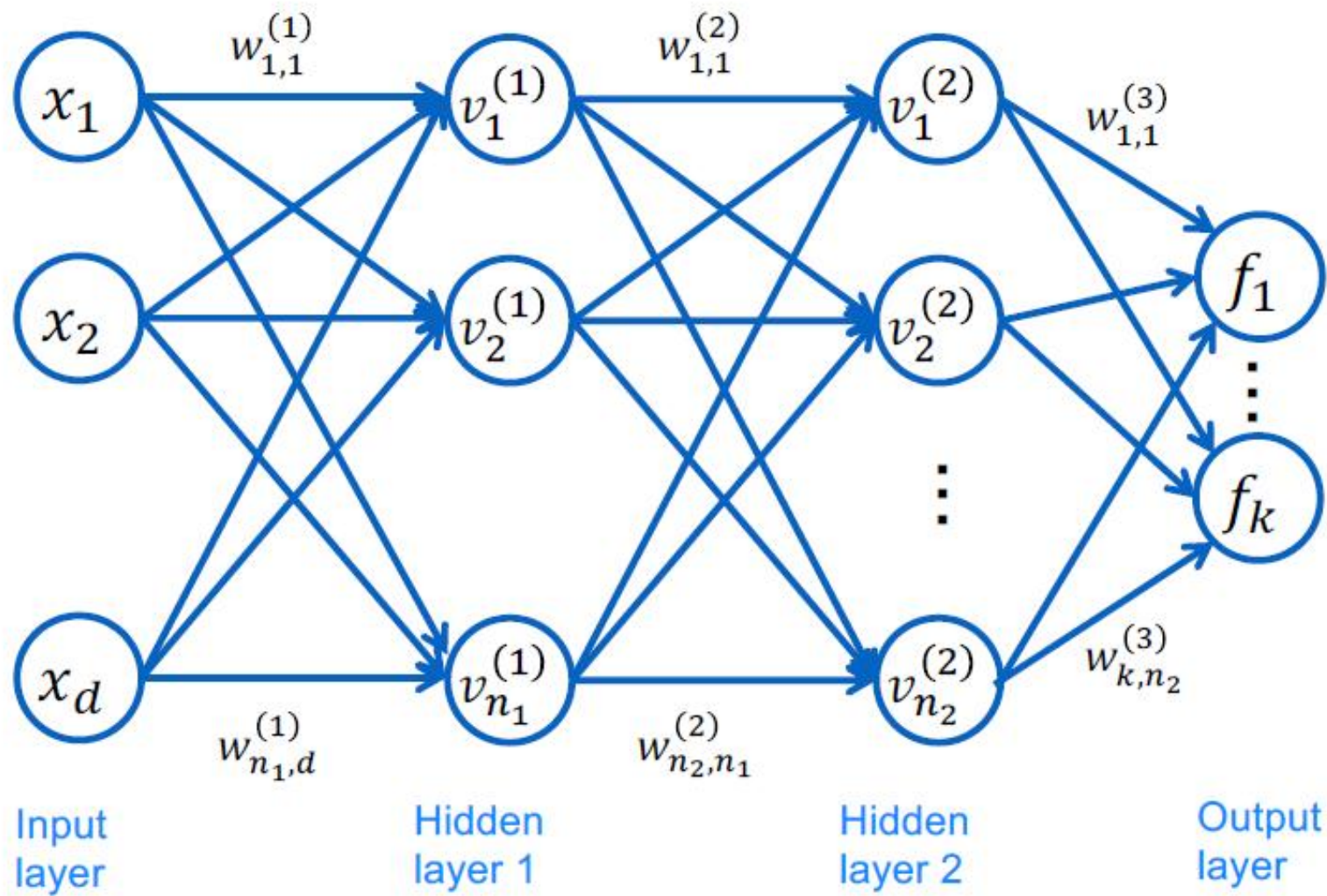


Introduction to Deep Learning

Calculus behind Backpropagation

Lecturer: Dongyu Yao



Forward propagation (short notation, with bias nodes)

- For input layer: $\mathbf{v}^{(0)} = [\mathbf{x}; 1]$

- For each hidden layer $l = 1:L - 1$

$$\mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{v}^{(l-1)} \quad \text{and} \quad \mathbf{v}^{(l)} = [\varphi(\mathbf{z}^{(l)}); 1]$$

- For output layer:

$$\mathbf{f} = \mathbf{W}^{(L)} \mathbf{v}^{(L-1)}$$

- Predict $\mathbf{y} = \mathbf{f}$ for regression, or $y = \arg \max_j f_j$ for classification

Computing the gradient

In order to apply SGD, need to compute

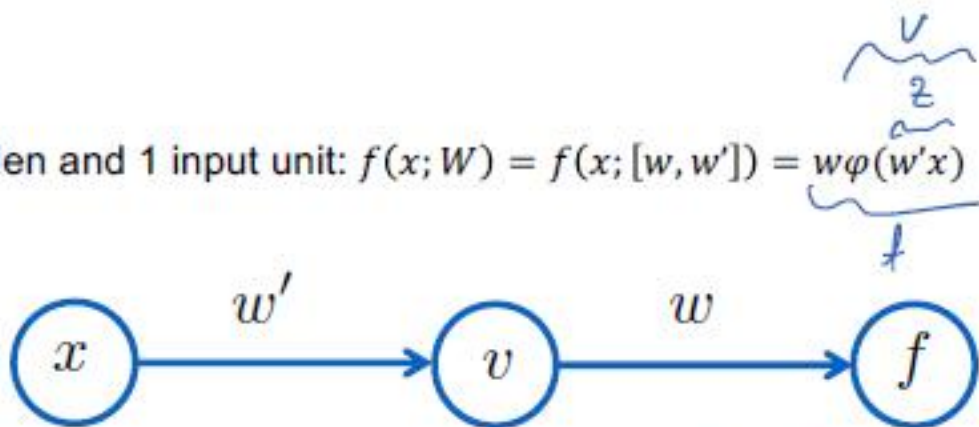
$$\nabla_{\mathbf{W}} \ell(\mathbf{W}; \mathbf{x}, y)$$

I.e., for each weight between any two connected units i and j on layer l need

$$\frac{\partial}{\partial w_{ij}^{(l)}} \ell(\mathbf{W}; \mathbf{x}, y)$$

Simple example

ANN with 1 output, 1 hidden and 1 input unit: $f(x; W) = f(x; [w, w']) = w\varphi(w'x)$



Loss: $l(w; x, y)$, e.g. $l(w; x, y) = (y - f(x; w))^2 =: l_y(f)$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial f} \frac{\partial f}{\partial w} = \underbrace{l'_y(f)}_{\delta} \cdot \underbrace{v}_{\delta'} \quad \leftarrow \text{Computed during forward pass}$$

$$\frac{\partial l}{\partial w'} = \underbrace{\frac{\partial l}{\partial f}}_{\delta} \cdot \underbrace{\frac{\partial f}{\partial v}}_{w} \cdot \underbrace{\frac{\partial v}{\partial z}}_{\varphi'(z)} \cdot \underbrace{\frac{\partial z}{\partial w'}}_x$$

Want:

$$\nabla_w l(w; x, y)$$

$$= \left[\frac{\partial l}{\partial w}, \frac{\partial l}{\partial w'} \right]^T$$

$$z = w' \cdot x$$

$$v = \varphi(z)$$

$$f = w \cdot v$$

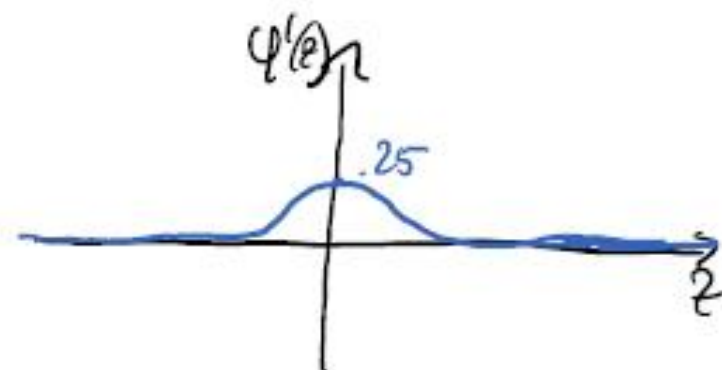
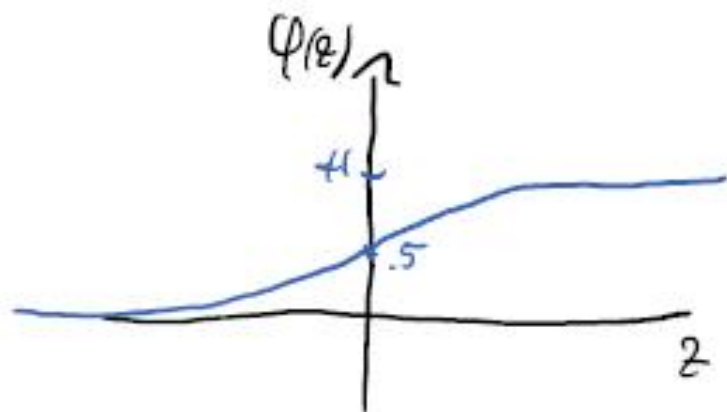
$$l'_y(f) = 2(f - y) =: \delta$$

Derivatives of activation functions

Sigmoid: $\varphi(z) = \frac{1}{1 + \exp(-z)}$

$$\begin{aligned}\varphi'(z) &= \frac{-1}{(1 + e^{-z})^2} \cdot e^{-z} \cdot (-1) = \frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \frac{e^{-z}}{(1 + e^{-z})} \cdot \frac{1}{(1 + e^{-z})} = \underbrace{(1 - \varphi(z))}_{(1-v)} \underbrace{\varphi(z)}_v\end{aligned}$$

-
- + $\varphi'(z)$ is easy to compute given $\varphi(z) = v$
 - + $\varphi'(z) \neq 0 \forall z$, defined everywhere
 - $\varphi'(z) \approx 0 \forall z$ (except $z \approx 0$)

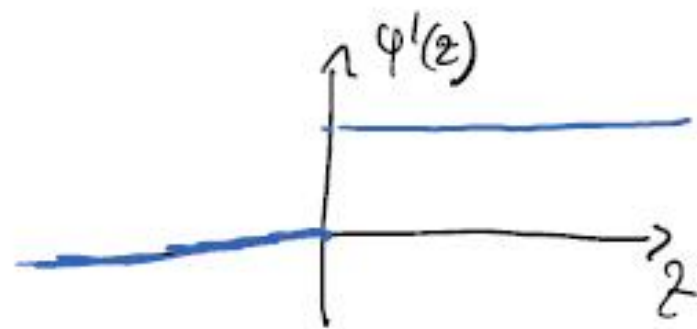
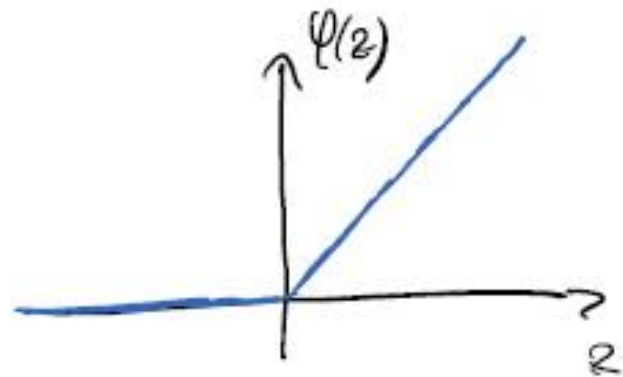


Derivatives of activation functions

Rectified Linear Unit (ReLU): $\varphi(z) = \max(z, 0)$

$$\varphi'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \end{cases}$$

-
- not differentiable at $z=0$
doesn't matter much in practice;
eg. simply define $\varphi'(0) = 0$
 - + even easier to compute than sigmoid
 - + $\neq 0 \quad \forall z > 0$



Recall: Jacobians & multivariate chain rule

For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ its Jacobian at \mathbf{x} is given by

$$\frac{d}{d\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} =: J_f(\mathbf{x}) =: \frac{\partial f}{\partial \mathbf{x}}$$

Multivariate Taylor series: Given differentiable f , it holds that

$$f(\mathbf{x}) = f(\mathbf{x}_0) + J_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|)$$

Multivariate chain rule: Given differentiable $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$, it holds that

$$\frac{d}{d\mathbf{x}} g(f(\mathbf{x})) = J_{g \circ f}(\mathbf{x}) = J_g(f(\mathbf{x}))J_f(\mathbf{x}) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}}$$

Examples

$$z = z(x) = Wx \quad \Rightarrow \quad \overset{\frac{\partial z}{\partial x}}{J_z(x)} = W$$

$$\varphi(x) = [\varphi(x_1), \dots, \varphi(x_n)]^T \quad \Rightarrow \quad J_\varphi(x) =$$

$$\begin{matrix} \text{diag}(\varphi'(x)) & \swarrow & \varphi(x) = [\varphi(x_1), \dots, \varphi(x_n)] \\ \left(\begin{array}{ccc} \varphi'(x_1) & 0 & \dots & 0 \\ & \varphi'(x_2) & & 0 \\ & & \dots & \\ 0 & & & \varphi'(x_n) \end{array} \right) \end{matrix}$$

$$f(x) = \varphi\left(\frac{Wx}{z}\right) \quad \Rightarrow \quad J_f(x) = J_\varphi(z) J_z(x) = \text{diag}(\varphi'(z)) \cdot W$$

Backpropagation

Consider neural network $f(x; W)$ with weights $W = [W^{(1)}, W^{(2)}, \dots, W^{(L)}]$

Want to compute gradient of loss $\ell(W; x, y)$ w.r.t., weights $W^{(i)}$

$$\begin{aligned}(\nabla_{W^{(L)}} \ell)^T &= \frac{\partial \ell}{\partial W^{(L)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial W^{(L)}} \\(\nabla_{W^{(L-1)}} \ell)^T &= \frac{\partial \ell}{\partial W^{(L-1)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial W^{(L-1)}} \\(\nabla_{W^{(L-2)}} \ell)^T &= \frac{\partial \ell}{\partial W^{(L-2)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \frac{\partial z^{(L-2)}}{\partial W^{(L-2)}} \\&\vdots \\(\nabla_{W^{(i)}} \ell)^T &= \frac{\partial \ell}{\partial W^{(i)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \dots \frac{\partial z^{(i+1)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial W^{(i)}}\end{aligned}$$

Key insight: Can reuse computations from **forward propagation** and from **layer $i+1$** to **compute $W^{(i)}$**

Backpropagation

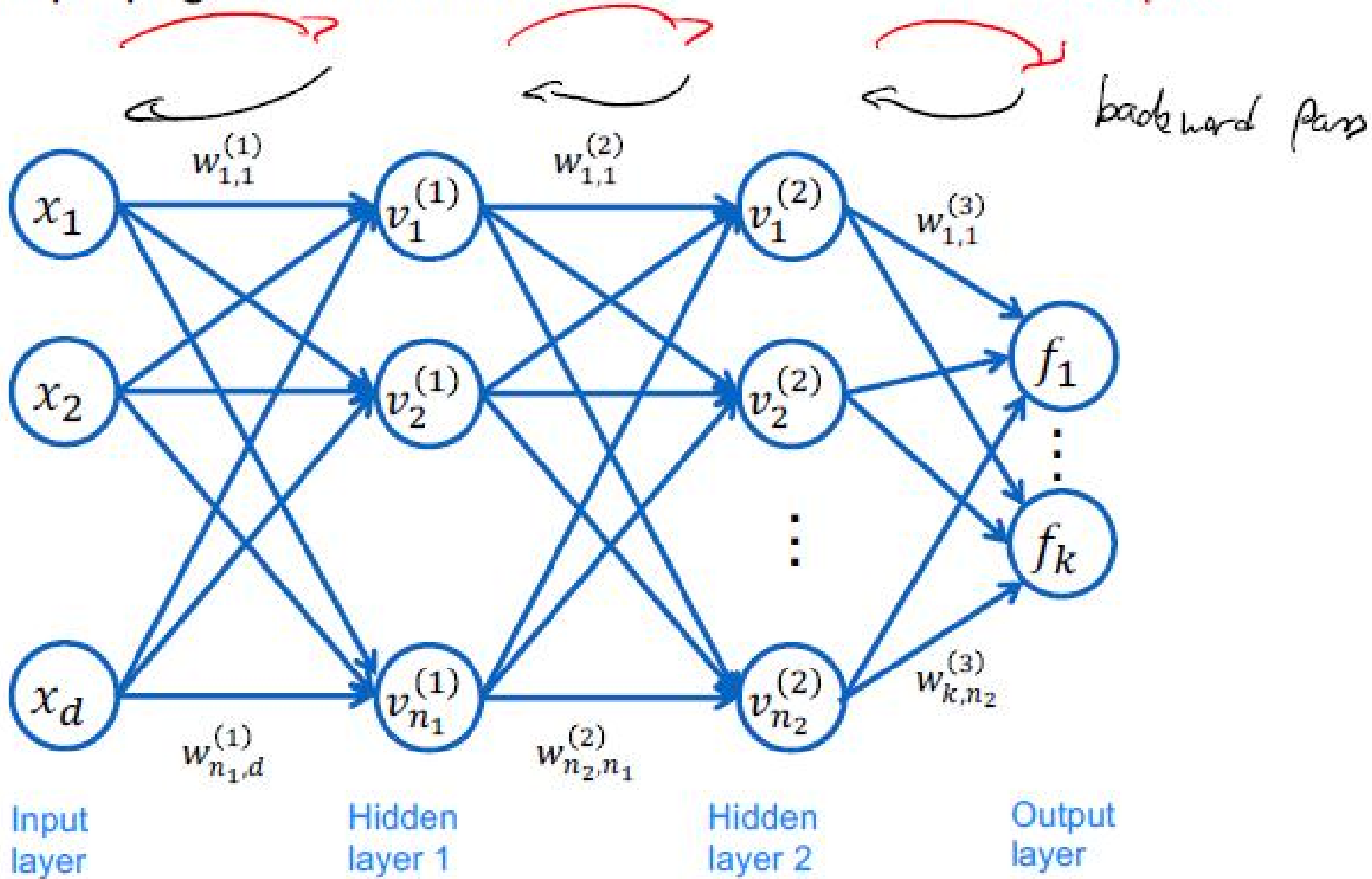
1. For each unit j on the output layer

- Compute error signal $\delta_j^{(L)} = \frac{\partial \ell}{\partial f_j}$
- For each unit i on layer L , $\frac{\partial}{\partial w_{ji}^{(L)}} \ell = \delta_j^{(L)} v_i^{(L-1)}$

2. For each unit j on hidden layer $l = L - 1: -1: 1$

- Compute error signal $\delta_j^{(l)} = \varphi'(z_j^{(l)}) \sum_{i=1}^{n_{l+1}} w_{i,j}^{(l+1)} \delta_i^{(l+1)}$
- For each unit i on layer l : $\frac{\partial}{\partial w_{ji}^{(l)}} \ell = \delta_j^{(l)} v_i^{(l-1)}$

Backpropagation illustration



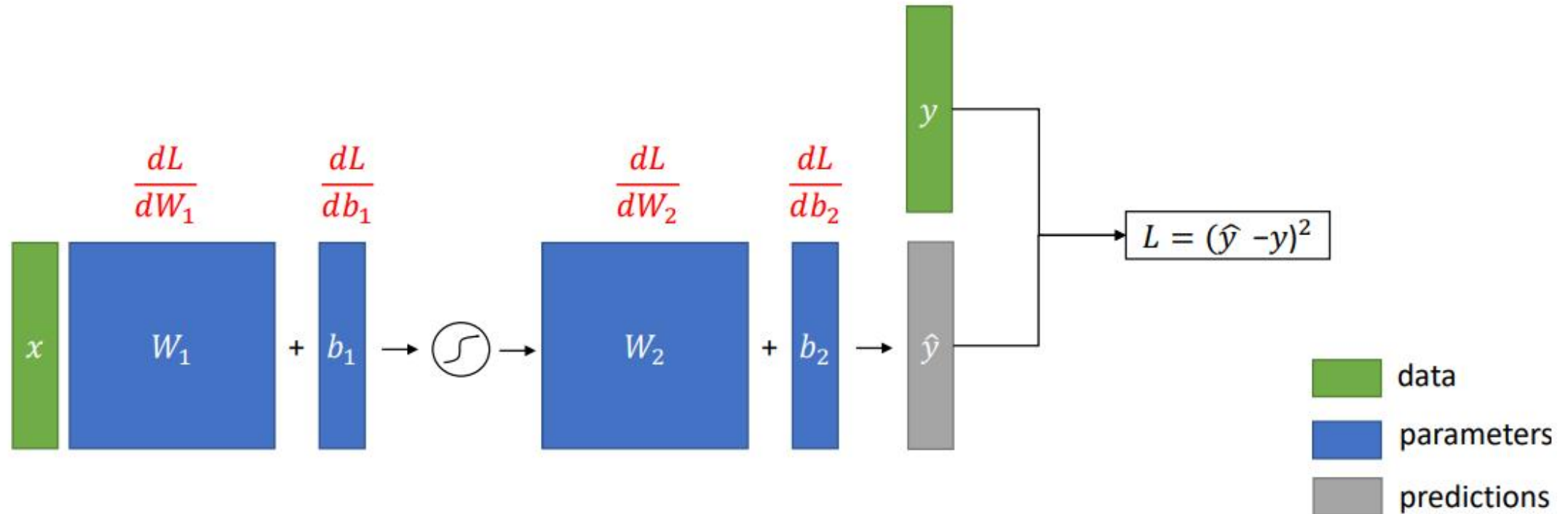
Autodifferentiation Example

Modern deep learning frameworks (eg, TensorFlow, PyTorch) allow to specify the computation graph (neural network architecture), and then automatically calculate gradients!

```
model = torch.nn.Sequential(          # define ANN with 5 input, 3 hidden
    torch.nn.Linear(5, 3),           # units and ReLU activations
    torch.nn.ReLU(),
    torch.nn.Linear(3, 1),
    torch.nn.Flatten(0, 1)
)
loss_fn = torch.nn.MSELoss()         # train with gradient descent on squared loss
optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
for t in range(1000):
    y_pred = model(x)                # forward pass
    loss = loss_fn(y_pred, y)         # compute and print loss
    optimizer.zero_grad()            # initialize gradient
    loss.backward()                   # compute gradient via backprop
    optimizer.step()                  # parameter update
```

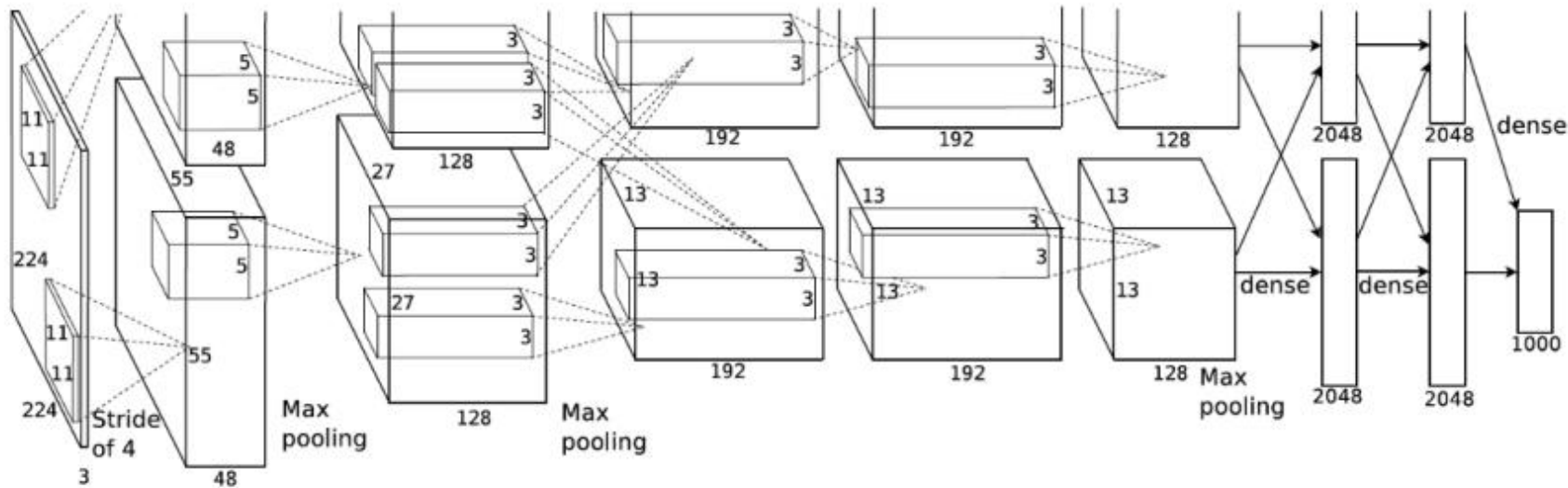
What is Autodiff?

- We only need to implement the **forward pass** (i.e., chain of operations from input to output)
- PyTorch automatically computes the **derivatives** for us → **backward pass**



Why Autodiff?

- Neural network architectures can be quite complicated
 - Deriving the gradients by hand is tedious, error prone and computationally less efficient than Autodiff



How does Autodiff work?

1) Forward Pass:

- Construct computation graph (DAG)
- Store intermediate computed values at each node of the graph

```
def f(x, y):  
    c = x + y  
    d = c**2  
    return x * d
```

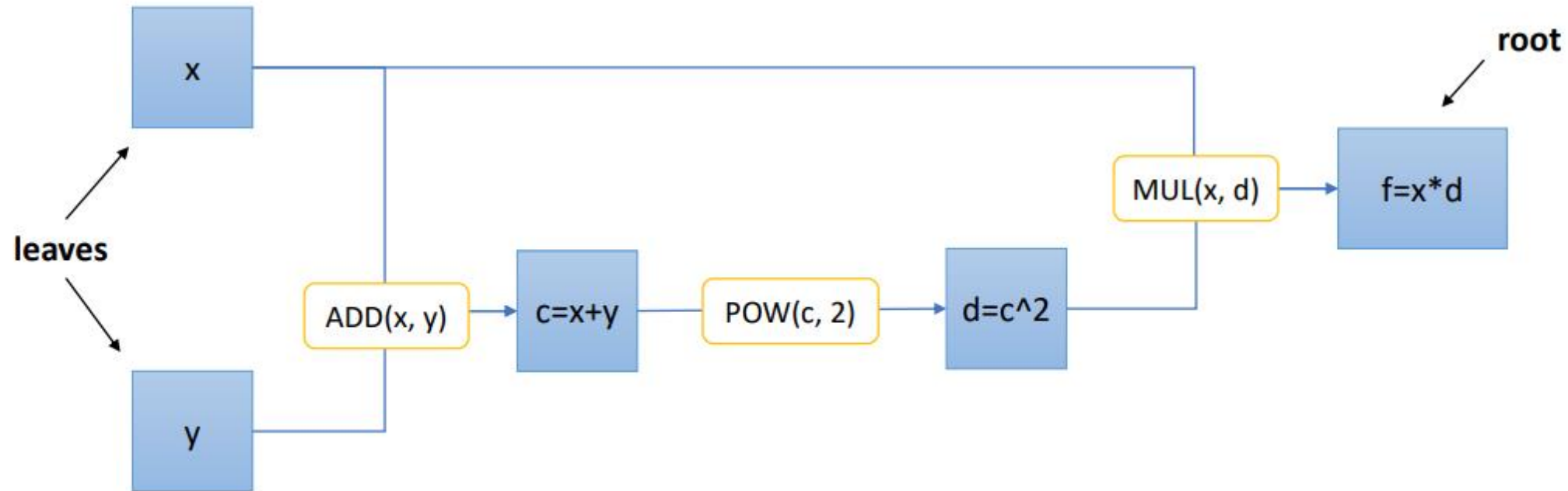
```
x = torch.tensor([3.0], requires_grad=True)  
y = torch.tensor([-1.0], requires_grad=True)  
f_val = f(x, y)
```


How does Autodiff work?

2) Backward Pass:

```
f_val.backward()
```

- Go back from the root (i.e. loss) to the leaves (i.e. parameters)



How does Autodiff work?

2) Backward Pass:

- Go back from the root (i.e. loss) to the leaves (i.e. parameters)
- Each elementary operation has a 'backward function' that computes its derivatives
- Apply **chain rule** along each path

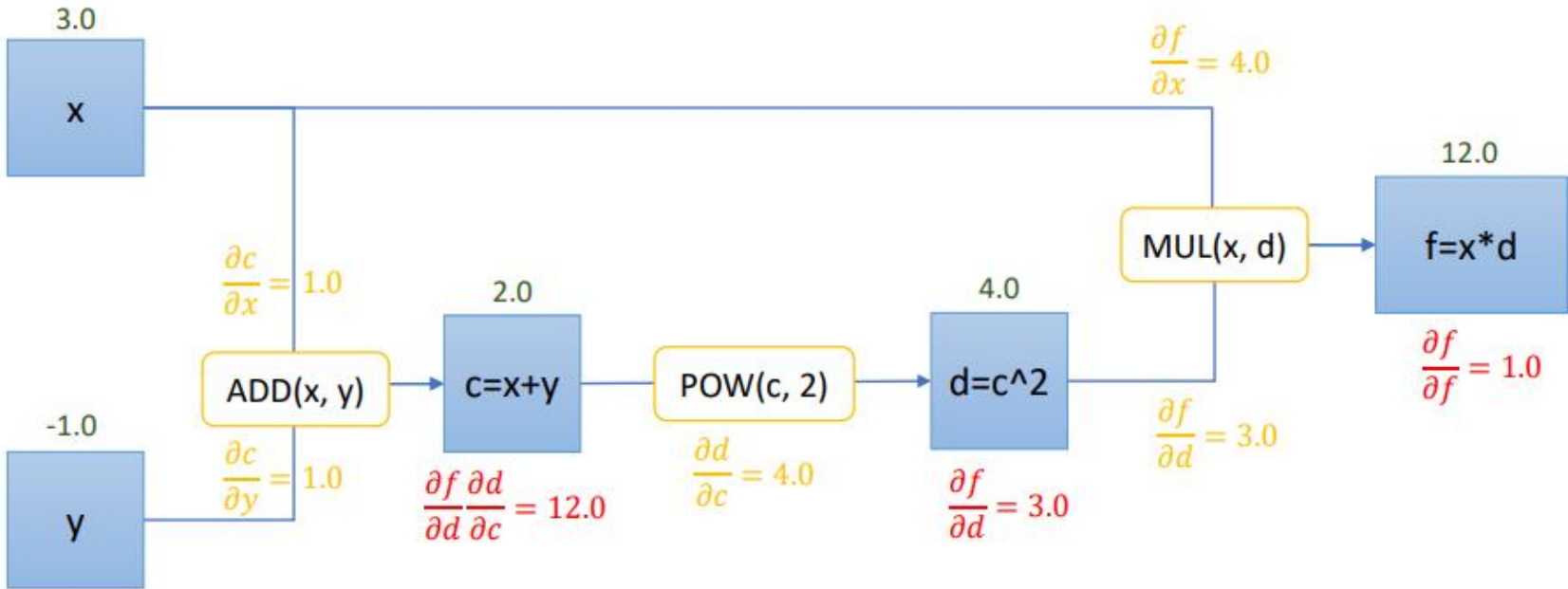
$$\frac{d f(g(x))}{dx} = \frac{df}{dg} \frac{dg}{dx} \quad \rightarrow \text{Make use of stored intermediate computation results from forward pass}$$

- Total derivative is **sum of all path derivatives** from the root to the leaf

Backward Pass

```
[31] f_val.backward()  
print('df/dx:', x.grad)  
print('df/dy:', y.grad)  
  
df/dx: tensor([16.])  
df/dy: tensor([12.])
```

$$\frac{df}{dx} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial x} + \frac{\partial f}{\partial x} = 16.0$$



$$\frac{df}{dy} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial y} = 12.0$$

Other Resource:

<https://towardsdatascience.com/understanding-backpropagation-algorithm-7bb3aa2f95fd>

Review on Backpropagation

Video tutorials:

<https://www.youtube.com/watch?v=llg3gGewQ5U>

<https://www.youtube.com/watch?v=tleHLnjs5U8>

<https://www.youtube.com/watch?v=i94OvYb6noo>

Thank You!