + b$$

$$\hookrightarrow df = \frac{\partial f}{\partial a} da + \frac{\partial f}{\partial b} db = da + db$$



Let's implement autodiff!

```
x1, x2 = Scalar(2), Scalar(3)
y = x1 * (x1 + x2)

y.backward()
print("(1):", x1.grad, x2.grad)
```

Output:

```
(1): 0.0, 0.0
```

```
1 class Scalar:
2     def __init__(self, val):
3         self.val = float(val)
4         self.grad = 0.0
5         self.backward = lambda: None
6
7     def __mul__(self, other):
8         res = Scalar(self.val * other.val)
9
10        def backward():
11            self.grad += other.val * res.grad
12            other.grad += self.val * res.grad
13            self.backward()
14            other.backward()
15
16        res.backward = backward
17        return res
18
19    def __add__(self, other):
20        res = Scalar(self.val + other.val)
21
22        def backward():
23            self.grad += res.grad
24            other.grad += res.grad
25            self.backward()
26            other.backward()
27
28        res.backward = backward
29        return res
```



Let's implement autodiff!

```
x1, x2 = Scalar(2), Scalar(3)
y = x1 * (x1 + x2)

y.backward()
print("(1):", x1.grad, x2.grad)

y.grad = 1.0
y.backward()
print("(2):", x1.grad, x2.grad)
```

Output:

```
(1) : 0.0, 0.0
(2) : 7.0, 2.0
```

```
1 class Scalar:
2     def __init__(self, val):
3         self.val = float(val)
4         self.grad = 0.0
5         self.backward = lambda: None
6
7     def __mul__(self, other):
8         res = Scalar(self.val * other.val)
9
10        def backward():
11            self.grad += other.val * res.grad
12            other.grad += self.val * res.grad
13            self.backward()
14            other.backward()
15
16        res.backward = backward
17        return res
18
19    def __add__(self, other):
20        res = Scalar(self.val + other.val)
21
22        def backward():
23            self.grad += res.grad
24            other.grad += res.grad
25            self.backward()
26            other.backward()
27
28        res.backward = backward
29        return res
```



Let's implement autodiff!

```
x1, x2 = Scalar(2), Scalar(3)
y = x1 * (x1 + x2)

y.backward()
print("(1):", x1.grad, x2.grad)

y.grad = 1.0
y.backward()
print("(2):", x1.grad, x2.grad)

y.grad = 1.0
y.backward()
print("(3):", x1.grad, x2.grad)
```

Output:

```
(1): 0.0, 0.0
(2): 7.0, 2.0
(3): 16.0, 6.0
```

```
1  class Scalar:
2      def __init__(self, val):
3          self.val = float(val)
4          self.grad = 0.0
5          self.backward = lambda: None
6
7      def __mul__(self, other):
8          res = Scalar(self.val * other.val)
9
10         def backward():
11             self.grad += other.val * res.grad
12             other.grad += self.val * res.grad
13             self.backward()
14             other.backward()
15
16         res.backward = backward
17         return res
18
19     def __add__(self, other):
20         res = Scalar(self.val + other.val)
21
22         def backward():
23             self.grad += res.grad
24             other.grad += res.grad
25             self.backward()
26             other.backward()
27
28         res.backward = backward
29         return res
```



Let's implement autodiff!

```
x1, x2 = Scalar(2), Scalar(3)
y = x1 * (x1 + x2)

y.backward()
print("(1):", x1.grad, x2.grad)

y.grad = 1.0
y.backward()
print("(2):", x1.grad, x2.grad)

y.grad = 1.0
y.backward()
print("(3):", x1.grad, x2.grad)

x1.grad = 0.0
x2.grad = 0.0
y.grad = 1.0
y.backward()
print("(4):", x1.grad, x2.grad)
```

Output:

```
(1): 0.0, 0.0
(2): 7.0, 2.0
(3): 16.0, 6.0
(4): 11.0, 6.0
```

```
1 class Scalar:
2     def __init__(self, val):
3         self.val = float(val)
4         self.grad = 0.0
5         self.backward = lambda: None
6
7     def __mul__(self, other):
8         res = Scalar(self.val * other.val)
9
10        def backward():
11            self.grad += other.val * res.grad
12            other.grad += self.val * res.grad
13            self.backward()
14            other.backward()
15
16        res.backward = backward
17        return res
18
19    def __add__(self, other):
20        res = Scalar(self.val + other.val)
21
22        def backward():
23            self.grad += res.grad
24            other.grad += res.grad
25            self.backward()
26            other.backward()
27
28        res.backward = backward
29        return res
```



Multivariate autodiff

$$y = \mathcal{L}(\mathbf{x}) = (A \circ B \circ C \circ D)(\mathbf{x}) \equiv A(B(C(D(\mathbf{x}))))$$

$$\mathcal{L} : \begin{array}{c} \blacksquare \\ \blacksquare \\ \blacksquare \\ \blacksquare \\ \blacksquare \end{array} \mapsto \begin{array}{c} \blacksquare \\ y \in \mathbb{R} \end{array}$$

$\mathbf{x} \in \mathbb{R}^n$

$$\mathcal{L}'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right) \in \mathbb{R}^{1 \times n}$$

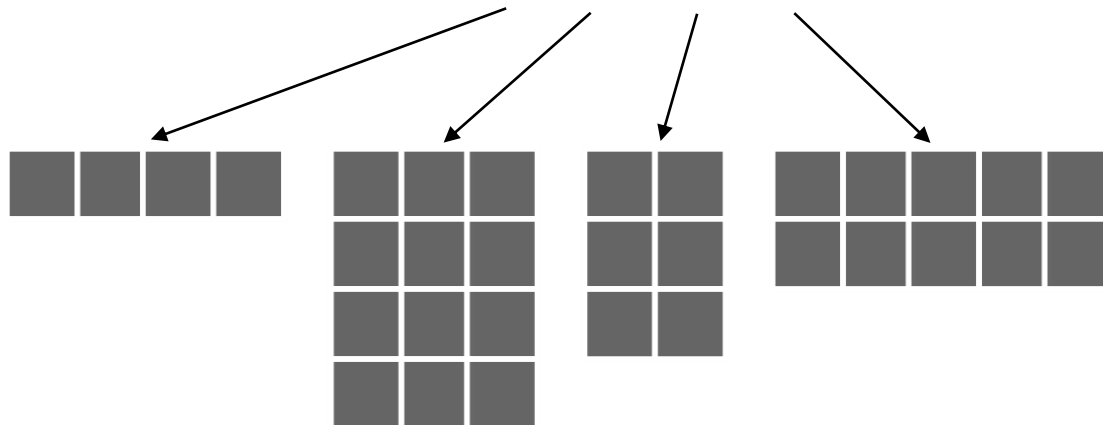


Multivariate autodiff

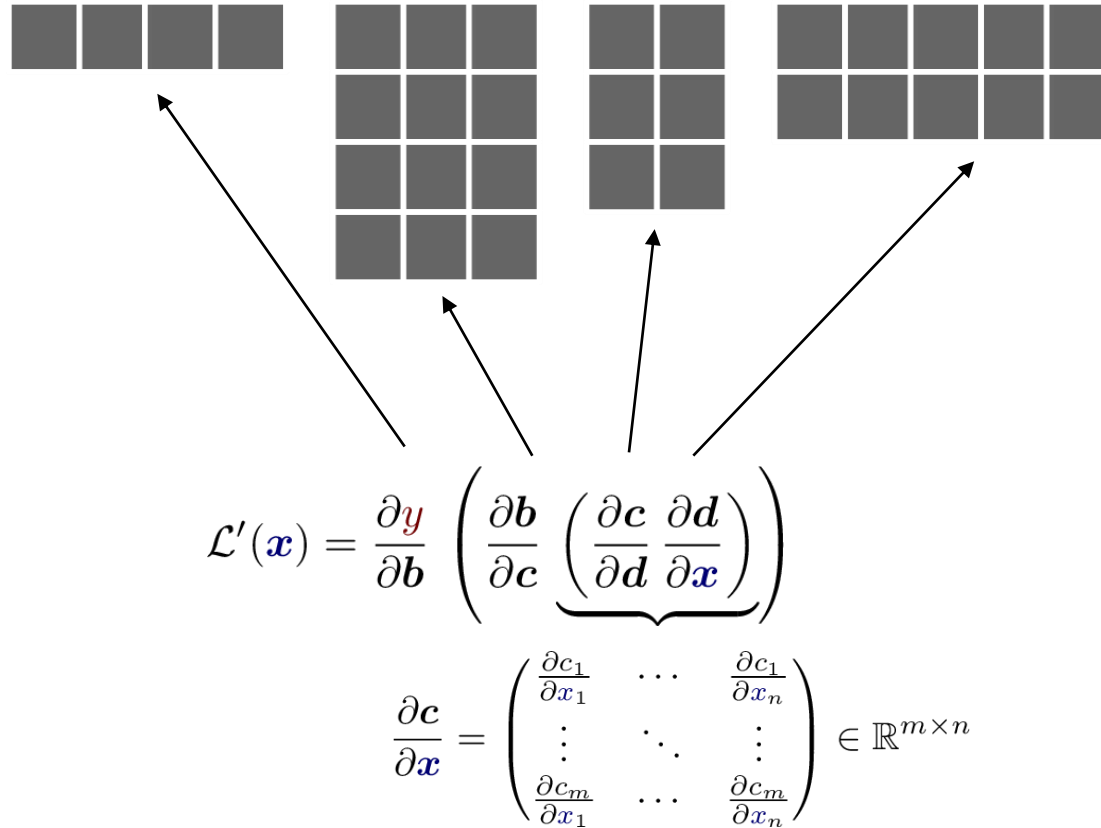
$$y = \mathcal{L}(\mathbf{x}) = (A \circ B \circ C \circ D)(\mathbf{x}) \equiv A(B(C(D(\mathbf{x}))))$$

$$y = A(\mathbf{b}), \quad \mathbf{b} = B(\mathbf{c}), \quad \mathbf{c} = C(\mathbf{d}), \quad \mathbf{d} = D(\mathbf{x})$$

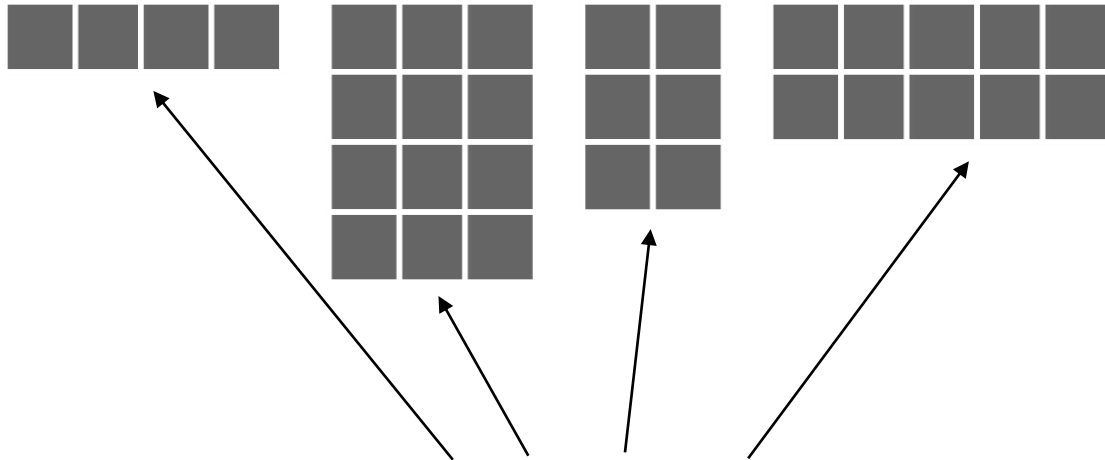
$$\mathcal{L}'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{d}} \frac{\partial \mathbf{d}}{\partial \mathbf{x}}$$



Multivariate autodiff (*forward accumulation*)



Multivariate autodiff (*reverse accumulation*)

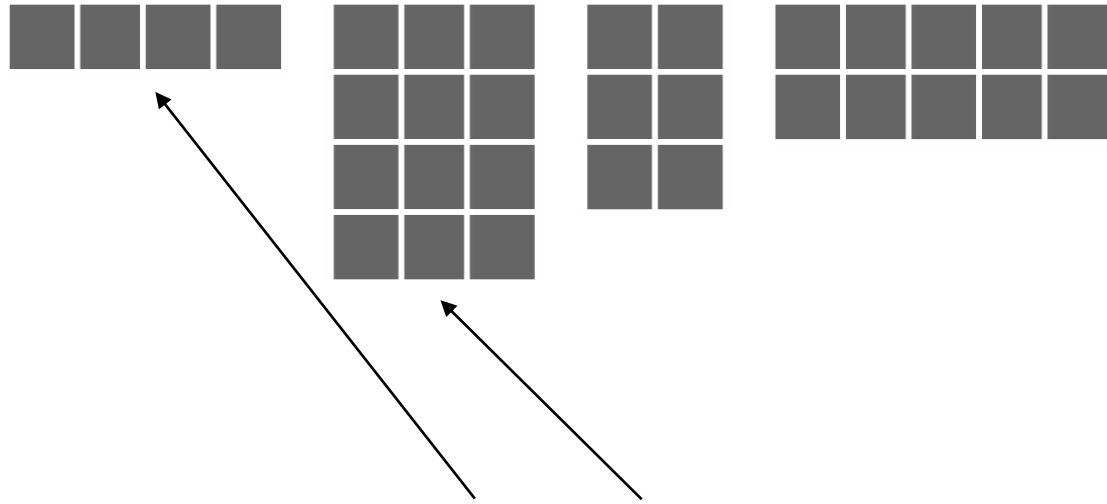


$$\mathcal{L}'(\mathbf{x}) = \left(\underbrace{\begin{pmatrix} \frac{\partial y}{\partial \mathbf{b}} & \frac{\partial \mathbf{b}}{\partial \mathbf{c}} \end{pmatrix}}_{\text{Jacobian of } y \text{ w.r.t. } \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{d}} \right) \frac{\partial \mathbf{d}}{\partial \mathbf{x}}$$

$$\frac{\partial y}{\partial \mathbf{c}} = \left(\frac{\partial y}{\partial c_1}, \dots, \frac{\partial y}{\partial c_{m'}} \right) \in \mathbb{R}^{1 \times m'}$$



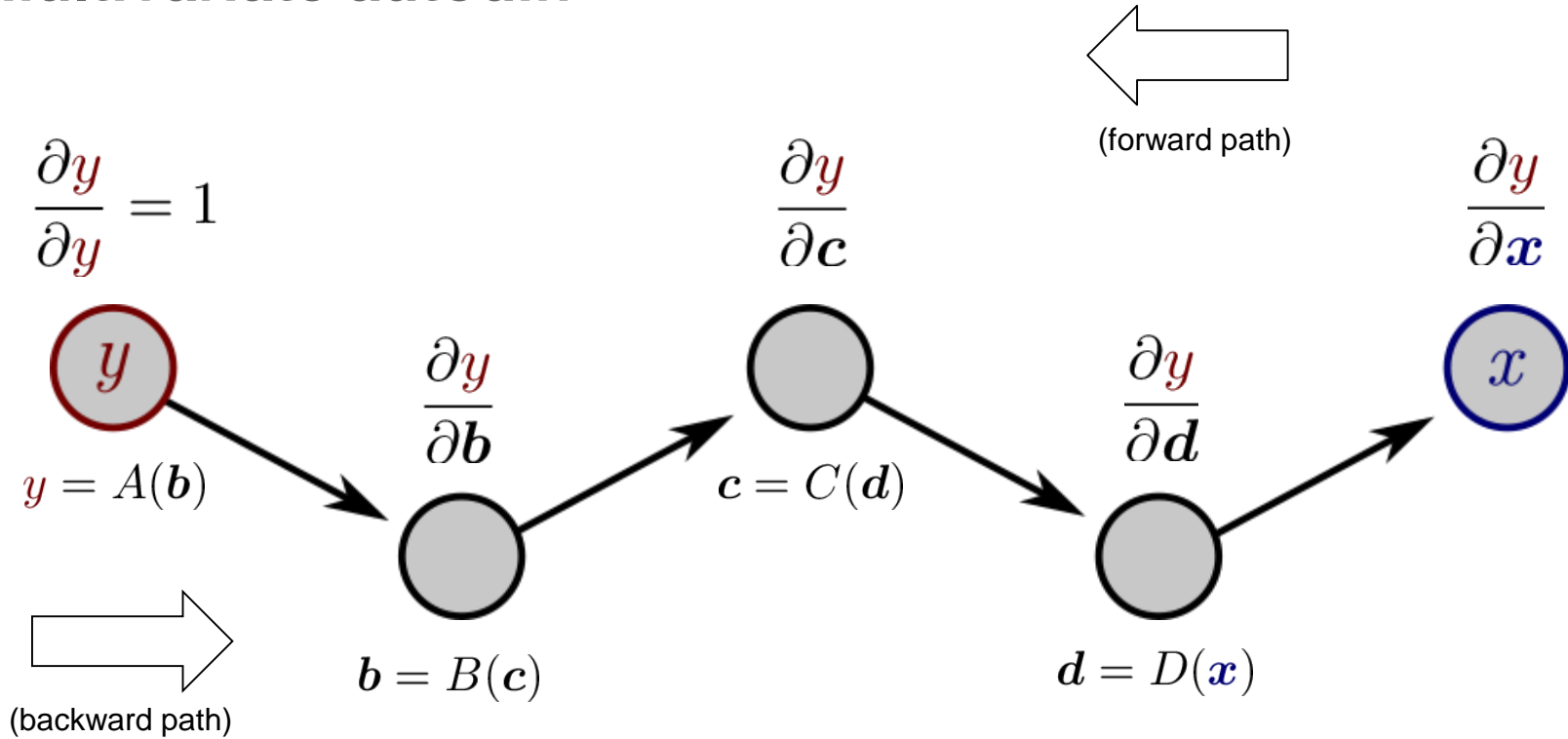
Multivariate autodiff (*reverse accumulation*)



More generally: this is a **vector-Jacobian** product (v_jp)



Multivariate autodiff



$$\mathcal{L}'(\mathbf{x}) = \frac{\partial y}{\partial b} \frac{\partial b}{\partial c} \frac{\partial c}{\partial d} \frac{\partial d}{\partial x}$$

$$\frac{\partial c}{\partial d} \equiv C'(\mathbf{d}) \equiv \left. \frac{\partial C(\tilde{\mathbf{d}})}{\partial \tilde{\mathbf{d}}} \right|_{\tilde{\mathbf{d}}=\mathbf{d}}$$



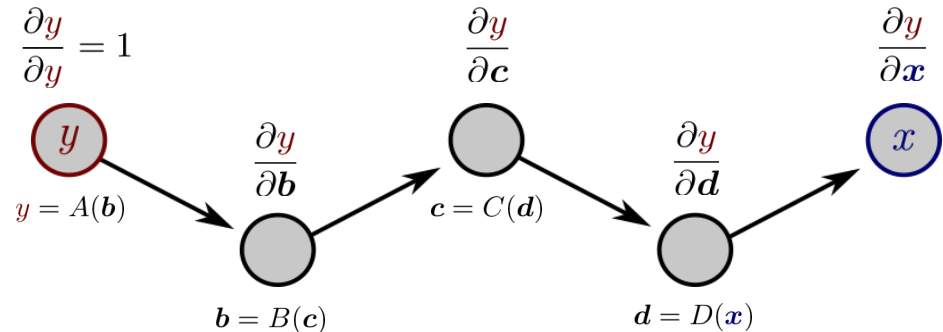
Multivariate autodiff

Forward (at each node):

- Evaluate function to find y
- Save intermediate results

Backward (at each node):

- Evaluate derivative using results from forward path
- Multiply with gradient of previous node



$$\mathcal{L}'(\mathbf{x}) = \frac{\partial y}{\partial b} \frac{\partial b}{\partial c} \frac{\partial c}{\partial d} \frac{\partial d}{\partial x}$$

$$\frac{\partial c}{\partial d} \equiv C'(d) \equiv \left. \frac{\partial C(\tilde{d})}{\partial \tilde{d}} \right|_{\tilde{d}=d}$$



Question

What about loops and branches?

They are **unrolled!**



PyTorch

An engine for computing $\nabla_j p_s$

PyTorch

- Library around a fork of autograd (wrapper of `numpy` function)
- Core data structure: `torch.tensor`
 - Akin to `np.array`
 - Default data type: `float32`
 - Different devices

```
x1 = torch.tensor(2.0, requires_grad=True)
x2 = torch.tensor(3.0, requires_grad=True)
```

```
y = x1 * (x1 + x2)
y.backward()
print("(1):", x1.grad, x2.grad)
```

Output:

```
(1): tensor(7.), tensor(2.)
```



```
>>> x = np.random.rand(4, 2, 3)
>>> sel = [False, False, True]
>>> x[0, :, sel]
array([[0.7545524 , 0.85684194]])
>>> torch.tensor(x)[0, :, sel]
tensor([[0.7546],
        [0.8568]], dtype=torch.float64)
```



PyTorch

```
x1 = torch.tensor(2.0, requires_grad=True)
x2 = torch.tensor(3.0, requires_grad=True)
```

```
y = x1 * (x1 + x2)
y.backward()
print("(1):", x1.grad, x2.grad)
```

```
y.backward()
print("(2):", x1.grad, x2.grad)
```



RuntimeError: Trying to backward through the graph a second time (or directly access saved variables after they have already been freed). Saved intermediate values of the graph are freed when you call `.backward()` or `autograd.grad()`. **Specify `retain_graph=True` if you need to backward through the graph a second time or if you need to access saved variables after calling backward.**



PyTorch

```
x1 = torch.tensor(2.0, requires_grad=True)
x2 = torch.tensor(3.0, requires_grad=True)
```

```
y = x1 * (x1 + x2)
y.backward()
print("(1):", x1.grad, x2.grad)
```

```
y = x1 * (x1 + x2)
y.backward()
print("(2):", x1.grad, x2.grad)
```

Output:

```
(1): tensor(7.), tensor(2.)
(2): tensor(14.), tensor(4.)
```



PyTorch

```
x1 = torch.tensor(2.0, requires_grad=True)
x2 = torch.tensor(3.0, requires_grad=True)
```

```
y = x1 * (x1 + x2)
y.backward()
print("(1):", x1.grad, x2.grad)
```

```
x1.grad.zero_() ←
x2.grad.zero_() ← Code smell!
y = x1 * (x1 + x2)
y.backward()
print("(2):", x1.grad, x2.grad)
```

Output:

```
(1): tensor(7.), tensor(2.)
(2): tensor(7.), tensor(2.)
```



A new PyTorch module

```
1 import torch
2 import torch.nn as nn
3
4
5 class NN(nn.Module):
6     def __init__(self, n: int):
7         super().__init__()
8         self.weights = nn.Parameter(torch.rand(n)) ** 2
9
10    def forward(self, x):
11        return x * self.weights
```

```
14 if __name__ == "__main__":
15     f = NN(5)
16     x = torch.rand(5)
17
18     f.zero_grad()
19     y = f(x)
20     y.backward()
21
22     f.zero_grad()
23     y = f(x)
24     y.backward()
```

RuntimeError: grad can be implicitly created only for scalar outputs



A new PyTorch module

```
1 import torch
2 import torch.nn as nn
3
4
5 class NN(nn.Module):
6     def __init__(self, n: int):
7         super().__init__()
8         self.weights = nn.Parameter(torch.rand(n)) ** 2
9
10    def forward(self, x):
11        return x * self.weights
```

```
14 if __name__ == "__main__":
15     f = NN(5)
16     x = torch.rand(5)
17
18     f.zero_grad()
19     y = torch.sum(f(x))
20     y.backward()
21
22     f.zero_grad()
23     y = torch.sum(f(x))
24     y.backward()
```

RuntimeError: Trying to backward through the graph a second time [...]

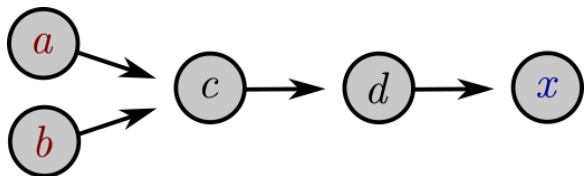


Excursus: backward's retain_graph

```

4 ▶ if __name__ == "__main__":
5     n = 4
6     x = torch.rand(n, requires_grad=True)
7     d = x ** 2
8     c = d * 2.0
9     b = c.mean()
10    a = c.sum()

```



$$a = \sum_{i=1}^n 2x_i^2 \equiv nb \quad \rightarrow \quad \frac{\partial a}{\partial x_j} = 4x_j$$

What is the output of this program?

```

14    b.backward(retain_graph=True)
15    grads_exp = 4.0 * x.detach() / n
16    print("(1)", torch.allclose(x.grad, grads_exp))
17
18    a.backward()
19    print("(2)", torch.allclose(x.grad, n * grads_exp))

```

Output:

- (1) True
- (2) False

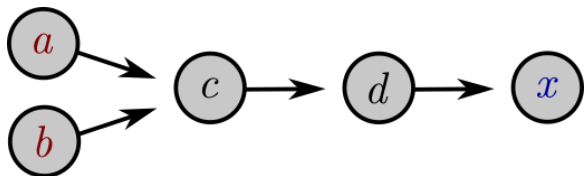


Excursus: backward's retain_graph

```

4 ▶ if __name__ == "__main__":
5     n = 4
6     x = torch.rand(n, requires_grad=True)
7     d = x ** 2
8     c = d * 2.0
9     b = c.mean()
10    a = c.sum()

```



$$a = \sum_{i=1}^n 2x_i^2 \equiv n b \quad \rightarrow \quad \frac{\partial a}{\partial x_j} = 4x_j$$

What is the output of this program?

```

14    b.backward(retain_graph=True)
15    grads = x.grad.clone()
16    grads_exp = 4.0 * x.detach() / n
17    print("(1)", torch.allclose(grads, grads_exp))
18
19    a.backward()
20    print("(2)", torch.allclose(x.grad - grads, n * grads_exp))

```

Output:

- (1) True
- (2) True

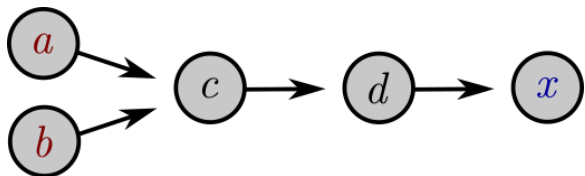


Excursus: backward's retain_graph

```

4 ▶ if __name__ == "__main__":
5     n = 4
6     x = torch.rand(n, requires_grad=True)
7     d = x ** 2
8     c = d * 2.0
9     b = c.mean()
10    a = c.sum()

```



$$a = \sum_{i=1}^n 2x_i^2 \equiv n b \quad \rightarrow \quad \frac{\partial a}{\partial x_j} = 4x_j$$

What is the output of this program?

```

14    b.backward(retain_graph=True)
15    grads_exp = 4.0 * x.detach() / n
16    print("(1)", torch.allclose(x.grad, grads_exp))
17
18    x.grad.zero_()
19    a.backward()
20    print("(2)", torch.allclose(x.grad, n * grads_exp))

```

Output:

- (1) True
- (2) True

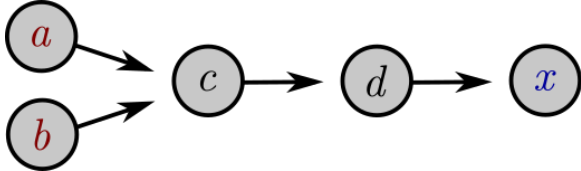


Excursus: backward's retain_graph

```

4 ▶ if __name__ == "__main__":
5     n = 4
6     x = torch.rand(n, requires_grad=True)
7     d = x ** 2
8     c = d * 2.0
9     b = c.mean()
10    a = c.sum()

```



$$a = \sum_{i=1}^n 2x_i^2 \equiv n b \quad \rightarrow \quad \frac{\partial a}{\partial x_j} = 4x_j$$

Still, retain_graph is a code smell...

```

14    b.backward(retain_graph=True)
15    a.backward(retain_graph=True)
16    grads = x.grad.clone()
17
18    x.grad.zero_()
19    loss = a + b
20    loss.backward()
21    assert torch.allclose(x.grad, grads)

```

... and there are alternatives

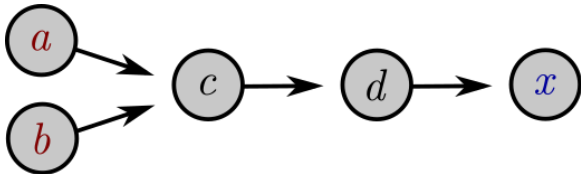


Excursus: backward's retain_graph

```

4 ▶ if __name__ == "__main__":
5     n = 4
6     x = torch.rand(n, requires_grad=True)
7     d = x ** 2
8     c = d * 2.0
9     b = c.mean()
10    a = c.sum()

```



$$a = \sum_{i=1}^n 2x_i^2 \equiv n b \quad \rightarrow \quad \frac{\partial a}{\partial x_j} = 4x_j$$

$$a = \sum_i c_i \quad \rightarrow \quad \left(\frac{\partial a}{\partial c} \right)_i \equiv \frac{\partial a}{\partial c_i} = 1$$

... and there are alternatives

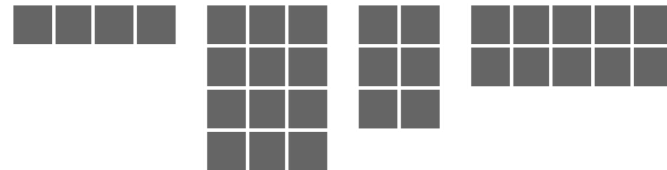
```

14    a.backward(retain_graph=True)
15    grads_exp = x.grad.clone()
16
17    x.grad.zero_()
18    c.backward(gradient=torch.ones_like(x))
19    assert torch.allclose(x.grad, grads_exp)

```

$$\frac{\partial a}{\partial c} = (1, \dots, 1)$$

(In theory: Build Jacobian one row at a time with one-hot vectors)



A new PyTorch module

```
1 import torch
2 import torch.nn as nn
3
4
5 class NN(nn.Module):
6     def __init__(self, n: int):
7         super().__init__()
8         self.weights = nn.Parameter(torch.rand(n)) ** 2
9
10    def forward(self, x):
11        return x * self.weights
```

```
14 if __name__ == "__main__":
15     f = NN(5)
16     x = torch.rand(5)
17
18     f.zero_grad()
19     y = torch.sum(f(x))
20     y.backward()
21
22     f.zero_grad()
23     y = torch.sum(f(x))
24     y.backward()
```

RuntimeError: Trying to backward through the graph a second time [...]



A new PyTorch module

```
1 import torch
2 import torch.nn as nn
3
4
5 class NN(nn.Module):
6     def __init__(self, n: int):
7         super().__init__()
8         self.weights = nn.Parameter(torch.rand(n))
9
10    def forward(self, x):
11        return x * (self.weights ** 2)
```

```
14 if __name__ == "__main__":
15     f = NN(5)
16     x = torch.rand(5)
17
18     f.zero_grad()
19     y = torch.sum(f(x))
20     y.backward()
21
22     f.zero_grad()
23     y = torch.sum(f(x))
24     y.backward()
```



A new PyTorch module

```
1  import torch
2  import torch.nn as nn
3
4
5  class NN(nn.Module):
6      def __init__(self, n: int):
7          super().__init__()
8          self.weights = nn.Parameter(torch.rand(n) ** 2)
9
10     def forward(self, x):
11         return x * self.weights
```

```
14  if __name__ == "__main__":
15     f = NN(5)
16     x = torch.rand(5)
17
18     f.zero_grad()
19     y = torch.sum(f(x))
20     y.backward()
21
22     f.zero_grad()
23     y = torch.sum(f(x))
24     y.backward()
```

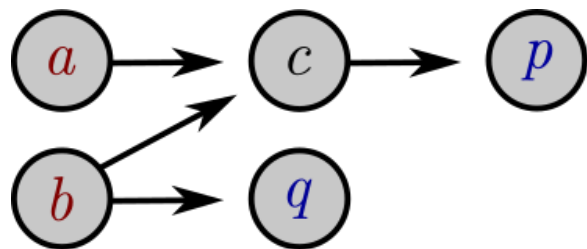


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



```

6
7
8
9
10
11
12
13
14
15
16
17
18
p_val, q_val = 2.0, 3.0
p = torch.tensor(p_val, requires_grad=True)
q = torch.tensor(q_val, requires_grad=True)

c = 2 * p
b = c * q
a = c ** 2

a.backward()
b.backward()

assert np.isclose(p.grad.item(), 8 * p_val)
assert np.isclose(q.grad.item(), 2 * p_val)

```

RuntimeError: Trying to backward through the graph a second time [...]

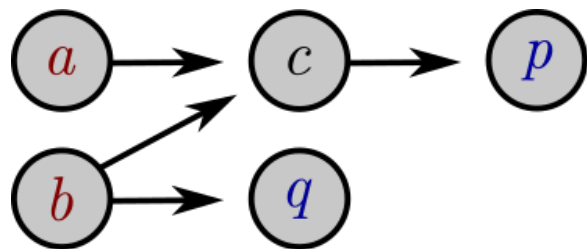


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

```

p_val, q_val = 2.0, 3.0
p = torch.tensor(p_val, requires_grad=True)
q = torch.tensor(q_val, requires_grad=True)

c = 2 * p
c_cpy = c.detach()
b = c_cpy * q
a = c ** 2

a.backward()
b.backward()

assert np.isclose(p.grad.item(), 8 * p_val)
assert np.isclose(q.grad.item(), 2 * p_val)

print("(1)", c_cpy)
c.mul_(2)
print("(2)", c_cpy)
  
```

Output:

```

(1) tensor(4.)
(2) tensor(8.)
  
```

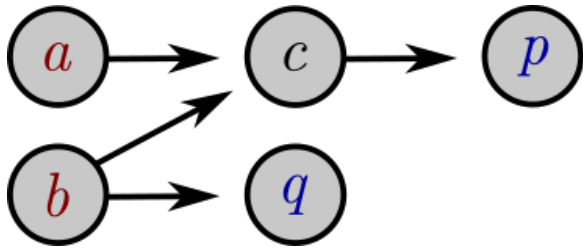


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



```

6   p_val, q_val = 2.0, 3.0
7   p = torch.tensor(p_val, requires_grad=True)
8   q = torch.tensor(q_val, requires_grad=True)
9
10  c = 2 * p
11  c_cpy = c.detach().clone()
12  b = c_cpy * q
13  a = c ** 2
14
15  a.backward()
16  b.backward()
17
18  assert np.isclose(p.grad.item(), 8 * p_val)
19  assert np.isclose(q.grad.item(), 2 * p_val)
20
21  print("(1)", c_cpy)
22  c.mul_(2)
23  print("(2)", c_cpy)

```

Output:

```

(1) tensor(4.)
(2) tensor(4.)

```

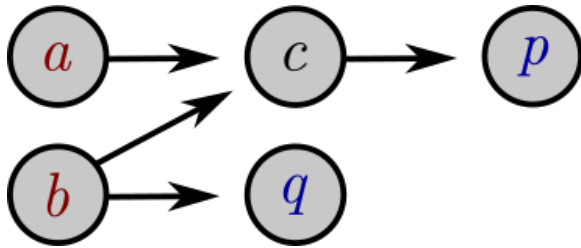


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

```

p_val, q_val = 2.0, 3.0
p = torch.tensor(p_val, requires_grad=True)
q = torch.tensor(q_val, requires_grad=True)

c = 2 * p
c.mul_(1.0)
c_cpy = c.detach()
b = c_cpy * q
a = c ** 2

a.backward()
b.backward()

assert np.isclose(p.grad.item(), 8 * p_val)
assert np.isclose(q.grad.item(), 2 * p_val)
  
```

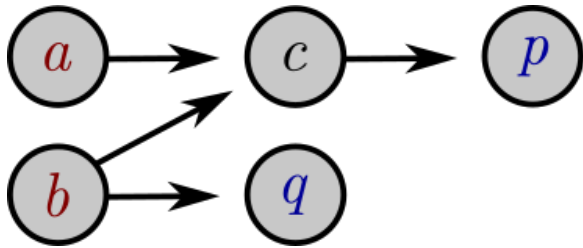


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



```

6     p_val, q_val = 2.0, 3.0
7     p = torch.tensor(p_val, requires_grad=True)
8     q = torch.tensor(q_val, requires_grad=True)
9
10    c = 2 * p
11    c_cpy = c.detach()
12    b = c_cpy * q
13    c.mul_(1.0)
14    a = c ** 2
15
16    a.backward()
17    b.backward()
18
19    assert np.isclose(p.grad.item(), 8 * p_val)
20    assert np.isclose(q.grad.item(), 2 * p_val)
  
```

RuntimeError: one of the variables needed for gradient computation has been modified by an inplace operation [...]

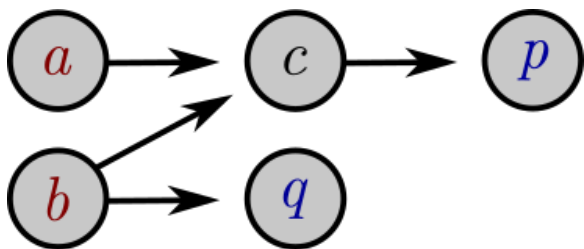


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



```

6   p_val, q_val = 2.0, 3.0
7   p = torch.tensor(p_val, requires_grad=True)
8   q = torch.tensor(q_val, requires_grad=True)
9
10  c = 2 * p
11  c_cpy = c.detach()
12  b = c_cpy * q
13  c *= 1.0
14  a = c ** 2
15
16  a.backward()
17  b.backward()
18
19  assert np.isclose(p.grad.item(), 8 * p_val)
20  assert np.isclose(q.grad.item(), 2 * p_val)

```

RuntimeError: one of the variables needed for gradient computation has been modified by an inplace operation [...]

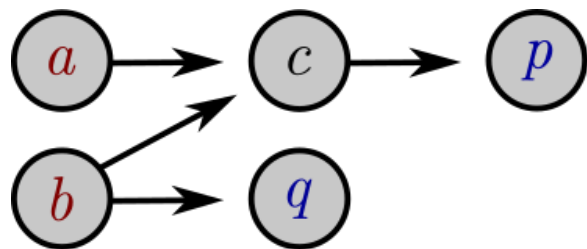


detach() & clone()

$$c = 2p \quad \frac{\partial a}{\partial p} = 8p$$

$$b = qc \quad \frac{\partial b}{\partial q} = 2p$$

$$a = c^2$$



```

6   p_val, q_val = 2.0, 3.0
7   p = torch.tensor(p_val, requires_grad=True)
8   q = torch.tensor(q_val, requires_grad=True)
9
10  c = 2 * p
11  c_cpy = c.detach()
12  b = c_cpy * q
13  c = c * 1.0
14  a = c ** 2
15
16  a.backward()
17  b.backward()
18
19  assert np.isclose(p.grad.item(), 8 * p_val)
20  assert np.isclose(q.grad.item(), 2 * p_val)

```

Rule of thumb: avoid in-place operations and be explicit when doing overwrites.



!!! WARNING: indexing is also in-place !!!

- Avoid indexing whenever possible:
 - Select and indexing yield views
 - There are in-place correctness checks but they are not perfect
- That is:
 - `... = x[i]` # **code smell**
 - `x[i] = ...` # **code smell**
 - Be explicit and use `torch.cat` if necessary

```

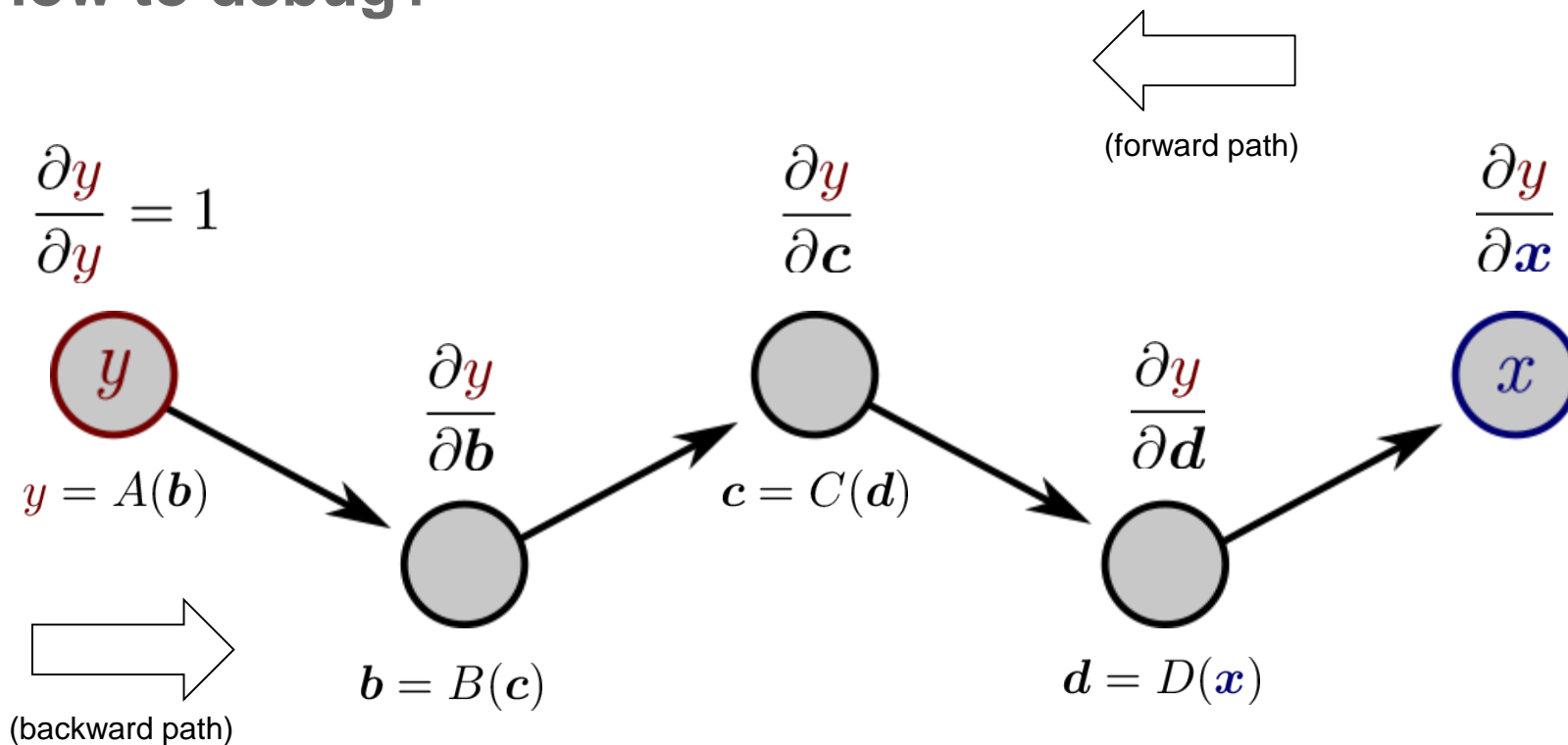
5 class MyFancyModel(nn.Module):
6     def __init__(self):
7         super().__init__()
8
9     def forward(self, x):
10        y = torch.empty(3)
11        y[:2] = x
12        y[2] = x[0] * y[0]
13        return y[0]
14
15
16 if __name__ == "__main__":
17     x = torch.rand(2, requires_grad=True)
18
19     model = MyFancyModel()
20     y = model(x)
21     y.backward()

```

RuntimeError: one of the variables needed for gradient computation has been modified by an inplace operation [...]



How to debug?



Forward path: `print`
 Backward path: **hooks**



Inspecting the backward path

$$\frac{\partial y}{\partial b_i} = \frac{\partial}{\partial b_i} \left(\sum_i b_i \right) = 1$$

$$\frac{\partial y}{\partial \mathbf{b}} = (1, \dots, 1)$$

$$\frac{\partial b_i}{\partial a_j} = \frac{\partial}{\partial a_j} (2a_i) = 2 \delta_{ij}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{a}} = \begin{pmatrix} 2 & & \\ & \ddots & \\ & & 2 \end{pmatrix}$$

$$\frac{\partial a_i}{\partial x_j} = \frac{\partial}{\partial x_j} (x_i^2) = 2x_i \delta_{ij}$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} = \begin{pmatrix} 2x_1 & & \\ & \ddots & \\ & & 2x_n \end{pmatrix}$$

```
grad of y: tensor(1.)
grad of b: tensor([1., 1., 1., 1.])
grad of a: tensor([2., 2., 2., 2.])
grad of x: tensor([2.28, 3.58, 2.06, 3.20])
Exp. grad of x: tensor([2.28, 3.58, 2.06, 3.20])
```

```
17 ▶ if __name__ == "__main__":
18     n = 4
19     x = torch.rand(n, requires_grad=True)
20     a = x ** 2
21     b = a * 2.0 ←
22     y = b.sum()
23
24     a.retain_grad()
25     b.retain_grad()
26     y.retain_grad()
27
28     y.backward()
29     print("grad of y:", y.grad)
30     print("grad of b:", b.grad)
31     print("grad of a:", a.grad)
32     print("grad of x:", x.grad)
33     print("Exp. grad of x:", x * 4)
```

$$a_i = x_i$$

$$b_i = 2a_i$$

$$y = \sum_i b_i$$

$$\frac{\partial y}{\partial \mathbf{x}} = \frac{\partial y}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}} = (1, \dots, 1) \begin{pmatrix} 2 & & \\ & \ddots & \\ & & 2 \end{pmatrix} \begin{pmatrix} 2x_1 & & \\ & \ddots & \\ & & 2x_n \end{pmatrix} = (4x_1, \dots, 4x_n)$$



Inspecting the backward path

```

17 ▶ if __name__ == "__main__":
18     n = 4
19     x = torch.rand(n, requires_grad=True)
20     a = x ** 2
21     b = a * 2.0
22     y = b.sum()
23
24     a.retain_grad()
25     b.retain_grad()
26     y.retain_grad()
27
28     y.backward()
29     print("grad of y:", y.grad)
30     print("grad of b:", b.grad)
31     print("grad of a:", a.grad)
32     print("grad of x:", x.grad)
33     print("Exp. grad of x:", x * 4)

```



```

5 def hook(msg, grad):
6     print(msg, grad)
7
8
9 def register_hooks(**kwargs):
10     for k in kwargs:
11         v = kwargs[k]
12
13         msg = f"grad of {k}:"
14         v.register_hook(functools.partial(hook, msg))
15
16
17 ▶ if __name__ == "__main__":
18     n = 4
19     x = torch.rand(n, requires_grad=True)
20     a = x ** 2
21     b = a * 2.0
22     y = b.sum()
23
24     register_hooks(y=y, b=b, a=a, x=x)
25
26     y.backward()
27     print("Exp. grad of x:", x * 4)

```

- Hooks are called **each time** the node is called **during back propagation**
- Hooks are allow to change gradients



Inspecting the backward path

$$\frac{\partial y}{\partial b_i} = \frac{\partial}{\partial b_i} \left(\sum_i b_i \right) = 1$$

$$\frac{\partial y}{\partial \mathbf{b}} = (1, \dots, 1)$$

$$\frac{\partial b_i}{\partial a_j} = \frac{\partial}{\partial a_j} (2a_i) = 2 \delta_{ij}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{a}} = \begin{pmatrix} 2 & & \\ & \ddots & \\ & & 2 \end{pmatrix}$$

$$\frac{\partial a_i}{\partial x_j} = \frac{\partial}{\partial x_j} (x_i^2) = 2x_i \delta_{ij}$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} = \begin{pmatrix} 2x_1 & & \\ & \ddots & \\ & & 2x_n \end{pmatrix}$$

```

5 class MyFancyModel(nn.Module):
6     def __init__(self):
7         super().__init__()
8
9     def forward(self, x):
10        return x ** 2
11
12
13 def make_hook(msg):
14     return lambda *args: print(f"{msg}:", *args)
15
16
17 if __name__ == "__main__":
18     f = MyFancyModel()
19     f.register_forward_hook(make_hook("fwd"))
20     f.register_full_backward_hook(make_hook("bwd"))
21
22     x = torch.tensor([1.0, 2.0, 3.0], requires_grad=True)
23     a = f(x)
24     b = a * 2.0
25     y = b.sum()
26
27     y.backward()

```

Output:

```
fwd: MyFancyModel() (tensor([1., 2., 3.]),) tensor([1., 4., 9.],)
```

```
bwd: MyFancyModel() (tensor([ 4.,  8., 12.]),) (tensor([2., 2., 2.]),)
```



Custom autograd functions

$$a = 2x$$

$$b = 3x$$

$$f = a \sin b$$

$$y = f^2$$

```

6 class MyFancyFun(torch.autograd.Function):
7     @staticmethod
8     def forward(ctx, a, b):
9         ctx.save_for_backward(a, b)
10        return a * torch.sin(b)
11
12    @staticmethod
13    def backward(ctx, grad):
14        a, b = ctx.saved_tensors
15        return torch.sin(b) * grad, a * torch.cos(b) * grad

```

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial f} \begin{pmatrix} \frac{\partial f}{\partial a} & \frac{\partial f}{\partial b} \end{pmatrix} \begin{pmatrix} \frac{\partial a}{\partial x} \\ \frac{\partial b}{\partial x} \end{pmatrix}$$

$$= 2f (\sin b, a \cos b) \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

```

18 if __name__ == "__main__":
19     fun = MyFancyFun().apply
20
21     x = torch.tensor([0.5], requires_grad=True)
22     a = 2.0 * x
23     b = 3.0 * x
24     f = fun(a, b)
25     y = f ** 2
26     y.backward()
27
28     x_val = x.detach().clone()
29     s = torch.sin(3.0 * x_val)
30     c = torch.cos(3.0 * x_val)
31     grad_exp = 4 * x_val * s * (2.0 * s + 6.0 * x_val * c)
32     assert torch.isclose(x.grad, grad_exp)

```



Custom autograd functions

$$a = 2x$$

$$b = 3x$$

$$f = a \sin b$$

$$y = f^2$$

```

6 class MyFancyFun(torch.autograd.Function):
7     @staticmethod
8     def forward(ctx, a, b):
9         res = a * torch.sin(b)
10        ctx.save_for_backward(a, b, res)
11        return res
12
13    @staticmethod
14    def backward(ctx, grad):
15        a, b, res = ctx.saved_tensors
16        sin_b = res / a
17        sign = torch.where(torch.abs(b) < math.pi / 2.0, +1.0, -1.0)
18        return sin_b * grad, sign * a * torch.sqrt(1.0 - sin_b ** 2) * grad

```

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial f} \begin{pmatrix} \frac{\partial f}{\partial a} & \frac{\partial f}{\partial b} \end{pmatrix} \begin{pmatrix} \frac{\partial a}{\partial x} \\ \frac{\partial b}{\partial x} \end{pmatrix}$$

$$= 2f (\sin b, a \cos b) \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

```

18 if __name__ == "__main__":
19     fun = MyFancyFun().apply
20
21     x = torch.tensor([0.5], requires_grad=True)
22     a = 2.0 * x
23     b = 3.0 * x
24     f = fun(a, b)
25     y = f ** 2
26     y.backward()
27
28     x_val = x.detach().clone()
29     s = torch.sin(3.0 * x_val)
30     c = torch.cos(3.0 * x_val)
31     grad_exp = 4 * x_val * s * (2.0 * s + 6.0 * x_val * c)
32     assert torch.isclose(x.grad, grad_exp)

```



Thank you for your attention!

