

Deep Learning Crash Course

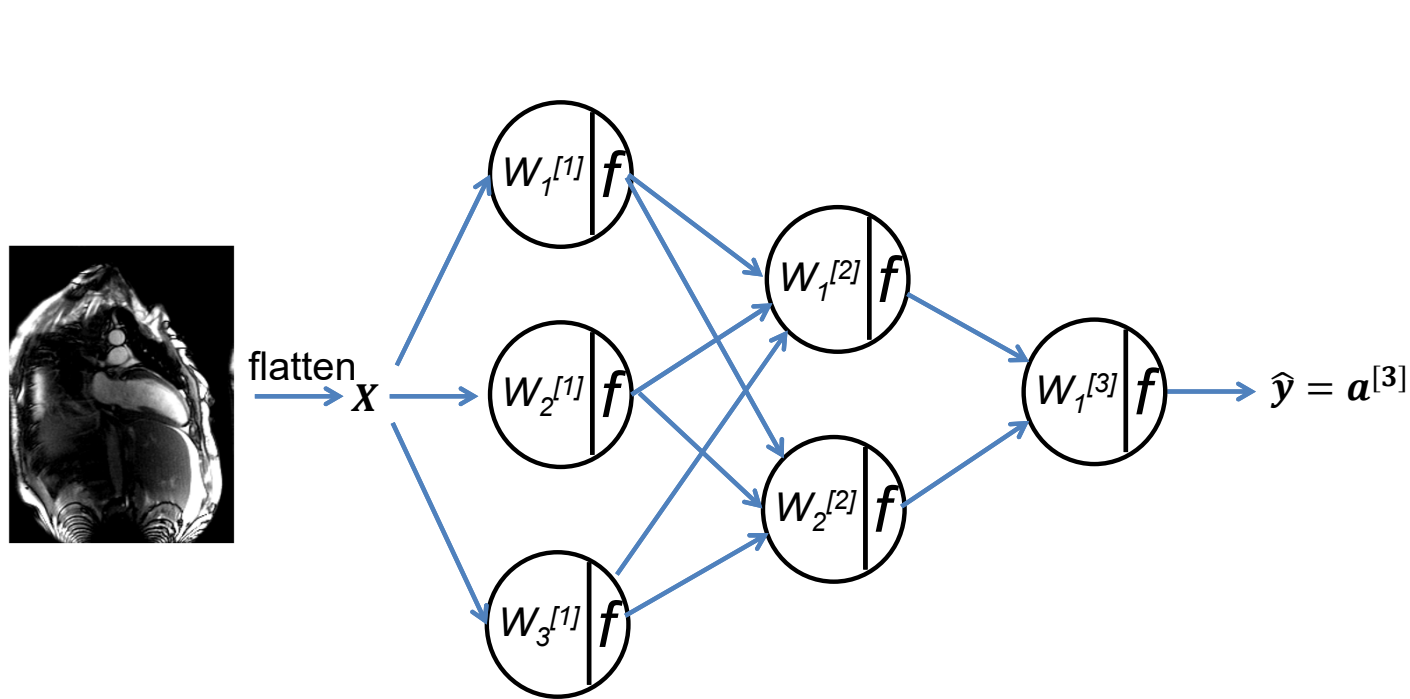


www.deeplearningcrashcourse.org

Hui Xue

Fall 2021

Recap : MLP



$W_1^{[2]}$ ← Index of layer
← Index of neuron

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ a^{[1]} &= f(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= f(Z^{[2]}) \\ Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} \\ \hat{y} &= a^{[3]} = f(Z^{[3]}) \end{aligned}$$

The i -th sample $(X^{(i)}, y^{(i)})$

Recap : Cross-Entropy Loss

Multi-class classification:

Given two probability distribution p and q , cross-entropy:

$$H(p, q) = - \sum_{\text{for all possible } x} p(x) \log[q(x)]$$

$$L_{CE_loss}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=0}^4 y[k] \log(\hat{\mathbf{y}}[k])$$

$$L_{CE_loss}(\mathbf{y}, \hat{\mathbf{y}}) = -\log(\hat{\mathbf{y}}[y_k])$$

Binary classification:

$$L_{BCE_loss}(\mathbf{y}, \hat{\mathbf{y}}) = -[y \cdot \log(\hat{\mathbf{y}}) + (1 - y) \cdot \log(1 - \hat{\mathbf{y}})]$$

Recap : Mini-batch Gradient Descent

Initialize weights and bias

Random shuffle dataset

BatchSize = 32

for epoch in range(E):

select #BatchSize samples

Evaluate loss function (forward pass) at this **batch**

$$L = \frac{1}{B} \sum_{i=0}^{B-1} L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}(\mathbf{W}, \mathbf{b}))$$

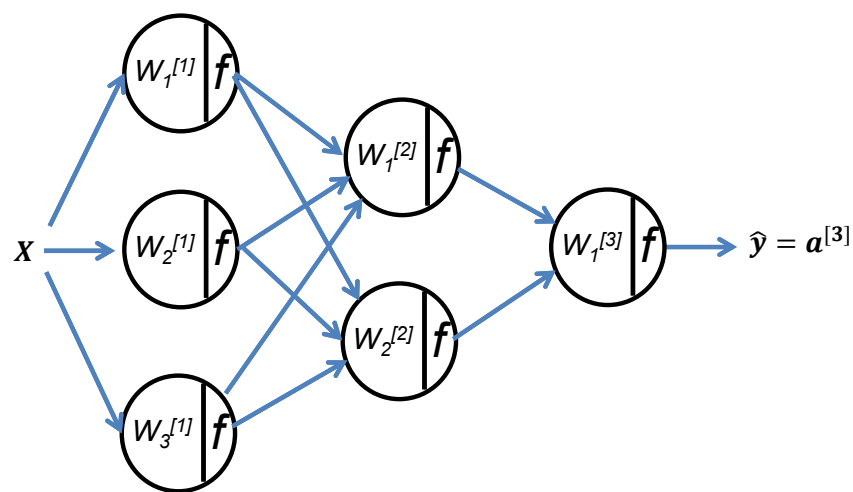
Compute gradient $\frac{\partial L}{\partial \mathbf{W}^{[l]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[l]}}$

Update parameter $\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{W}^{[l]}}$, $\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{b}^{[l]}}$

Lecture 3

- **Backprop**

To find the good model parameter



Evaluate loss function (forward pass) at this **batch**

$$L = \frac{1}{B} \sum_{i=0}^{B-1} L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}(\mathbf{W}, \mathbf{b}))$$

Compute gradient $\frac{\partial L}{\partial \mathbf{w}^{[l]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[l]}}$

Update parameter $\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{w}^{[l]}}$, $\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{b}^{[l]}}$

Evaluate loss function : A forward pass

Evaluate loss function (forward pass) at this **batch**

$$L = \frac{1}{B} \sum_{i=0}^{B-1} L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}(\mathbf{W}, \mathbf{b}))$$

Let's look at one sample and write $\mathbf{y}^{(i)}$ as \mathbf{y} :

$$L(y, \hat{y}(\mathbf{W}, \mathbf{b})) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

For every sample, we input (\mathbf{X}, y) and model outputs \hat{y}

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = \mathbf{f}(\mathbf{Z}^{[1]})$$

$$\mathbf{Z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$\mathbf{a}^{[2]} = \mathbf{f}(\mathbf{Z}^{[2]})$$

$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]} \mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\hat{y} = a^{[3]} = \mathbf{f}(\mathbf{Z}^{[3]})$$

Compute gradient : Is it easy to do?

Compute gradient $\frac{\partial L}{\partial \mathbf{W}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[3]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[3]}}$

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = f(\mathbf{Z}^{[1]})$$

$$\mathbf{Z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$\mathbf{a}^{[2]} = f(\mathbf{Z}^{[2]})$$

$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[3]} = f(\mathbf{Z}^{[3]})$$

First idea, compute analytical derivative

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{a}^{[3]} = f(\mathbf{Z}^{[3]}) \\ &= f(\mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]})\end{aligned}$$

$$= f(\mathbf{W}^{[3]}f(\mathbf{Z}^{[2]}) + \mathbf{b}^{[3]})$$

$$= f(\mathbf{W}^{[3]}f(\mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]})$$

$$= f(\mathbf{W}^{[3]}f(\mathbf{W}^{[2]}f(\mathbf{Z}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]})$$

$$= f(\mathbf{W}^{[3]}f(\mathbf{W}^{[2]}f(\mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]})$$

Compute analytical gradient : hard to do

Compute gradient $\frac{\partial L}{\partial \mathbf{W}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[3]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[3]}}$

So, let's compute:

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = f(\mathbf{Z}^{[1]})$$

$$\mathbf{Z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$\mathbf{a}^{[2]} = f(\mathbf{Z}^{[2]})$$

$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

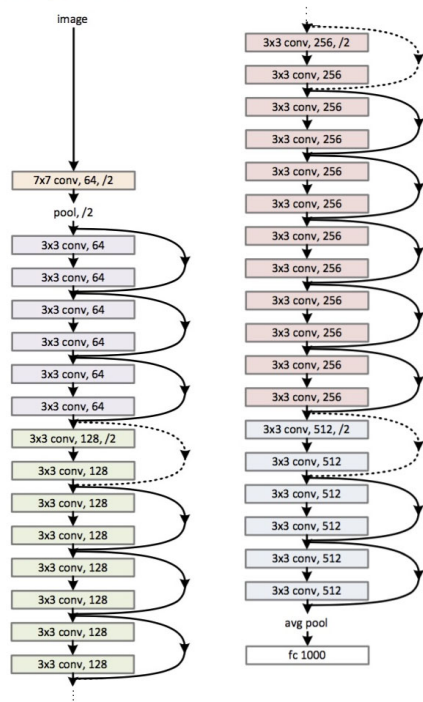
$$\hat{y} = \mathbf{a}^{[3]} = f(\mathbf{Z}^{[3]})$$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}^{[1]}} &= - \frac{\partial \{ [y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})] \}}{\partial \mathbf{W}^{[1]}} \\ &= - \frac{\partial \left\{ \left[\begin{aligned} &y \log \left(f(\mathbf{W}^{[3]} f(\mathbf{W}^{[2]} f(\mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) \right) \right. \right. \\ &\left. \left. + (1 - y) \log \left(1 - f(\mathbf{W}^{[3]} f(\mathbf{W}^{[2]} f(\mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) \right) \right] \right\}}{\partial \mathbf{W}^{[1]}} \end{aligned}$$

Compute derivative gradient is hard ... for this small MLP

Compute derivative gradient : cannot be done

34-layer residual



Deep Residual Learning for Image Recognition. <https://arxiv.org/pdf/1512.03385.pdf>



Compute analytical gradient becomes impossible for deep neural networks !

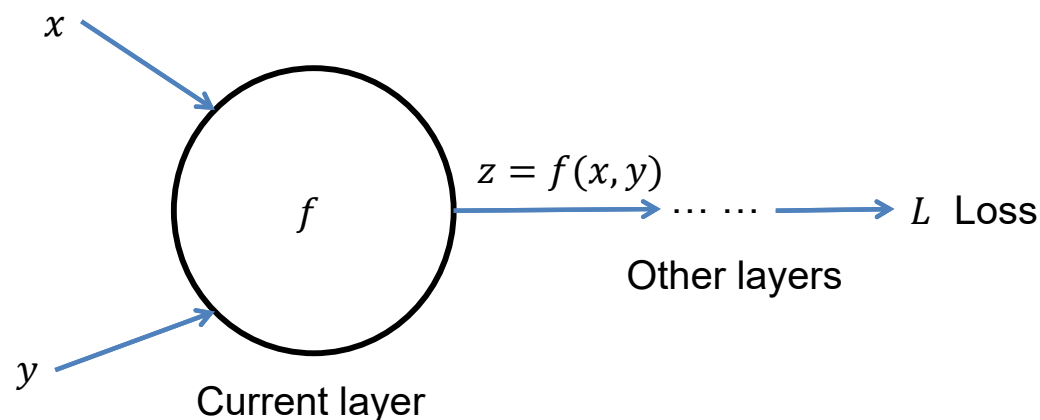
Don't try it!

Compute gradient : Back propagation

Compute derivative gradient becomes impossible for deep neural networks. We need back prop.

Three ideas behind BackProp:

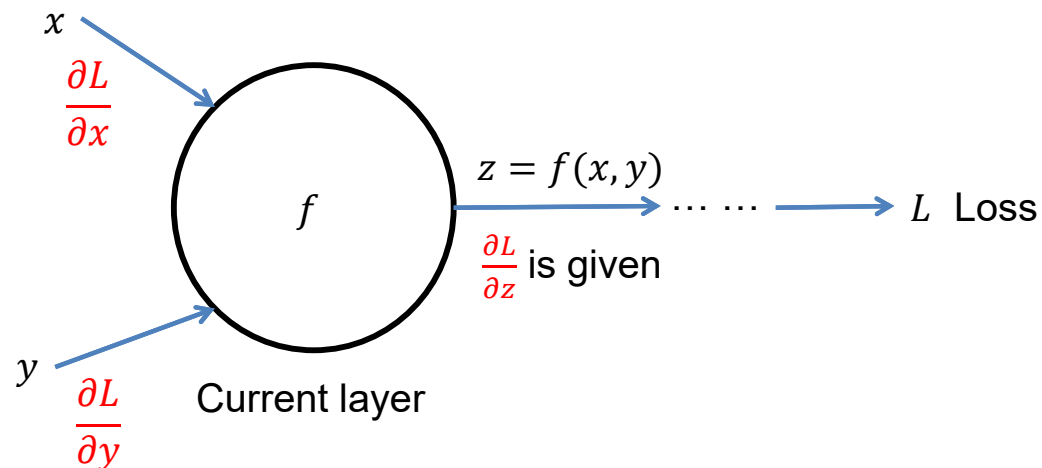
- Divide and conquer
- Chain-rule
- Let computer do the job - AutoDiff



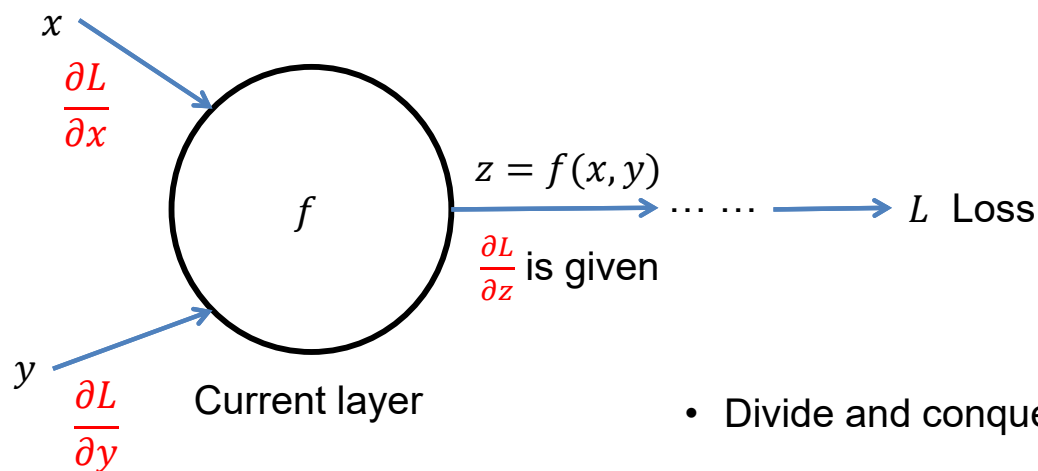
Compute gradient : Back propagation

For this current layer, what we want to do:

- Given the gradient from the loss to output parameter z
- Compute gradient from the loss to input x and y



Compute gradient : Back propagation



- Chain rule

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

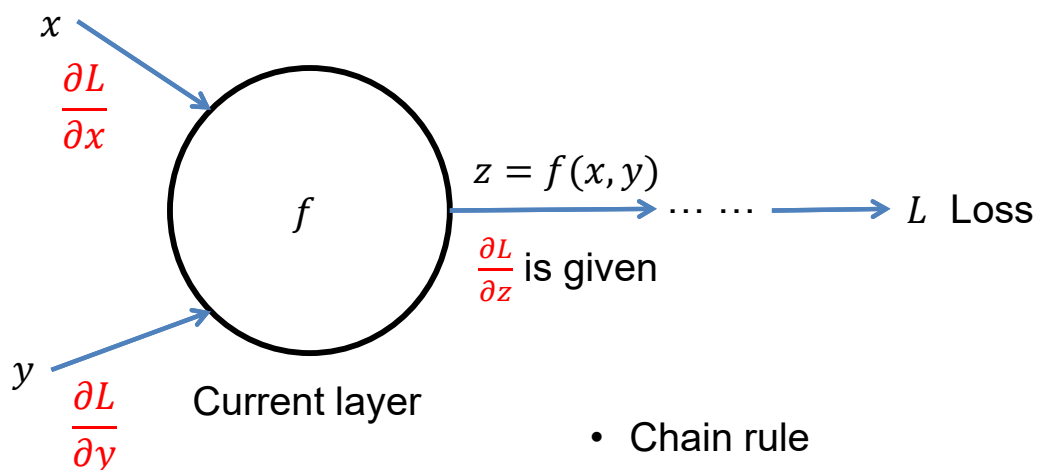
$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

- Divide and conquer

We only need to know:

- 1) the current layer to produce z from x and y
- 2) the upstream gradient

Compute gradient : Back propagation



- Chain rule

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial f(x,y)}{\partial x}$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial f(x,y)}{\partial y}$$

- Divide and conquer

We know the current **function** to compute z from x and y

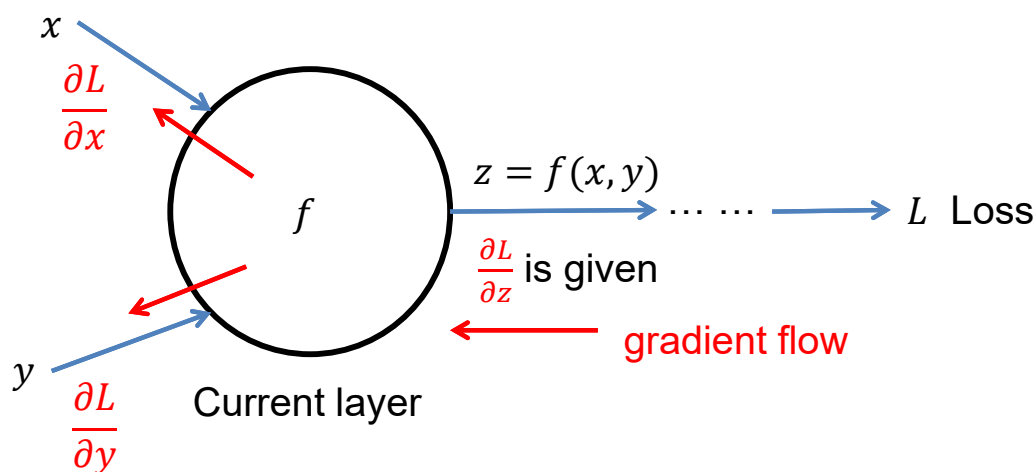
We know the gradient of loss to output z

Enough information to compute derivative of loss to input x and y

$$\frac{\partial z}{\partial x} = \frac{\partial f(x,y)}{\partial x} \quad \frac{\partial z}{\partial y} = \frac{\partial f(x,y)}{\partial y}$$

Compute gradient : Back propagation

To compute current gradient, only need **upstream gradient** and **local function**

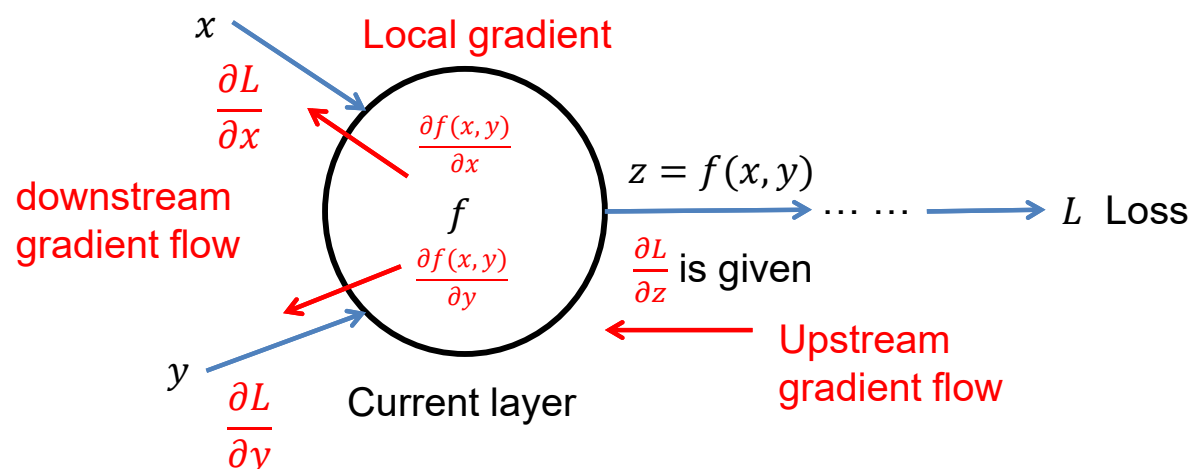


$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial f(x, y)}{\partial y}$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial f(x, y)}{\partial x}$$

Compute gradient : Back propagation

To compute current gradient, only need **upstream gradient** and **local gradient**



- 1) Inputs flow through the model to compute loss
- 2) Gradient of loss flows back to compute gradient to inputs (and all others ...)
- 3) Upstream gradient is combined with local gradient to compute downstream gradient

Computational Graph

$$L = x^2 + 2xy \quad x = 1, y = 4$$

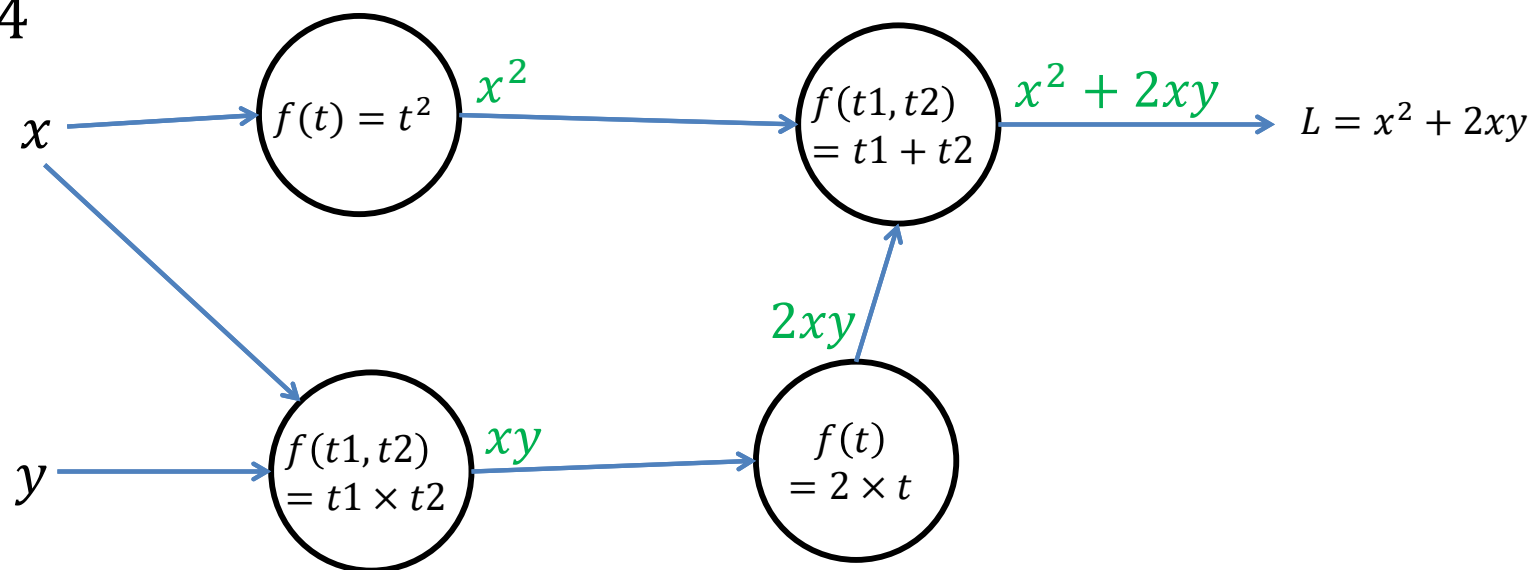
$$\frac{\partial L}{\partial x} = 2x + 2y = 2 + 8 = 10$$

$$\frac{\partial L}{\partial y} = 2x = 2$$

Computational Graph

$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

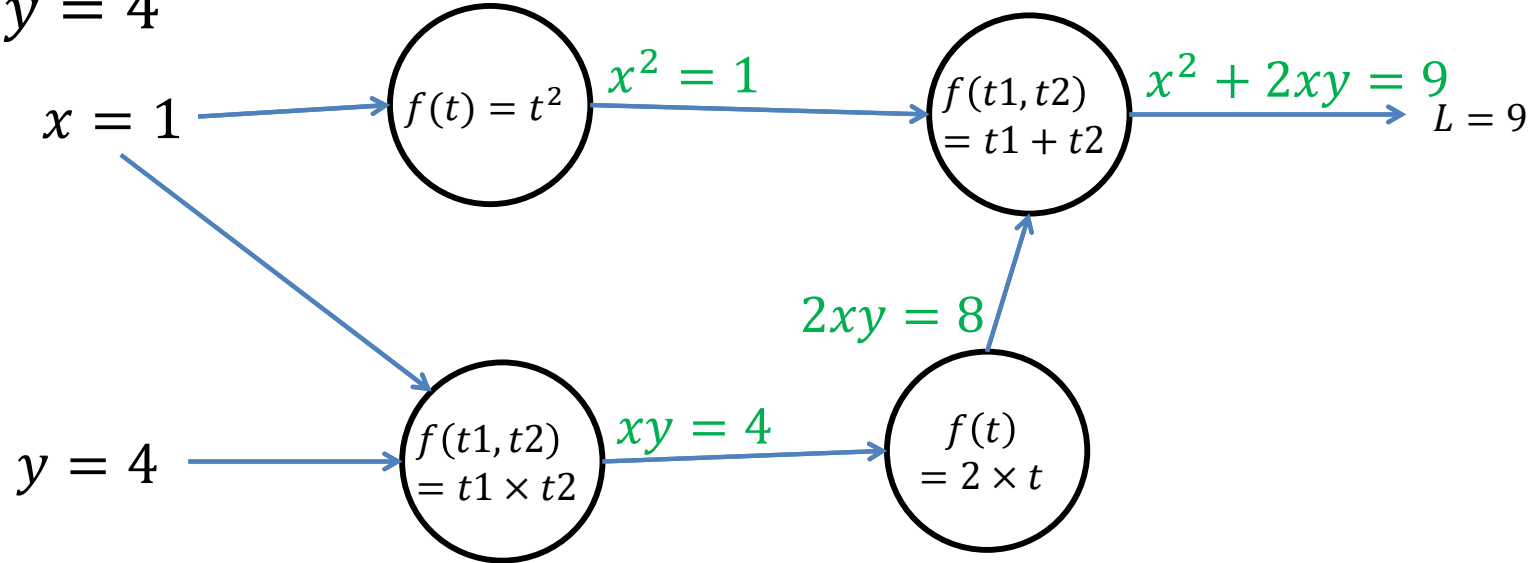


Computational Graph

$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

Forward pass

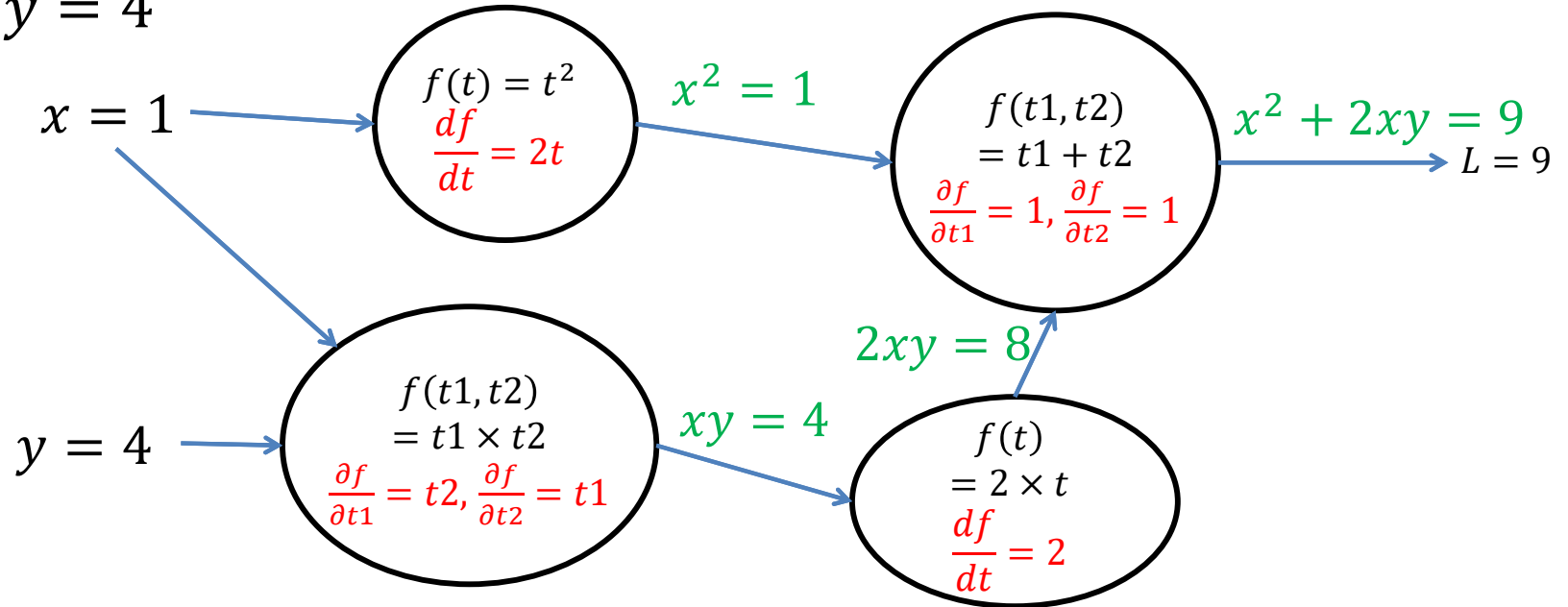


Computational Graph

$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

Local gradient

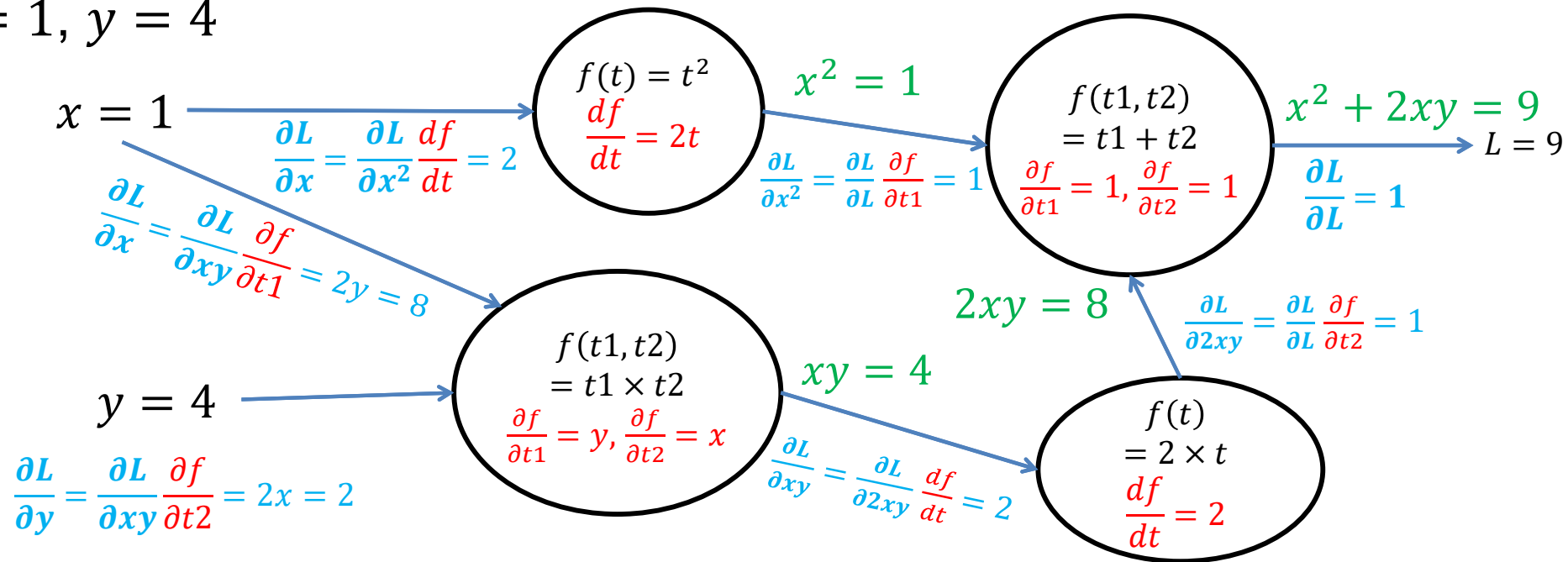


Computational Graph

$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

Back prop



Computational Graph

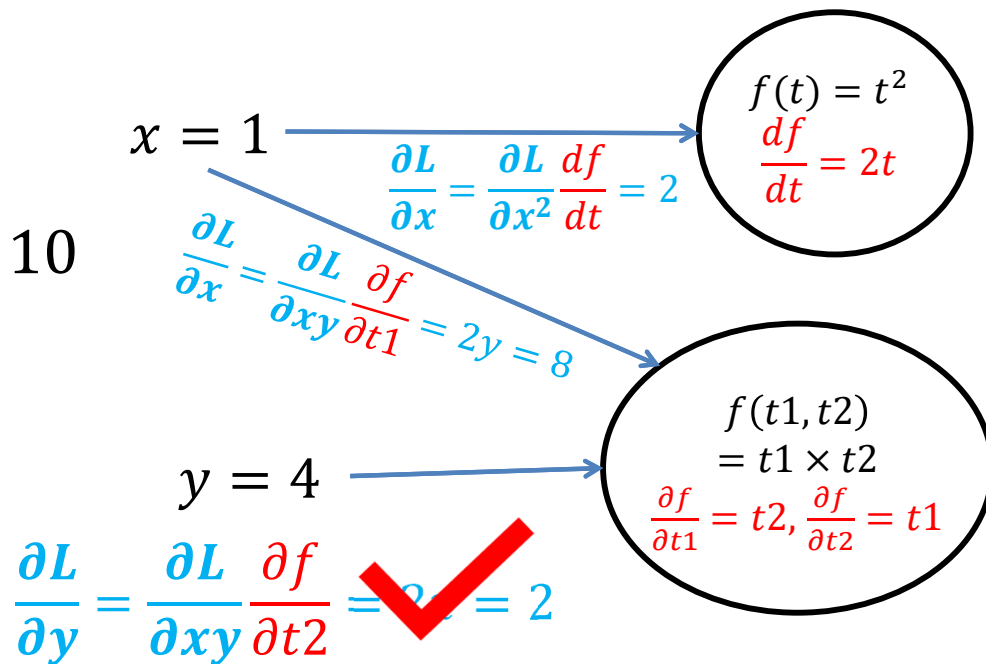
$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

$$\frac{\partial L}{\partial x} = 2x + 2y = 2 + 8 = 10$$

$$\frac{\partial L}{\partial y} = 2x = 2$$

Back prop



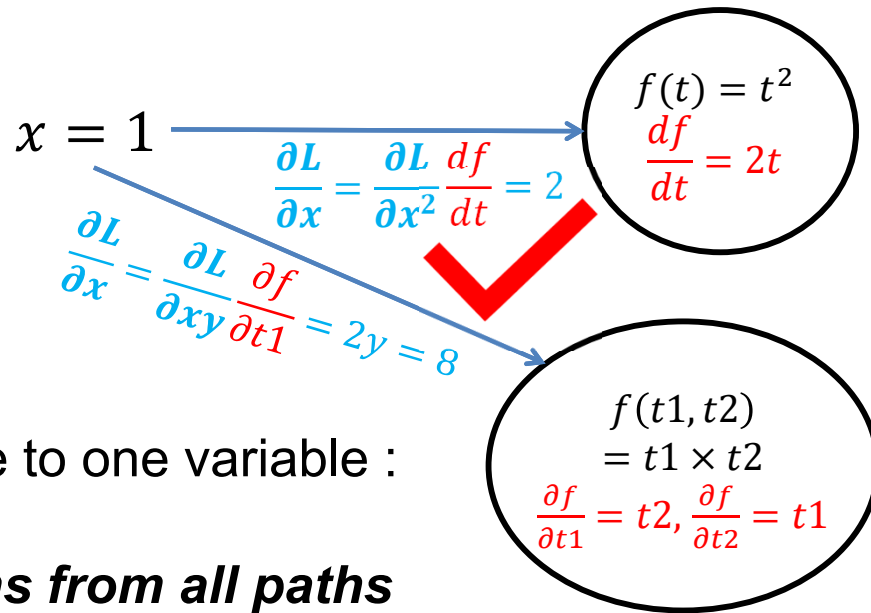
Computational Graph

$$L = x^2 + 2xy$$

$$x = 1, y = 4$$

$$\frac{\partial L}{\partial x} = 2x + 2y = 2 + 8 = 10$$

Back prop



Two paths contribute to one variable :

Sum contributions from all paths

Derivative of Scalar to Vector and Matrix

Loss is a scalar $L \in \mathbb{R}$

If x is a scalar, then

$\frac{\partial L}{\partial x}$ is a scalar

Derivative of a scalar function to a scalar variable is a scalar.

If x is a vector, $x \in \mathbb{R}^{N \times 1}$

$\frac{\partial L}{\partial x}$ is a vector in the same shape as x

$$x = [x_0 \quad x_1 \quad \dots \quad x_{N-2} \quad x_{N-1}]^T$$

$$\frac{\partial L}{\partial x} = \left[\frac{\partial L}{\partial x_0} \quad \frac{\partial L}{\partial x_1} \quad \dots \quad \frac{\partial L}{\partial x_{N-2}} \quad \frac{\partial L}{\partial x_{N-1}} \right]^T$$

If x is a matrix, $x \in \mathbb{R}^{N \times M}$

$\frac{\partial L}{\partial x}$ still has in the same shape as x

$$x = \begin{bmatrix} x_{00} & x_{01} & \dots & x_{0,M-1} \\ x_{10} & x_{11} & \dots & x_{1,M-1} \\ \dots & \dots & \dots & \dots \\ x_{N-1,0} & x_{N-1,M-1} & \dots & x_{N-1,M-1} \end{bmatrix}$$

$$\frac{\partial L}{\partial x} = \left[\frac{\partial L}{\partial x_{ij}} \right]_{i=0:N-1, j=0:M-1}$$

Vector to Vector : Jacobian matrix

$$f \in \mathbb{R}^N \rightarrow \mathbb{R}^M \quad \mathbf{y} = f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^N \quad \mathbf{y} \in \mathbb{R}^M$$

\mathbf{x} is a vector of length N $\mathbf{x} = [x_0 \quad x_1 \quad \dots \quad x_{N-2} \quad x_{N-1}]^T$

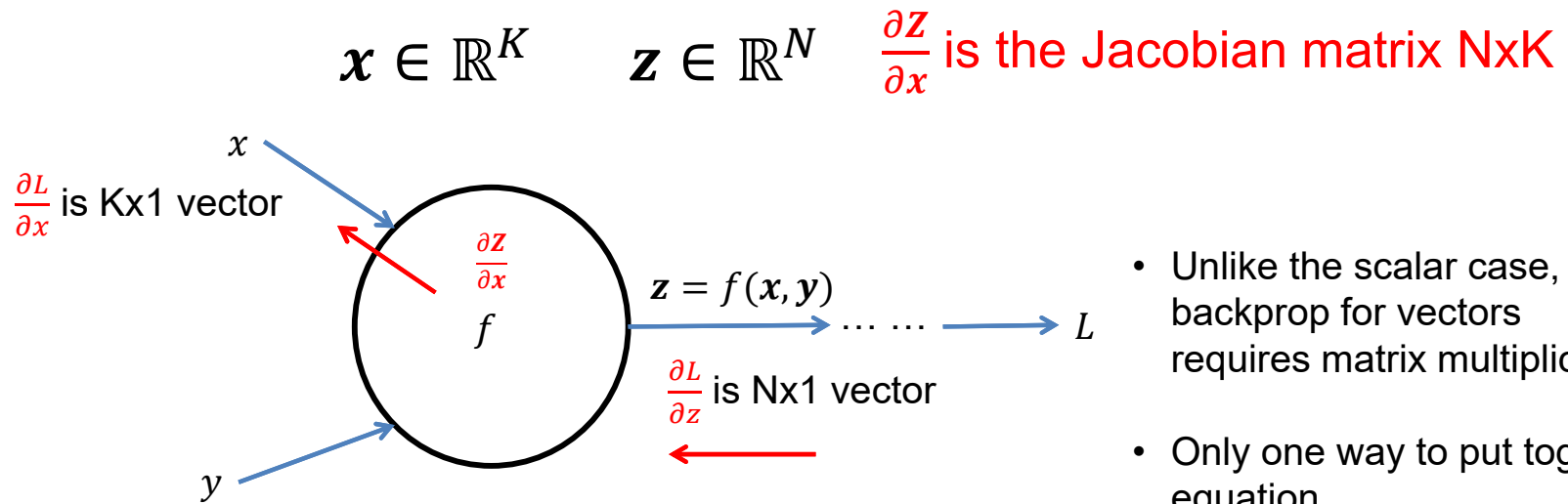
\mathbf{y} is a vector of length M $\mathbf{y} = [y_0 \quad y_1 \quad \dots \quad y_{M-2} \quad y_{M-1}]^T$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[\frac{\partial y_i}{\partial x_j} \right]_{i=0:M-1, j=0:N-1} \quad \begin{bmatrix} \frac{\partial y_0}{\partial x_0} & \dots & \frac{\partial y_0}{\partial x_{N-1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_{M-1}}{\partial x_0} & \dots & \frac{\partial y_{M-1}}{\partial x_{N-1}} \end{bmatrix} \quad \begin{array}{l} \text{Jacobian matrix:} \\ \\ \text{MxN} \\ \text{output length x input length} \end{array}$$

Derivative of every element of \mathbf{y} to every element of \mathbf{x}

Backprop for Vector

Apply Jacobian matrix as local gradient

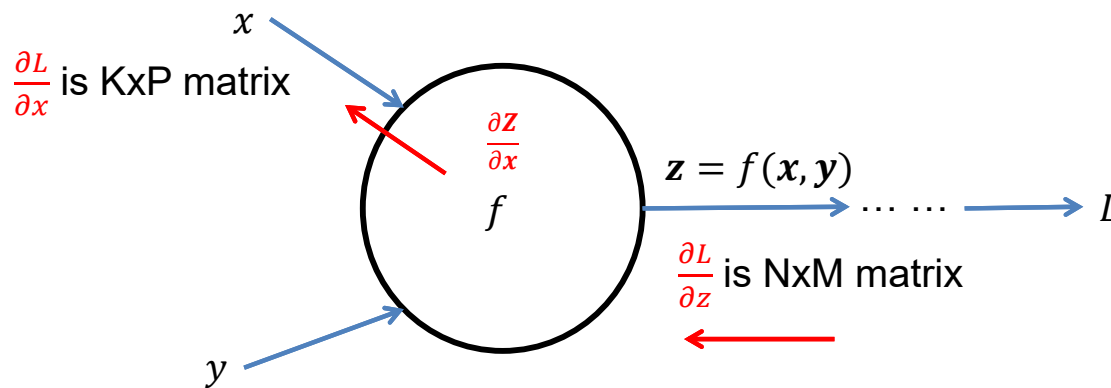


$$\frac{\partial L}{\partial \mathbf{x}} = \left[\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right]^T \cdot \frac{\partial L}{\partial \mathbf{z}} \quad \frac{\partial L}{\partial \mathbf{y}} = \left[\frac{\partial \mathbf{z}}{\partial \mathbf{y}} \right]^T \cdot \frac{\partial L}{\partial \mathbf{z}}$$

Backprop for Matrix and Tensors

$$\mathbf{x} \in \mathbb{R}^{K \times P} \quad \mathbf{z} \in \mathbb{R}^{N \times M}$$

$\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ is the Jacobian Tensor
(N x M) x (K x P)

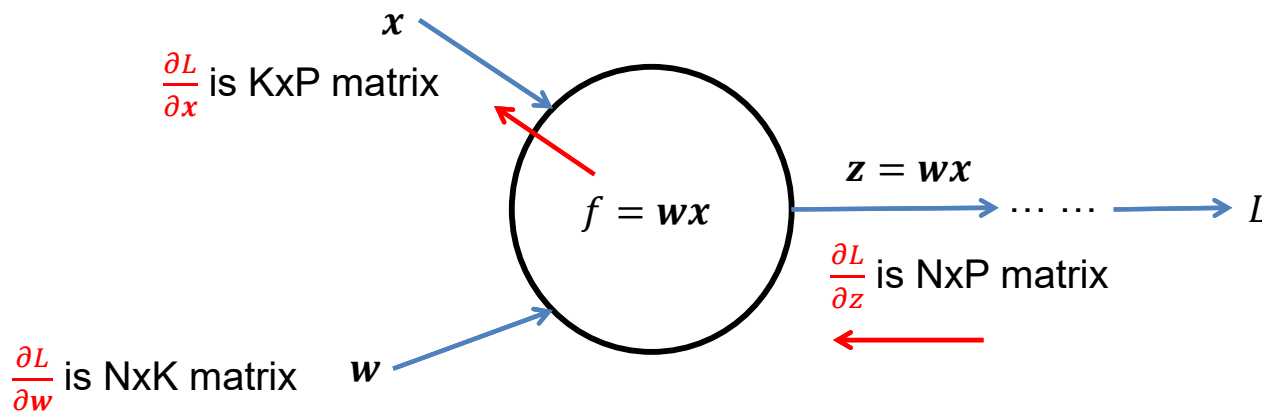


- Often use backprop to compute Jacobian Tensor
- Still only one way to put together equation

$$\frac{\partial L}{\partial \mathbf{x}} = \left[\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right]^T \cdot \frac{\partial L}{\partial \mathbf{z}} \quad \frac{\partial L}{\partial \mathbf{y}} = \left[\frac{\partial \mathbf{z}}{\partial \mathbf{y}} \right]^T \cdot \frac{\partial L}{\partial \mathbf{z}}$$

Backprop for Matrix multiplication

$$\mathbf{x} \in \mathbb{R}^{K \times P} \quad \mathbf{w} \in \mathbb{R}^{N \times K} \quad \mathbf{z} \in \mathbb{R}^{N \times P}$$



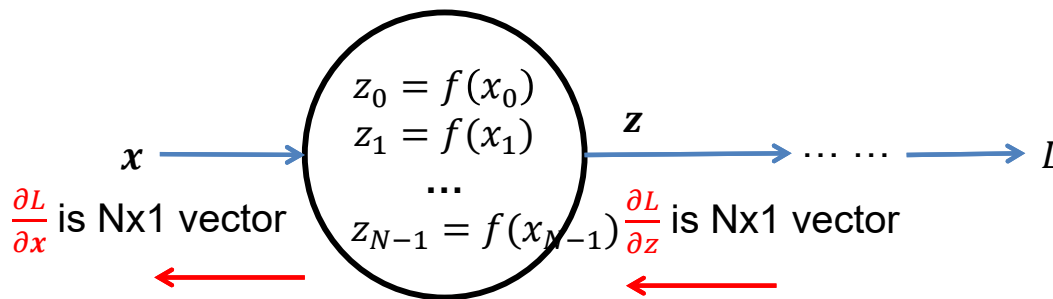
- Matrix multiplication formula is useful
- Only one way to put together them

$$\frac{\partial L}{\partial \mathbf{x}} = \mathbf{w}^T \cdot \frac{\partial L}{\partial \mathbf{z}}$$

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial L}{\partial \mathbf{z}} \cdot \mathbf{x}^T$$

Backprop for Element-wise operation

$$\mathbf{x} \in \mathbb{R}^{N \times 1} \quad \mathbf{z} \in \mathbb{R}^{N \times 1}$$

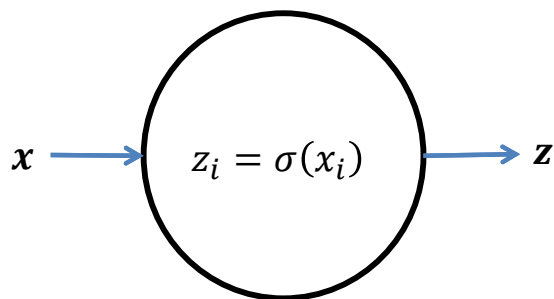


- Simply compute backprop element by element

$$\frac{\partial L}{\partial x_i} = \frac{df(x_i)}{dx_i} \frac{\partial L}{\partial z_i} \quad i = 0:N-1$$

Backprop for Sigmoid

$$\mathbf{x} \in \mathbb{R}^{N \times 1} \quad \mathbf{z} \in \mathbb{R}^{N \times 1} \quad \frac{\partial L}{\partial x_i} = \frac{df(x_i)}{dx_i} \frac{\partial L}{\partial z_i} \quad i = 0:N-1$$



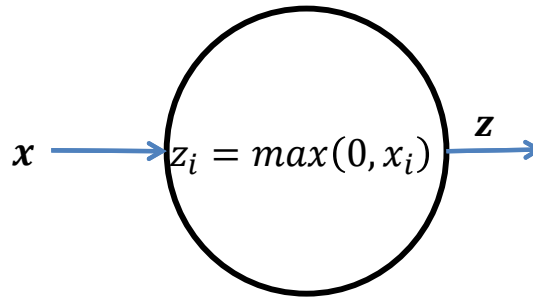
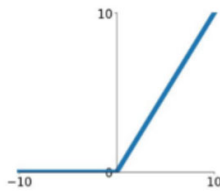
- Simply compute backprop element by element

$$\frac{d\sigma(t)}{dt} = \sigma(t)[1 - \sigma(t)] \quad \frac{\partial L}{\partial x_i} = \sigma(x_i)[1 - \sigma(x_i)] \frac{\partial L}{\partial z_i}$$

Backprop for ReLU activation function

$$\mathbf{x} \in \mathbb{R}^{N \times 1} \quad \mathbf{z} \in \mathbb{R}^{N \times 1} \quad \frac{\partial L}{\partial x_i} = \frac{df(x_i)}{dx_i} \frac{\partial L}{\partial z_i} \quad i = 0:N-1$$

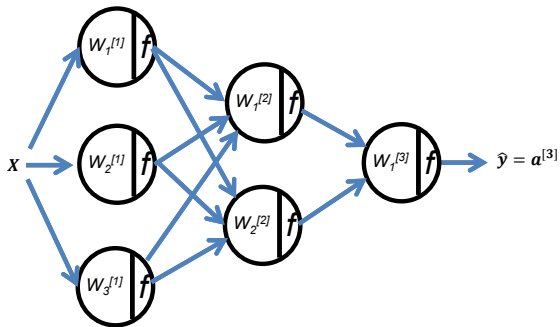
ReLU
 $\max(0, x)$



- compute backprop element by element

$$\frac{d\max(0, t)}{dt} = \begin{cases} 1, t > 0 \\ 0.5, t = 0 \\ 0, otherwise \end{cases} \quad \frac{\partial L}{\partial x_i} = \begin{cases} \frac{\partial L}{\partial z_i}, x_i > 0 \\ \frac{\partial L}{2\partial z_i}, x_i = 0 \\ 0, x_i < 0 \end{cases}$$

Computational Graph for MLP



$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = \mathbf{f}(\mathbf{Z}^{[1]})$$

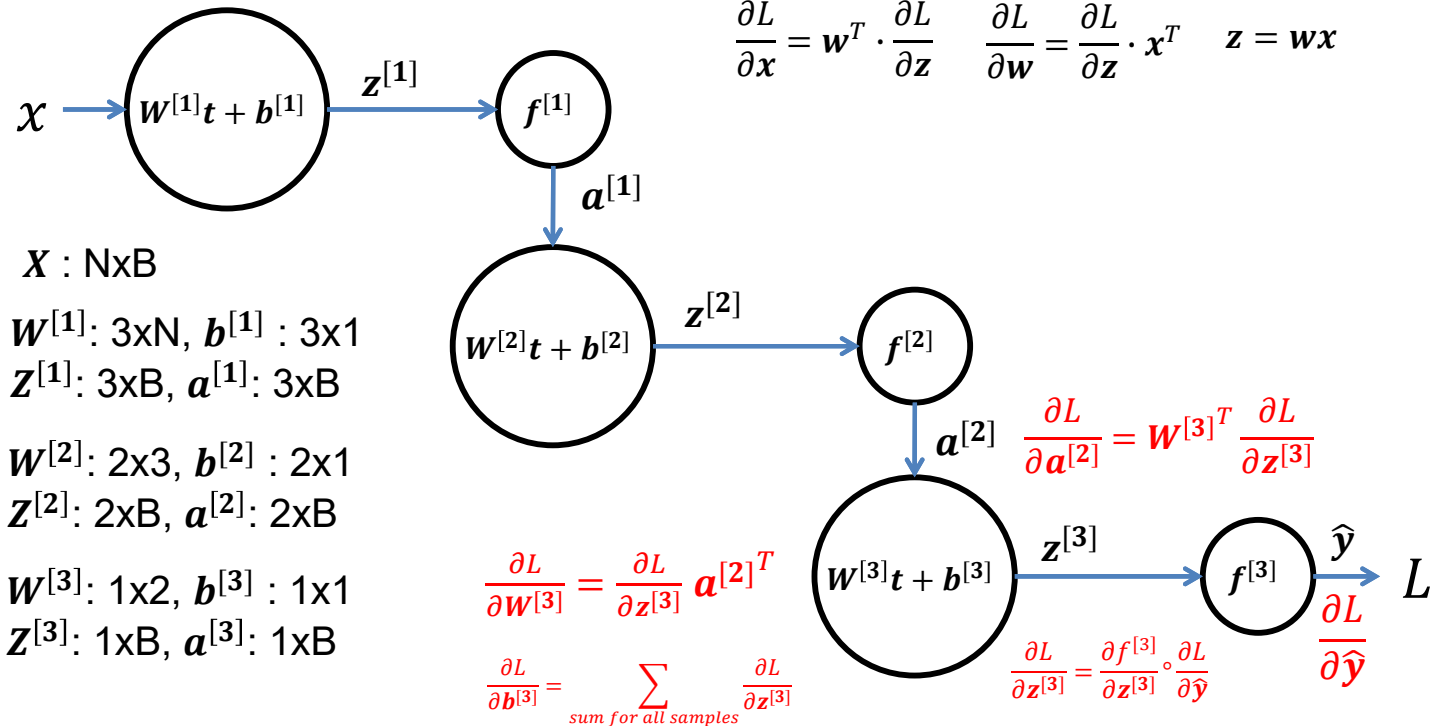
$$\mathbf{Z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$\mathbf{a}^{[2]} = \mathbf{f}(\mathbf{Z}^{[2]})$$

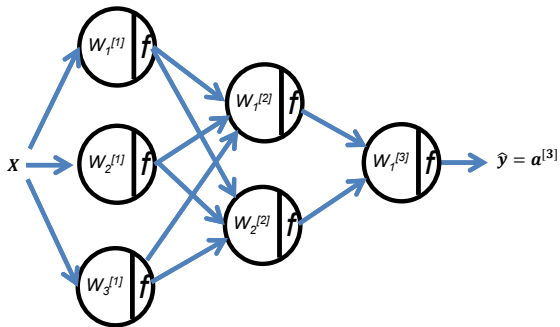
$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[3]} = \mathbf{f}(\mathbf{Z}^{[3]})$$

Compute gradient $\frac{\partial L}{\partial \mathbf{W}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[3]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[3]}}$

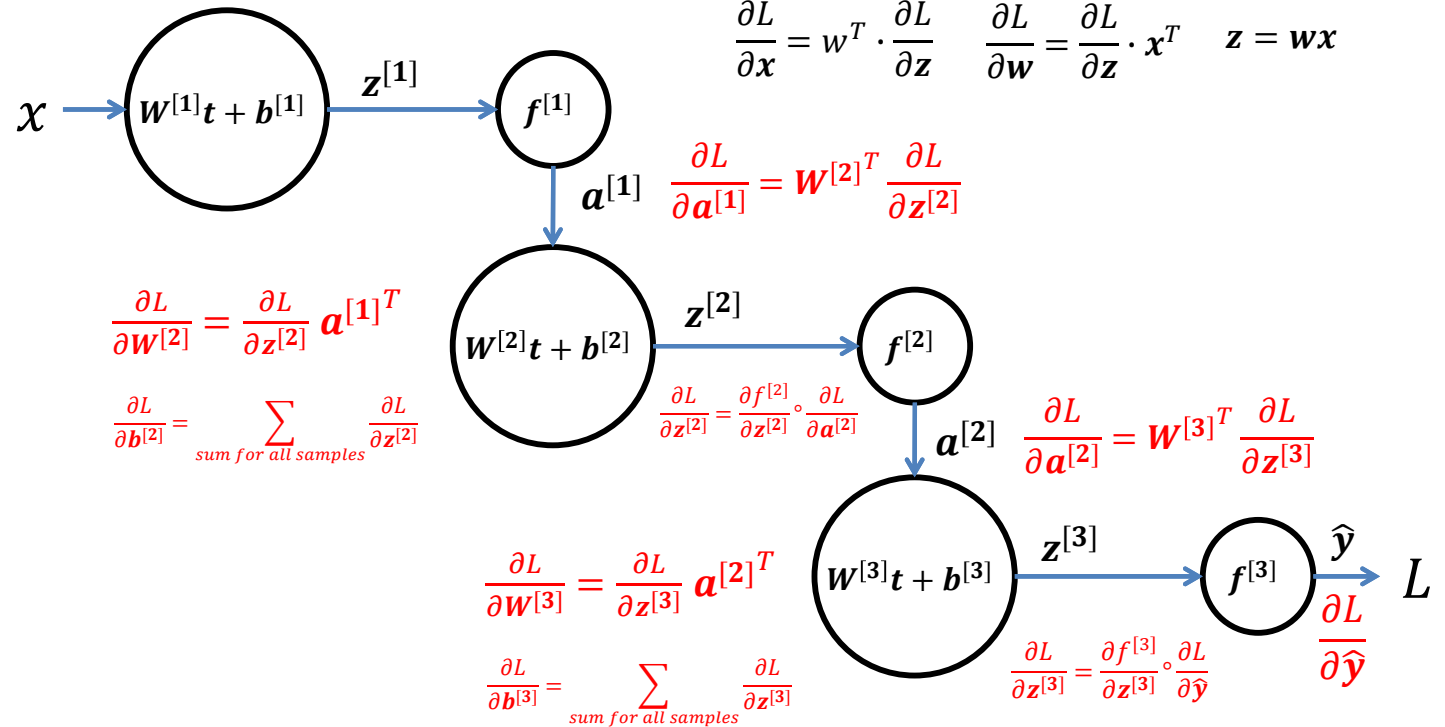


Computational Graph for MLP



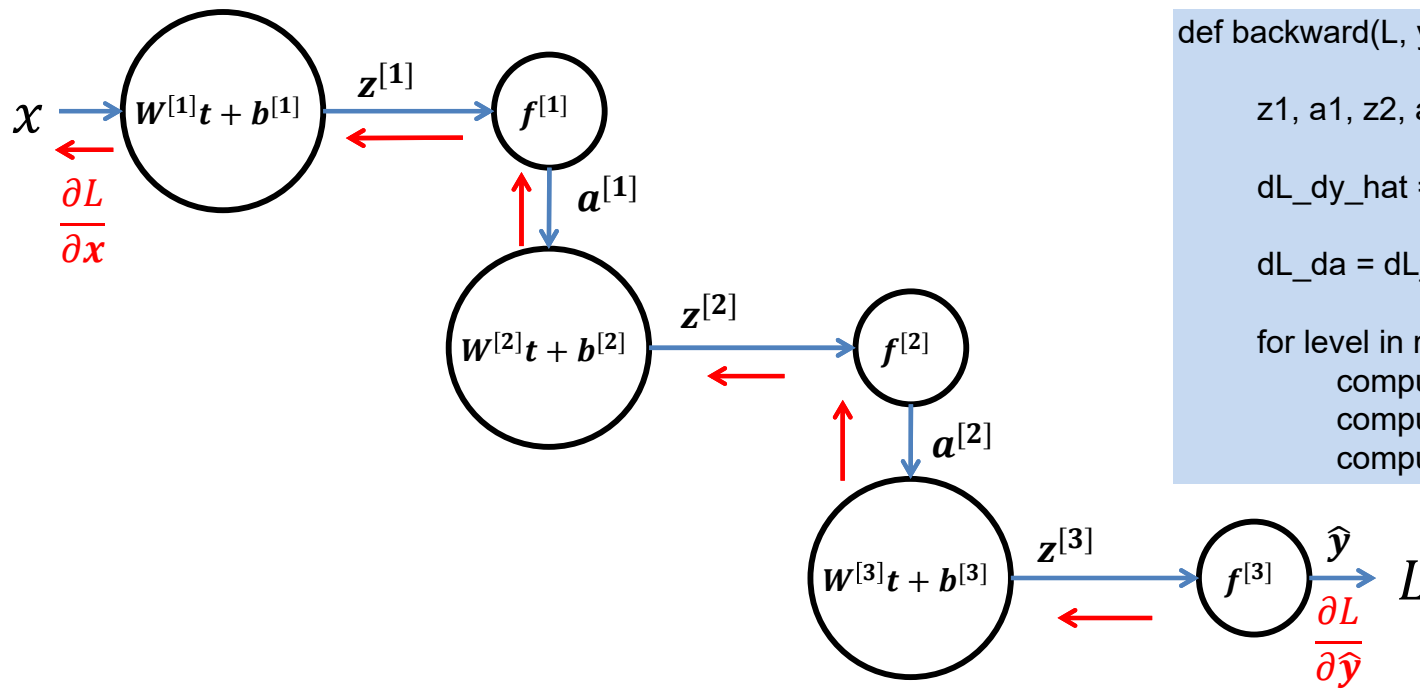
$$\begin{aligned} \mathbf{Z}^{[1]} &= \mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]} \\ \mathbf{a}^{[1]} &= \mathbf{f}(\mathbf{Z}^{[1]}) \\ \mathbf{Z}^{[2]} &= \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]} \\ \mathbf{a}^{[2]} &= \mathbf{f}(\mathbf{Z}^{[2]}) \\ \mathbf{Z}^{[3]} &= \mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]} \\ \hat{\mathbf{y}} &= \mathbf{a}^{[3]} = \mathbf{f}(\mathbf{Z}^{[3]}) \end{aligned}$$

Compute gradient $\frac{\partial L}{\partial \mathbf{W}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[1]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[2]}}$, $\frac{\partial L}{\partial \mathbf{W}^{[3]}}$, $\frac{\partial L}{\partial \mathbf{b}^{[3]}}$



Backprop for MLP

Compute gradient $\frac{\partial L}{\partial W^{[1]}}$, $\frac{\partial L}{\partial b^{[1]}}$, $\frac{\partial L}{\partial W^{[2]}}$, $\frac{\partial L}{\partial b^{[2]}}$, $\frac{\partial L}{\partial W^{[3]}}$, $\frac{\partial L}{\partial b^{[3]}}$



```
def backward(L, y, y_hat, buffer):
```

```
    z1, a1, z2, a2, z3 = buffer
```

```
    dL_dy_hat = compute_loss_gradient(L, y, y_hat)
```

```
    dL_da = dL_dy_hat
```

```
    for level in range(3,0, -1):
```

```
        compute dL_dz
```

```
        compute dL_dW, dL_db
```

```
        compute new dL_da
```

Start the backprop

Given a batch (\mathbf{x}, \mathbf{y}) , the mini-batch loss is:

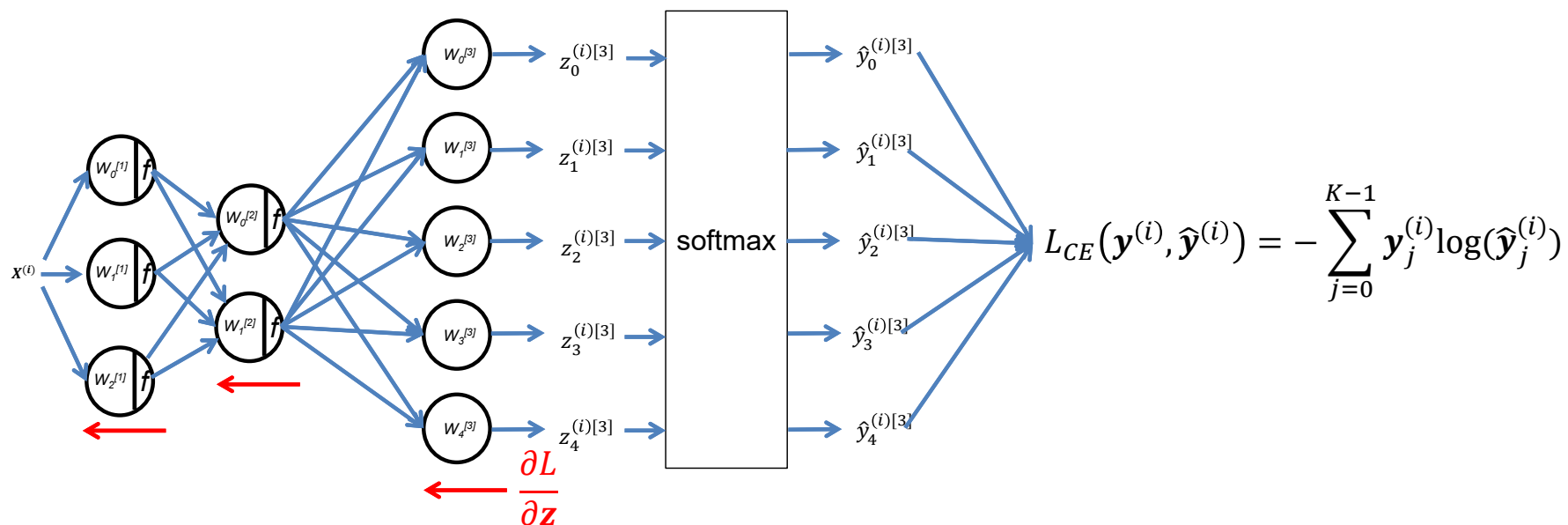
$$L = \frac{1}{B} \sum_{i=0}^{B-1} L^{(i)}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

$$L^{(i)}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) = -[\mathbf{y}^{(i)} \log(\hat{\mathbf{y}}^{(i)}) + (1 - \mathbf{y}^{(i)}) \log(1 - \hat{\mathbf{y}}^{(i)})]$$

$$\frac{\partial L}{\partial \hat{\mathbf{y}}} = \left[\frac{\partial L}{\partial \hat{\mathbf{y}}^{(i)}} \right]_{i=0:B-1} \quad \frac{\partial L}{\partial \hat{\mathbf{y}}^{(i)}} = \frac{1}{B} \frac{\partial L^{(i)}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})}{\partial \hat{\mathbf{y}}^{(i)}} = \frac{1}{B} \left[\frac{\mathbf{y}^{(i)}}{\hat{\mathbf{y}}^{(i)}} - \frac{1 - \mathbf{y}^{(i)}}{1 - \hat{\mathbf{y}}^{(i)}} \right]$$

Start the backprop

K=5, 5 class classification



$$\frac{\partial L}{\partial \mathbf{z}} = [\hat{y}_0 - y_0, \hat{y}_1 - y_1, \dots, \hat{y}_{k-1} - y_{k-1}]^T$$

We have all ingredients for MB-SGD

Initialize weights and bias

Random shuffle dataset

BatchSize = 32

for epoch in range(E):

select #BatchSize samples

1. Evaluate loss function (forward pass) at this **batch**

$$L = \frac{1}{B} \sum_{i=0}^{B-1} L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}(\mathbf{W}, \mathbf{b}))$$

2. Compute gradient $\frac{\partial L}{\partial \mathbf{w}^{[l]}}, \frac{\partial L}{\partial \mathbf{b}^{[l]}}$

3. Update parameter $\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{w}^{[l]}}$, $\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \frac{\partial L}{\partial \mathbf{b}^{[l]}}$

Forward pass

Run data through model to compute output

Backward pass or backprop

Starting from loss, compute gradient to all parameters

Update

Use gradient information to update parameters

