



# CSCE 585: Machine Learning Systems

## Lecture 9: Backpropagation and Automatic Differentiation

Pooyan Jamshidi





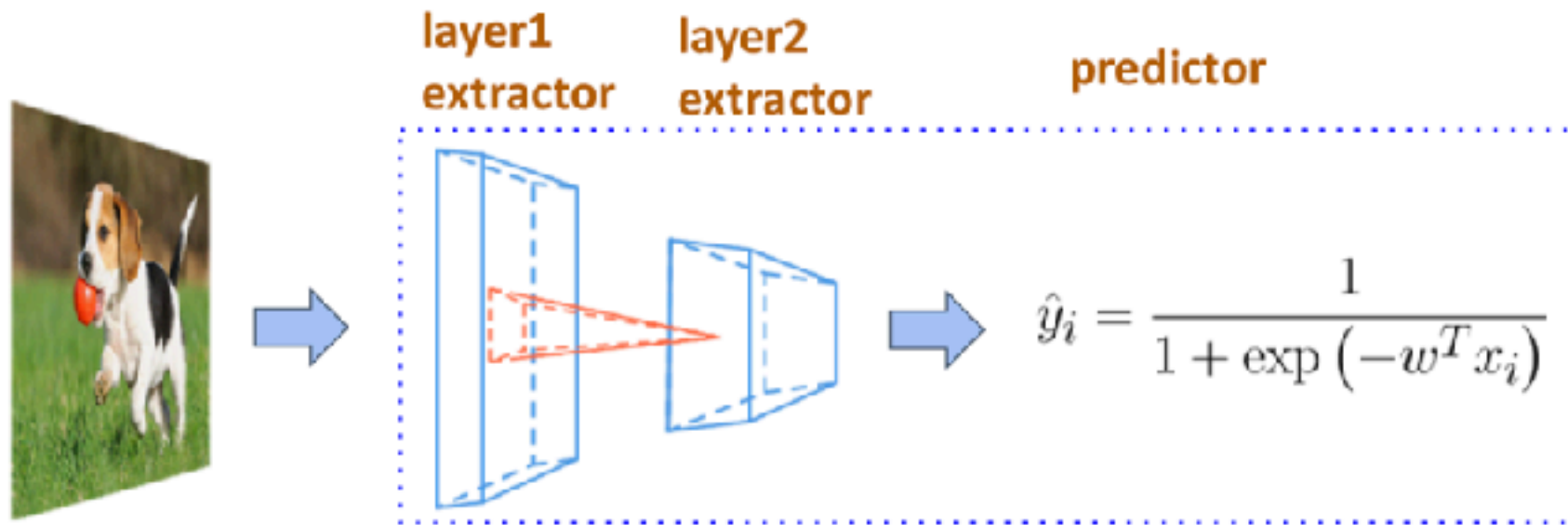
# Motivation: Why Gradients Matter

- Training = optimization.
- Optimization = gradients.
- Gradients tell us how to change parameters to reduce loss.

- Equation:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t)$$

# Model Training Overview



**Objective**

$$L(w) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \lambda \|w\|^2$$

**Training**

$$w \leftarrow w - \eta \nabla_w L(w)$$

# The Three Ways to Differentiate

Method	Type	Pros	Cons
Numerical	Approximate	Simple	Inaccurate, slow

# The Three Ways to Differentiate

Method	Type	Pros	Cons
Numerical	Approximate	Simple	Inaccurate, slow
Symbolic	Exact	Closed form	Expression swell

# The Three Ways to Differentiate

Method	Type	Pros	Cons
Numerical	Approximate	Simple	Inaccurate, slow
Symbolic	Exact	Closed form	Expression swell
Automatic	Exact	Efficient, flexible	Requires computational graph

# Numerical Differentiation

- We can approximate the gradient using

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$$

$$f(W, x) = W \cdot x$$
$$[-0.8 \quad 0.3] \cdot \begin{bmatrix} 0.5 \\ -0.2 \end{bmatrix}$$

# Numerical Differentiation

- We can approximate the gradient using

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$$

$$f(W, \mathbf{x}) = W \cdot \mathbf{x} \\ [-0.8 \quad 0.3] \cdot \begin{bmatrix} 0.5 \\ -0.2 \end{bmatrix}$$

$$f(W, \mathbf{x}) = W \cdot \mathbf{x} \\ [-0.8 + \varepsilon \quad 0.3] \cdot \begin{bmatrix} 0.5 \\ -0.2 \end{bmatrix}$$



# Numerical Differentiation

- We can approximate the gradient using

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$$

- Reduce the truncation error by using center difference

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h}$$

✗ Bad: rounding error, and slow to compute

✓ A powerful tool to check the correctness of implementation, usually use  $h = 1\text{e-}6$ .

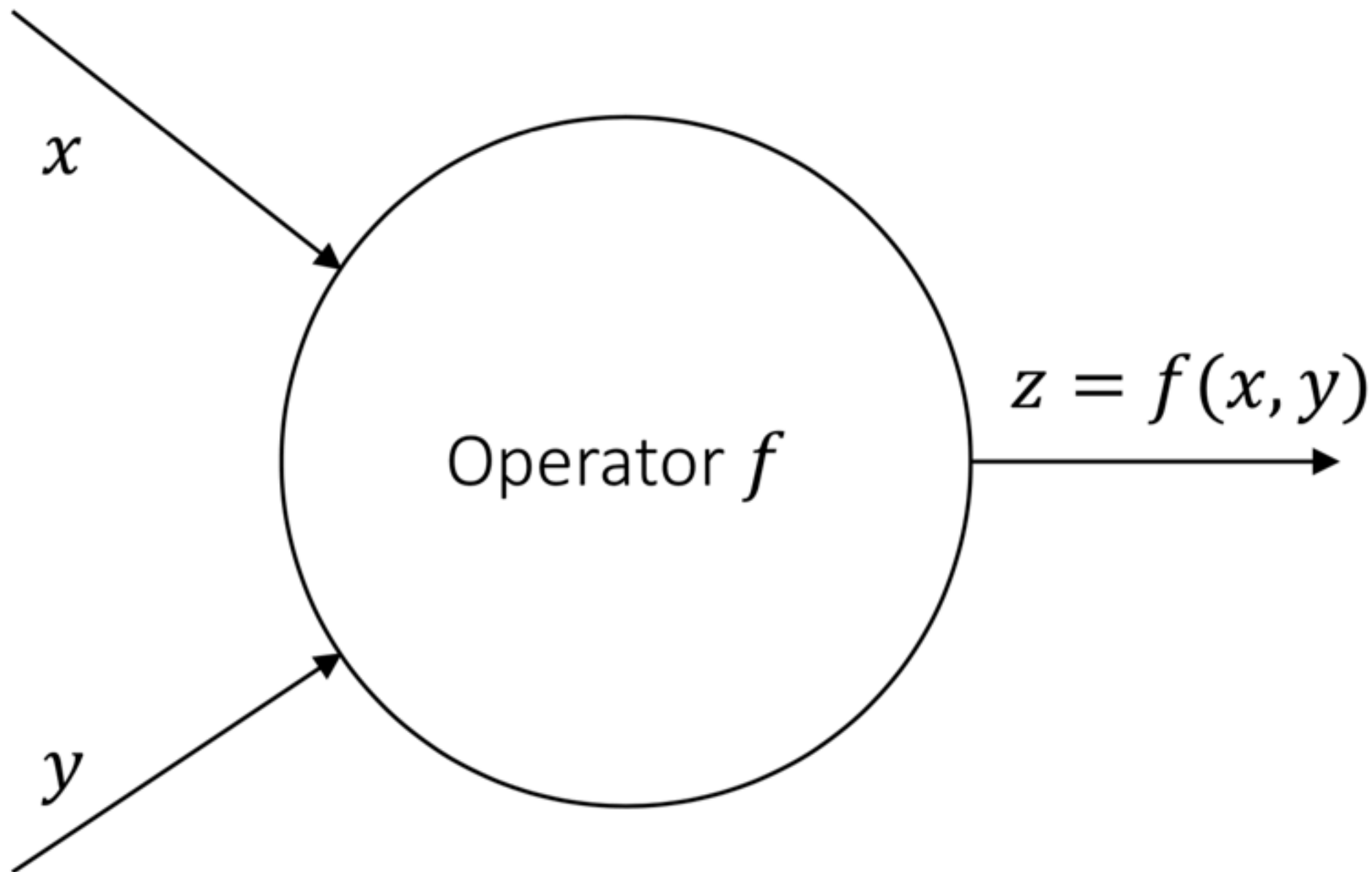
# Symbolic Differentiation

- Input formulae is a symbolic expression tree (computation graph).
- Implement differentiation rules, e.g., sum rule, product rule, chain rule

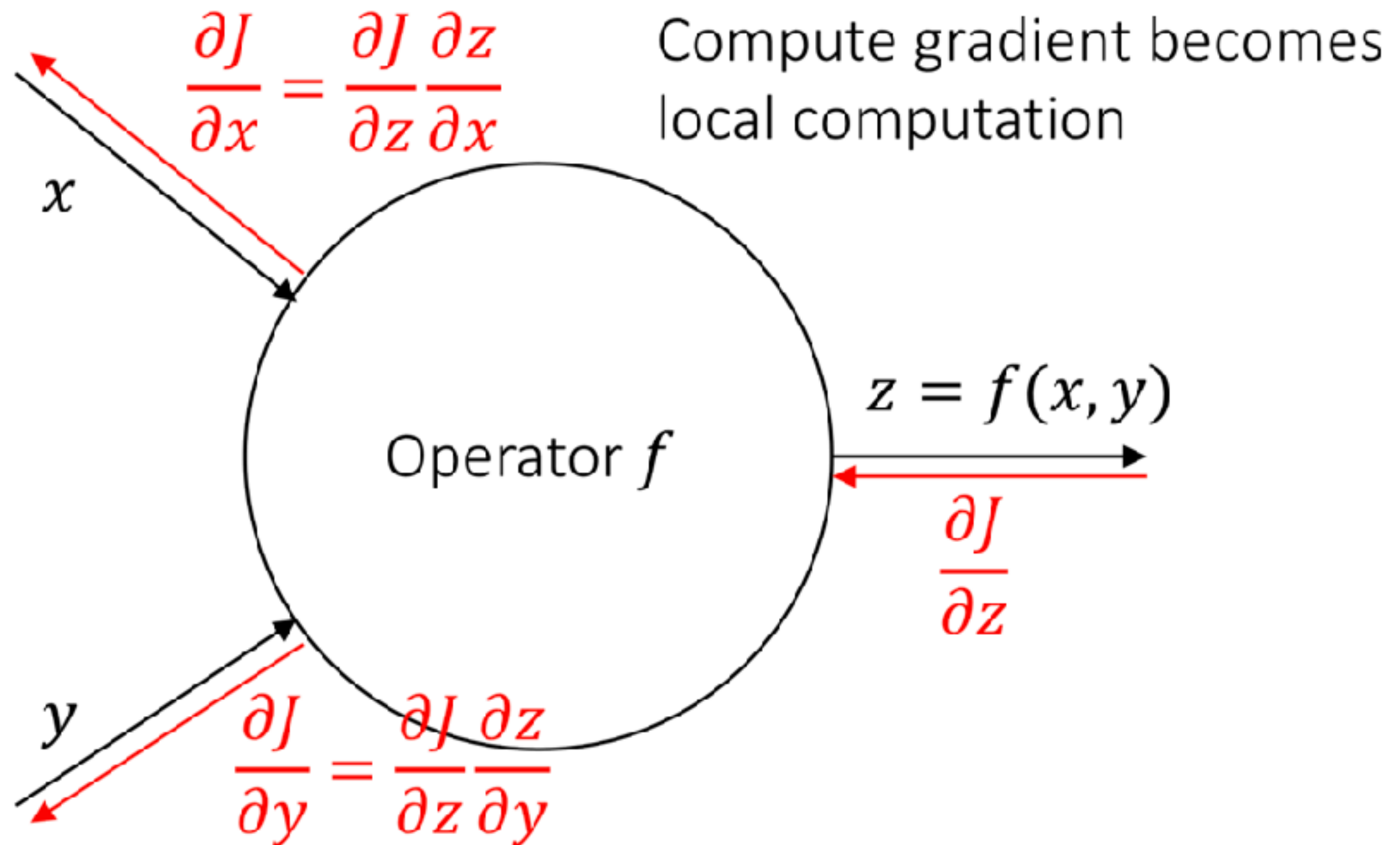
$$\frac{d(f + g)}{dx} = \frac{df}{dx} + \frac{dg}{dx} \quad \frac{d(fg)}{dx} = \frac{df}{dx}g + f\frac{dg}{dx} \quad \frac{d(h(x))}{dx} = \frac{df(g(x))}{dx} \cdot \frac{dg(x)}{dx}$$

- ✗ For complicated functions, the resultant expression can be exponentially large.
- ✗ Wasteful to keep around intermediate symbolic expressions if we only need a numeric value of the gradient in the end
- ✗ Prone to error

# Backpropagation



# Backpropagation

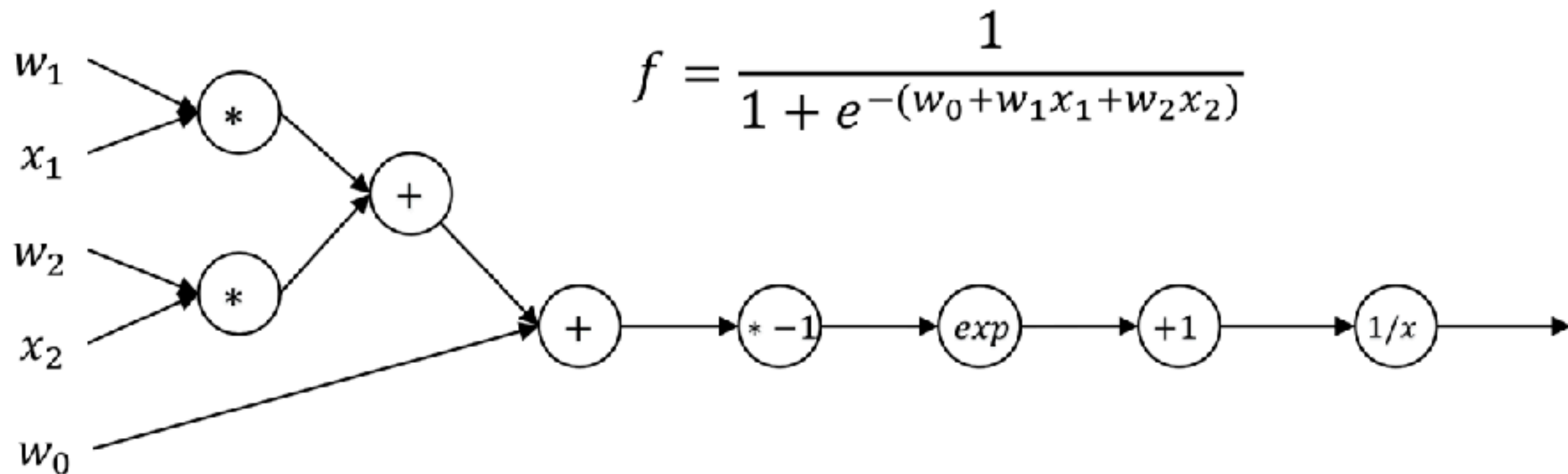


# Backpropagation simple example

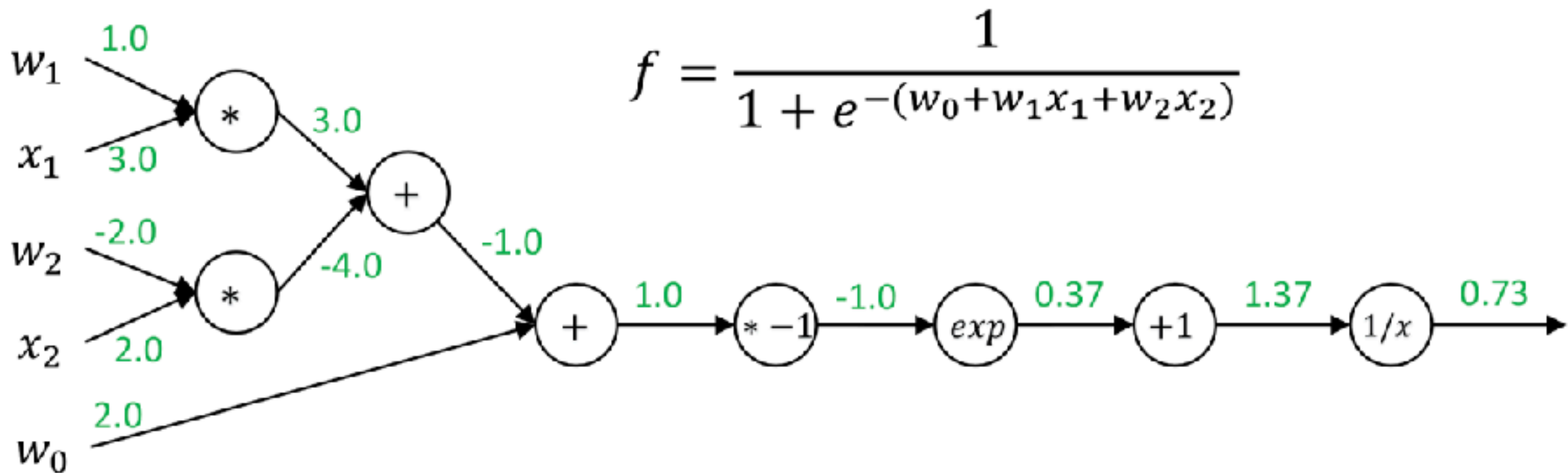
$$f = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$



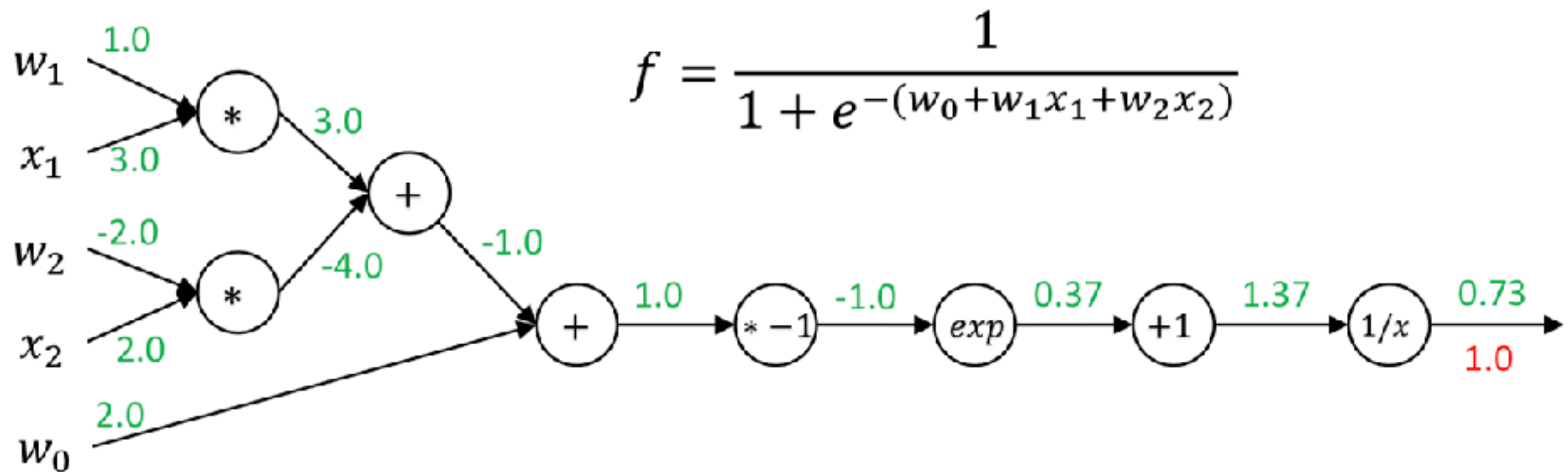
# Backpropagation simple example



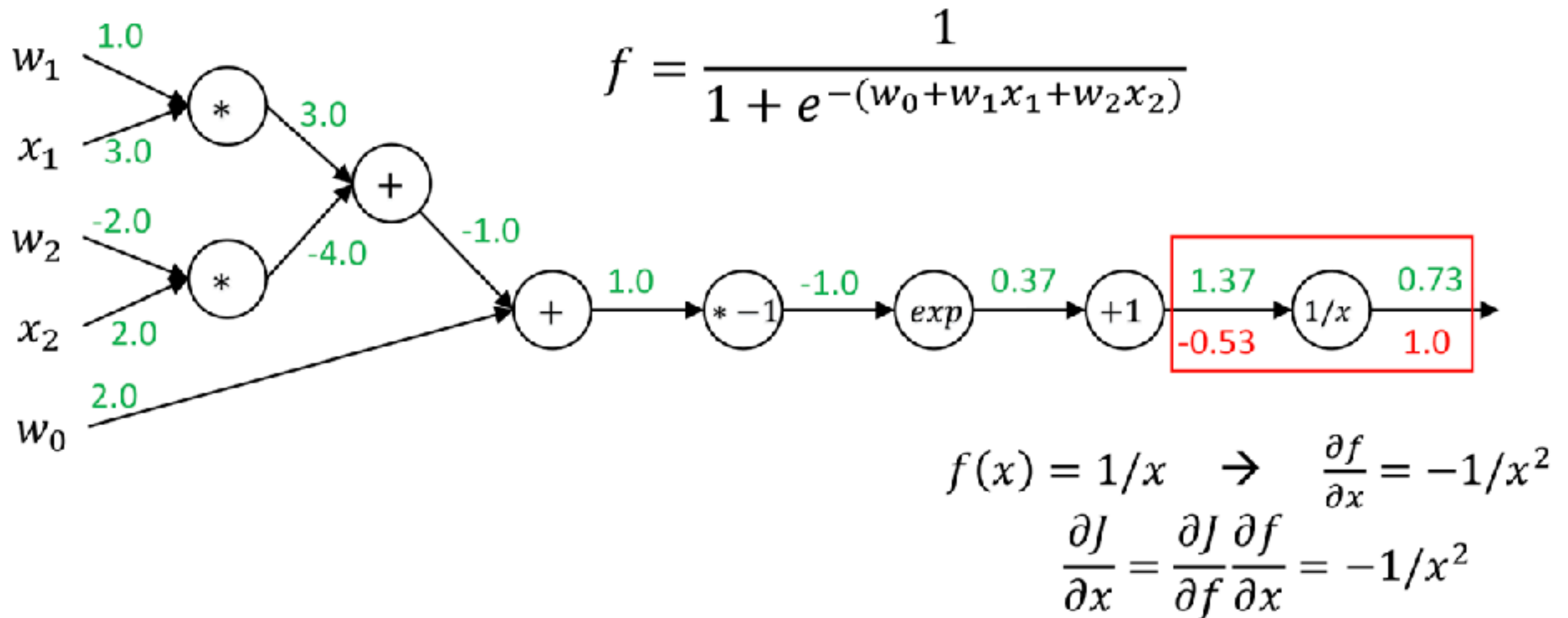
# Backpropagation simple example



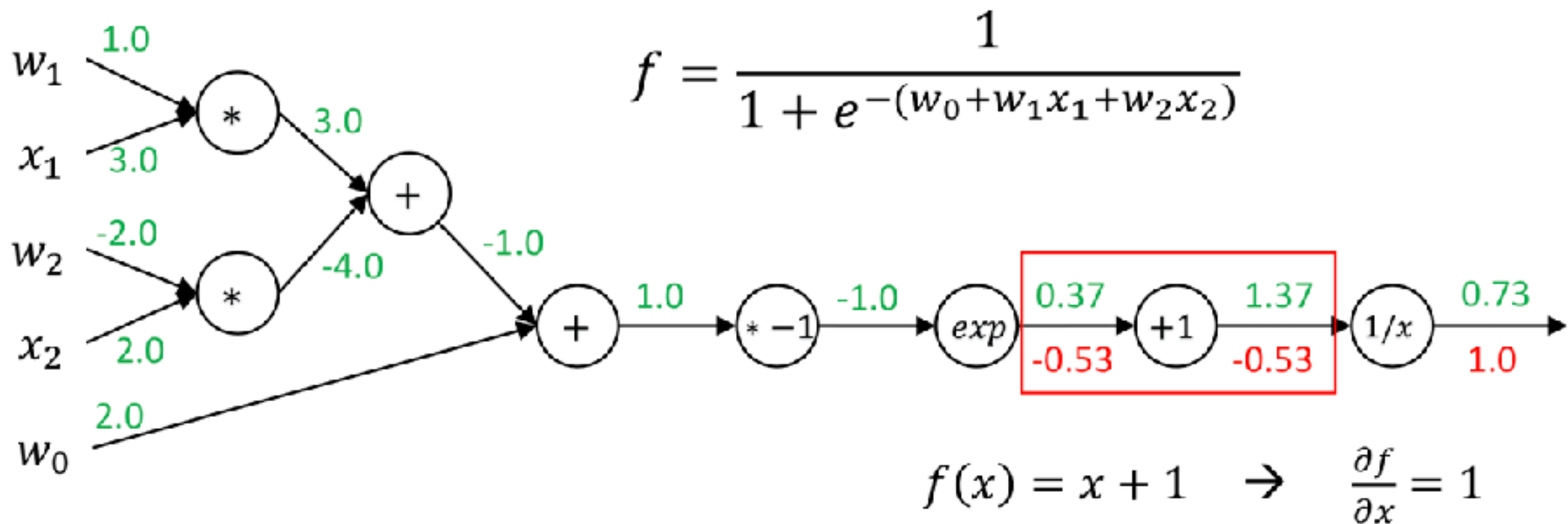
# Backpropagation simple example



# Backpropagation simple example

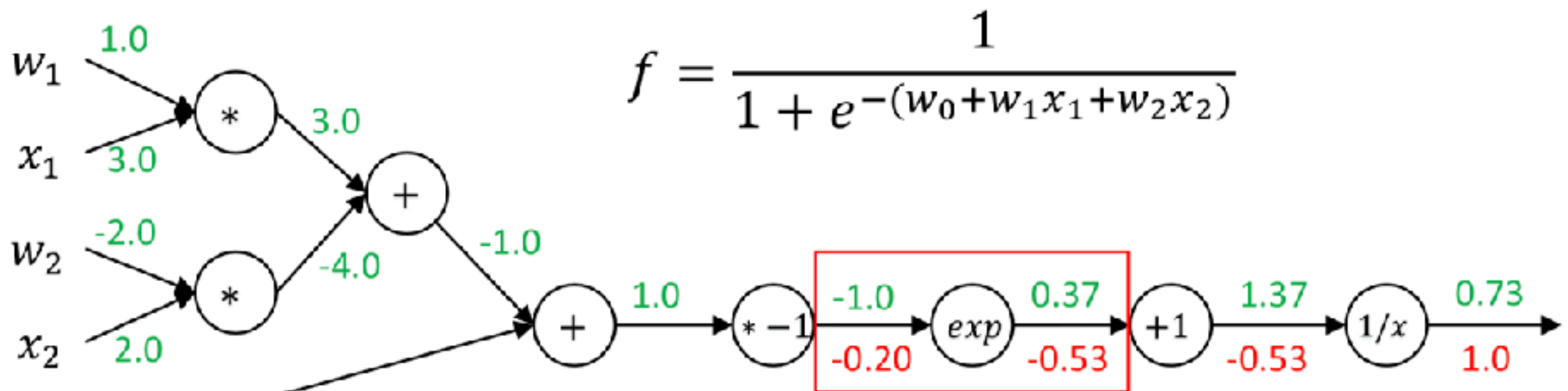


# Backpropagation simple example





# Backpropagation simple example

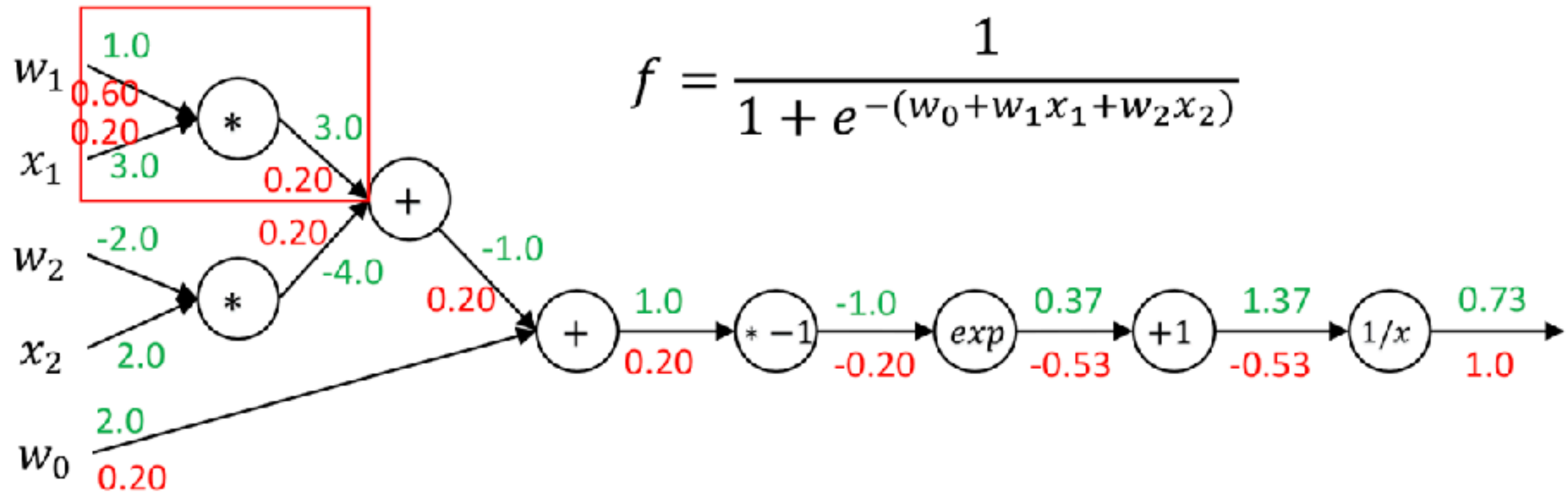


$$f = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial f} \frac{\partial f}{\partial x} = \frac{\partial J}{\partial f} \cdot e^x$$

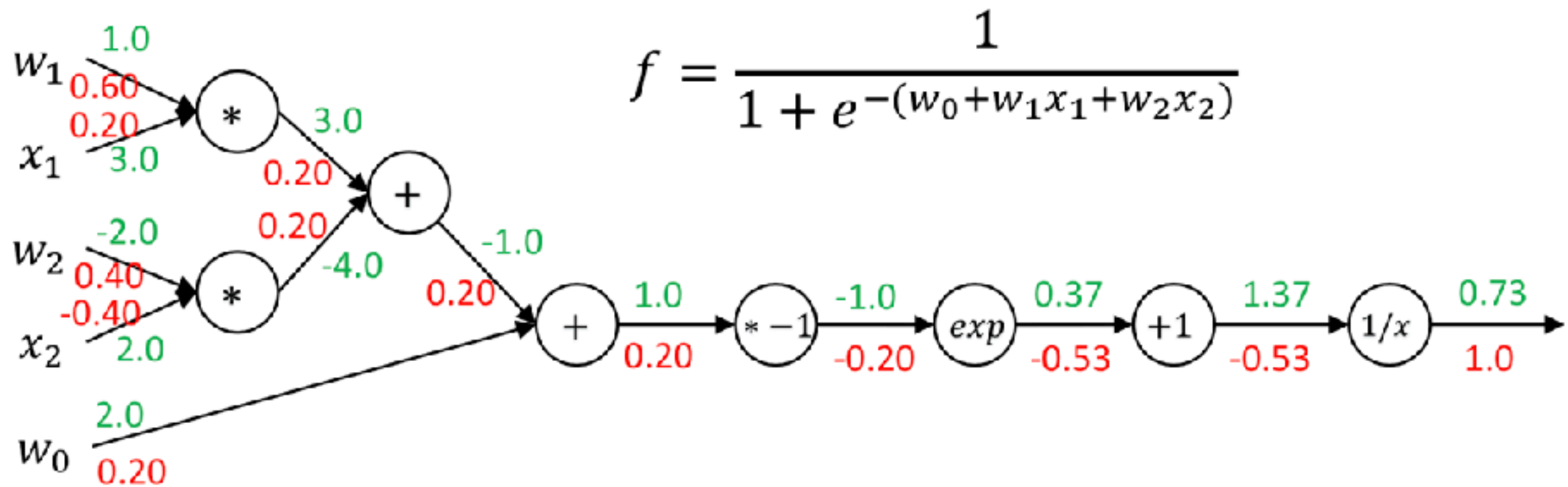
# Backpropagation simple example



$$f = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

$$f(x, w) = xw \rightarrow \frac{\partial f}{\partial x} = w, \frac{\partial f}{\partial w} = x$$

# Backpropagation simple example



$$f(x, y) = \frac{x + \sigma(y)}{\sigma(x) + (x + y)^2}$$

# Forward Pass

```
x = 3 # example values
y = -4

# forward pass
sigy = 1.0 / (1 + math.exp(-y)) # sigmoid in numerator    #(1)
num = x + sigy # numerator                                     #(2)
sigx = 1.0 / (1 + math.exp(-x)) # sigmoid in denominator  #(3)
xpy = x + y                                                  #(4)
xpysqr = xpy**2                                              #(5)
den = sigx + xpysqr # denominator                            #(6)
invden = 1.0 / den                                           #(7)
f = num * invden # done!                                     #(8)
```



# Backward Pass

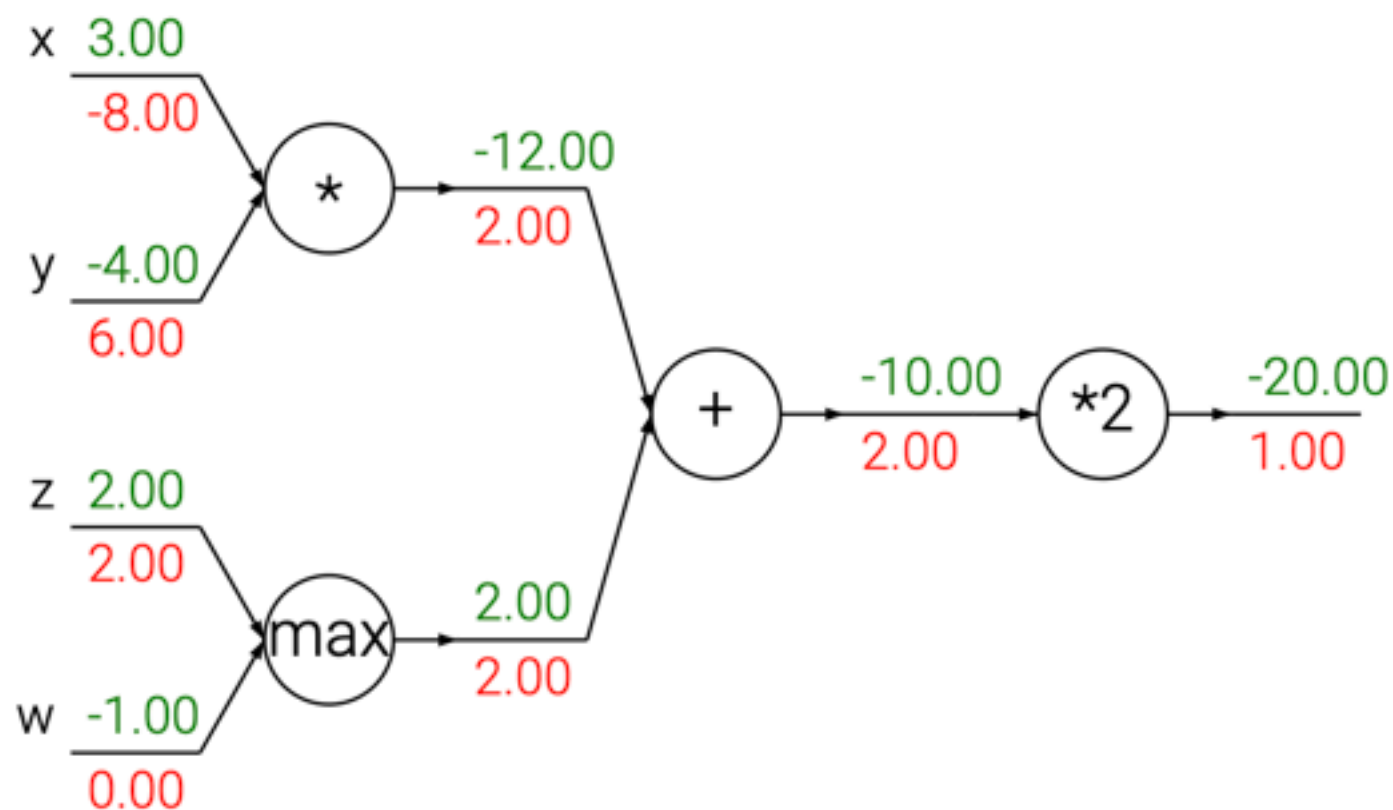
```
# backprop f = num * invden
dnum = invden # gradient on numerator                #(8)
dinvden = num                                         #(8)
# backprop invden = 1.0 / den
dden = (-1.0 / (den**2)) * dinvden                   #(7)
# backprop den = sigx + xpysqr
dsigx = (1) * dden                                    #(6)
dxpysqr = (1) * dden                                  #(6)
# backprop xpysqr = xpy**2
dxdpy = (2 * xpy) * dxpysqr                          #(5)
# backprop xpy = x + y
dx = (1) * dxdpy                                       #(4)
dy = (1) * dxdpy                                       #(4)
# backprop sigx = 1.0 / (1 + math.exp(-x))
dx += ((1 - sigx) * sigx) * dsigx # Notice += !! See notes below #(3)
# backprop num = x + sigy
dx += (1) * dnum                                       #(2)
dsigy = (1) * dnum                                     #(2)
# backprop sigy = 1.0 / (1 + math.exp(-y))
dy += ((1 - sigy) * sigy) * dsigy                     #(1)
# done! phew
```

**Cache forward pass variables.** To compute the backward pass it is very helpful to have some of the variables that were used in the forward pass. In practice you want to structure your code so that you cache these variables, and so that they are available during backpropagation. If this is too difficult, it is possible (but wasteful) to recompute them.

**Gradients add up at forks.** The forward expression involves the variables  $x, y$  multiple times, so when we perform backpropagation we must be careful to use `+=` instead of `=` to accumulate the gradient on these variables (otherwise we would overwrite it). This follows the *multivariable chain rule* in Calculus, which states that if a variable branches out to different parts of the circuit, then the gradients that flow back to it will add.

## Patterns in backward flow

It is interesting to note that in many cases the backward-flowing gradient can be interpreted on an intuitive level. For example, the three most commonly used gates in neural networks (*add,mul,max*), all have very simple interpretations in terms of how they act during backpropagation. Consider this example circuit:



An example circuit demonstrating the intuition behind the operations that backpropagation performs during the backward pass in order to compute the gradients on the inputs. Sum operation distributes gradients equally to all its inputs. Max operation routes the gradient to the higher input. Multiply gate takes the input activations, swaps them and multiplies by its gradient.

**Any problem? Can we  
do better?**

# Problems of backpropagation

- You always need to keep intermediate data in the memory during the forward pass in case it will be used in the backpropagation.
- Lack of flexibility, e.g., compute the gradient of gradient.

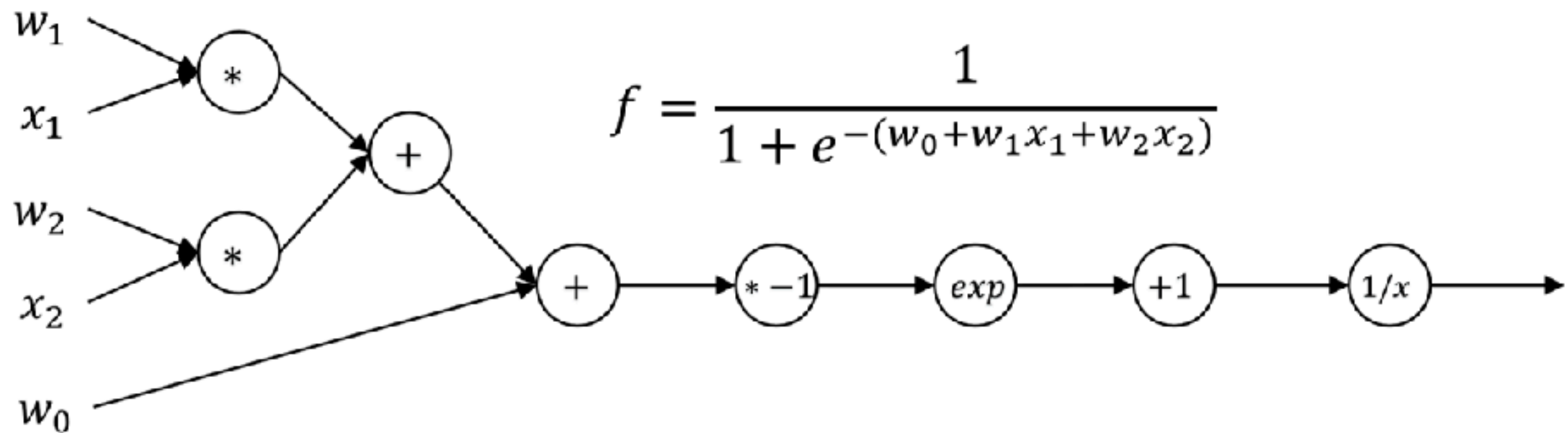


# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation

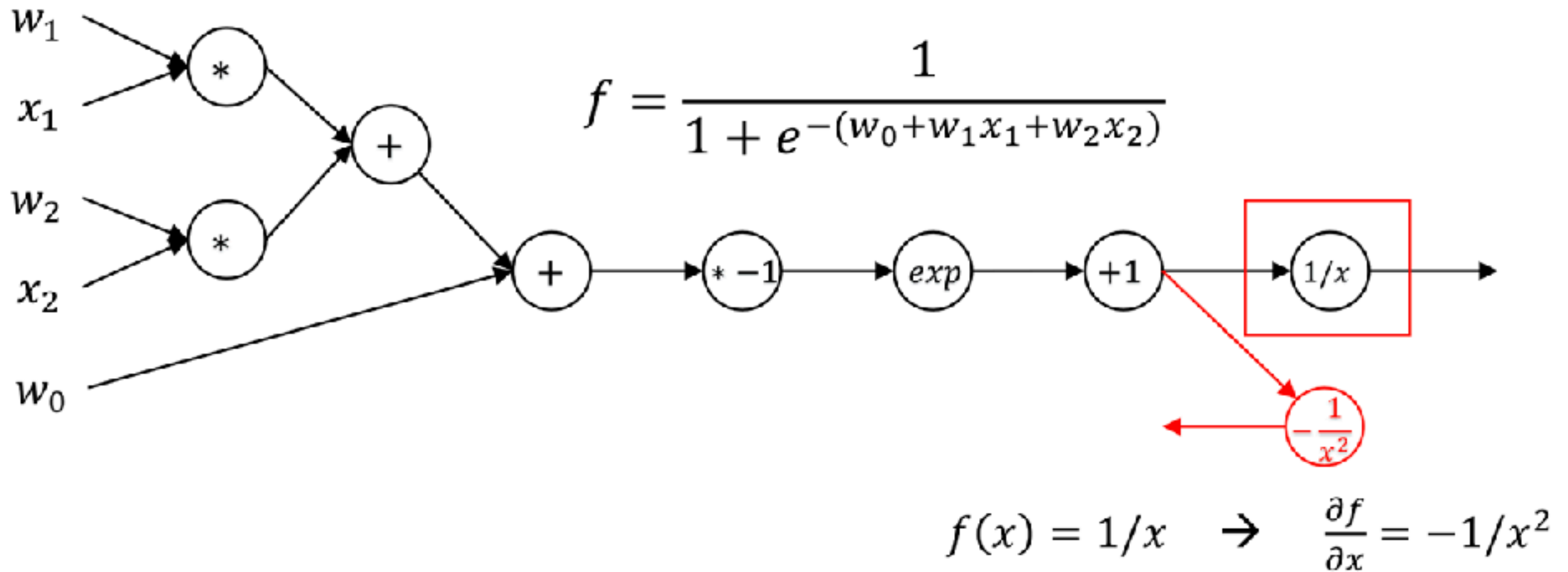
# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation



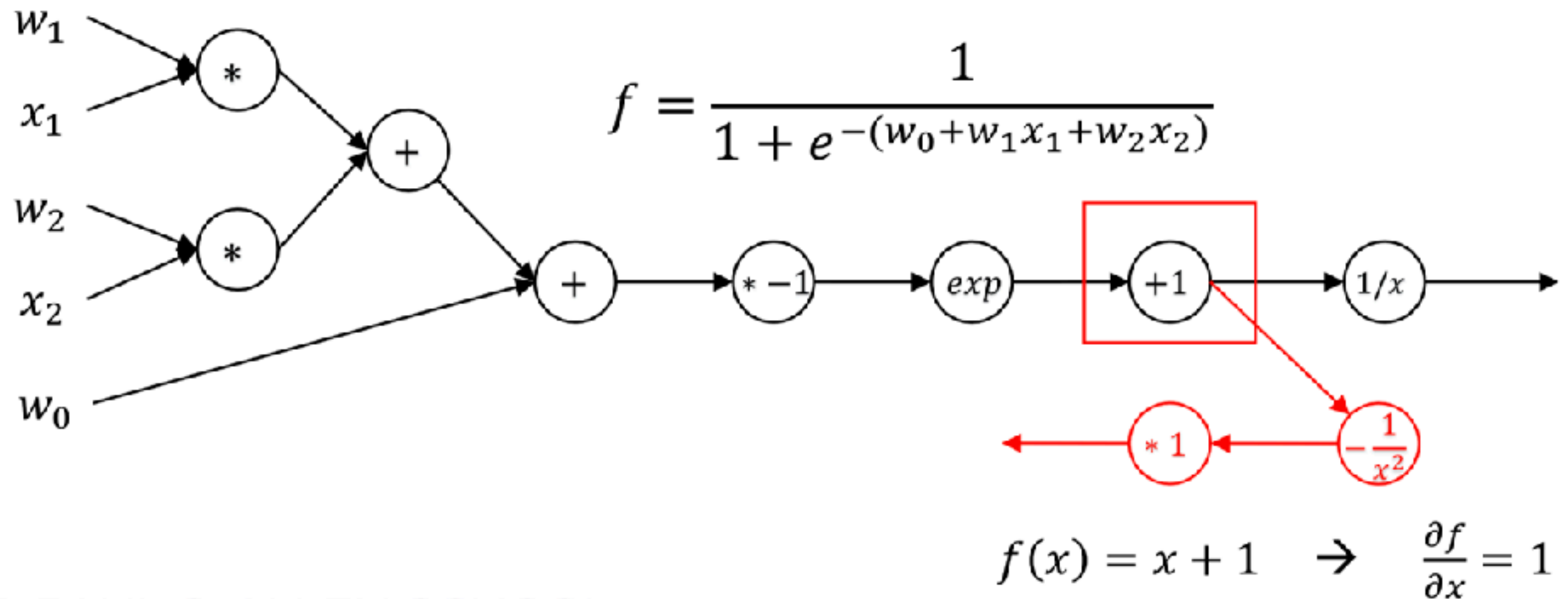
# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation



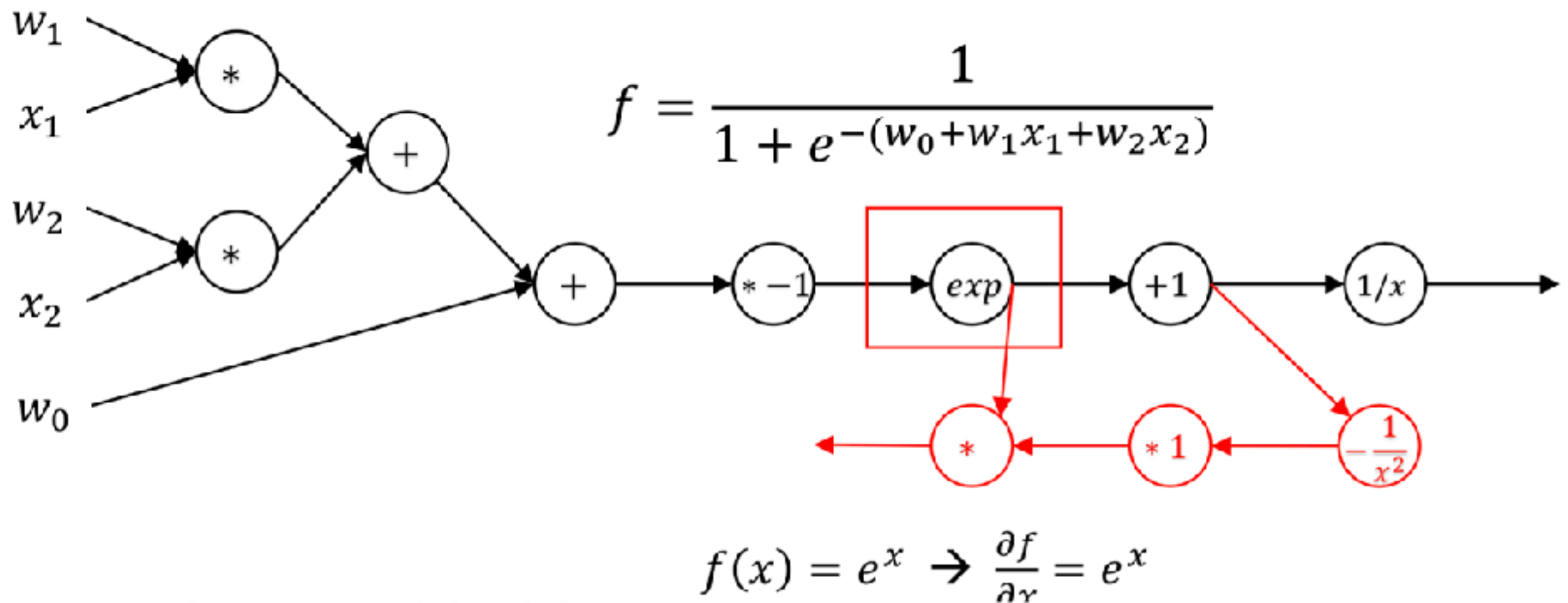
# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation



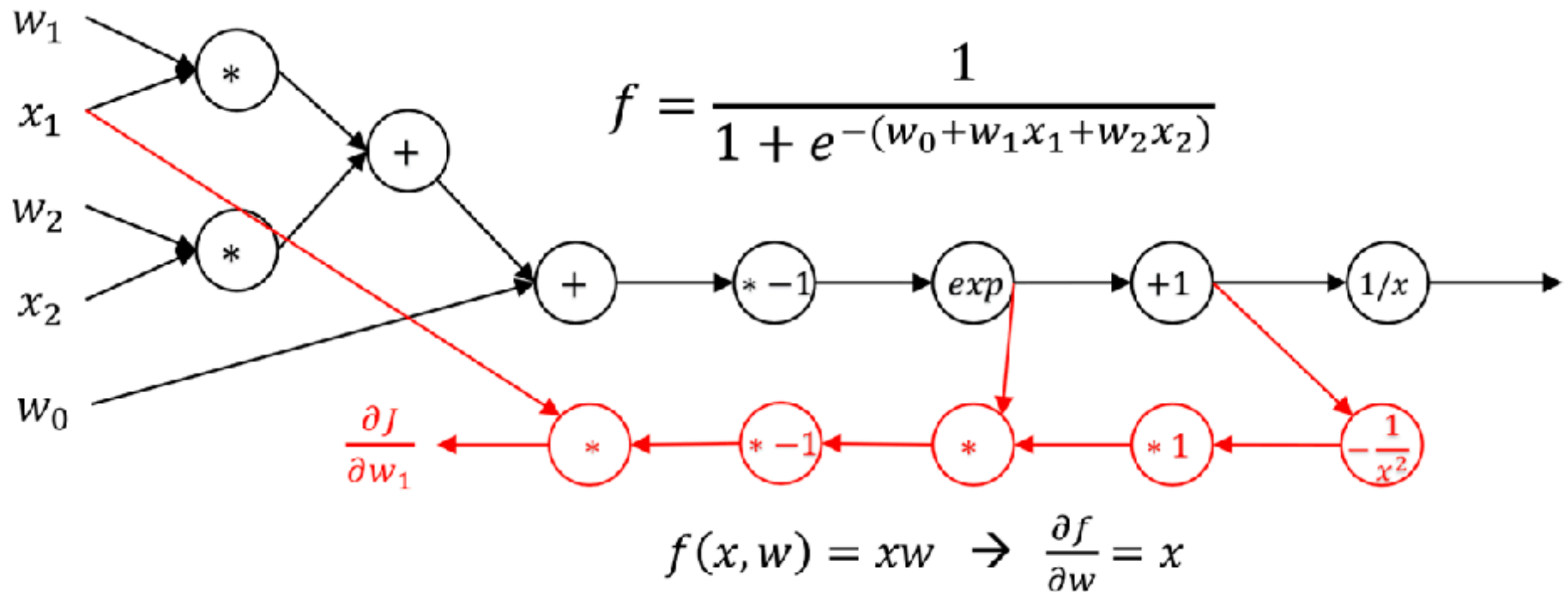
# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation



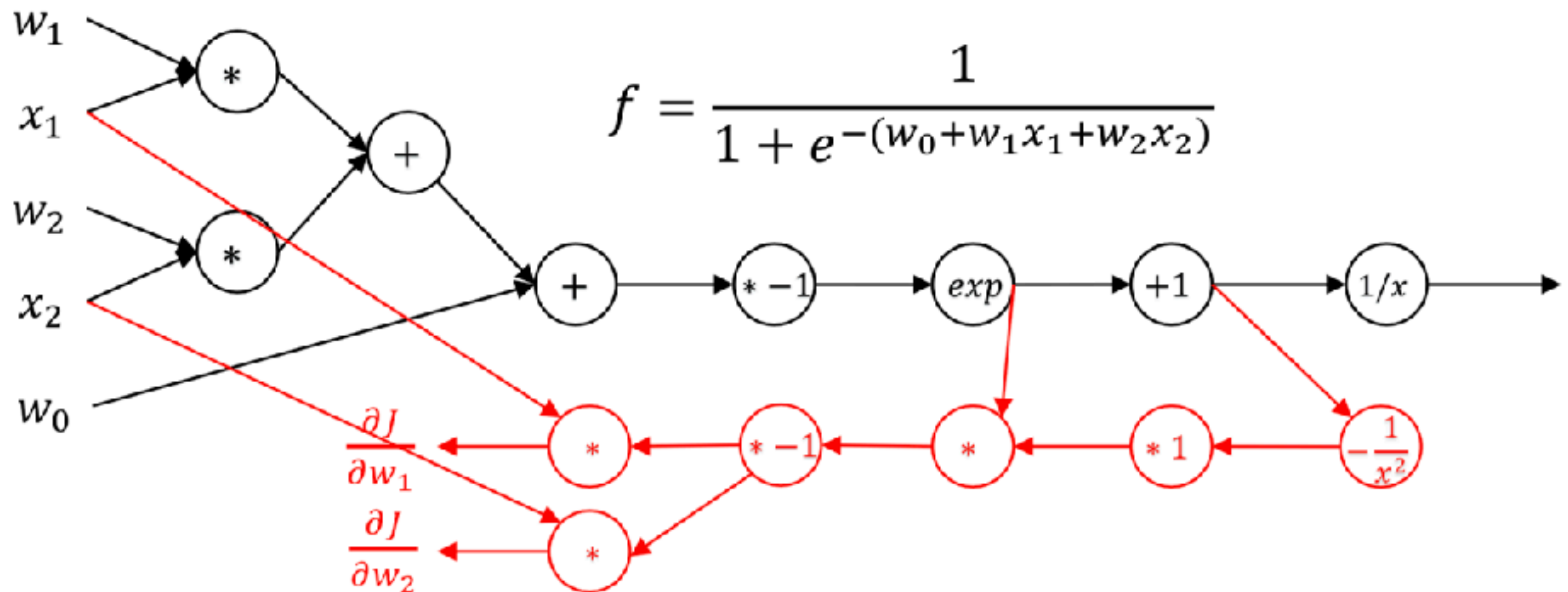
# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation

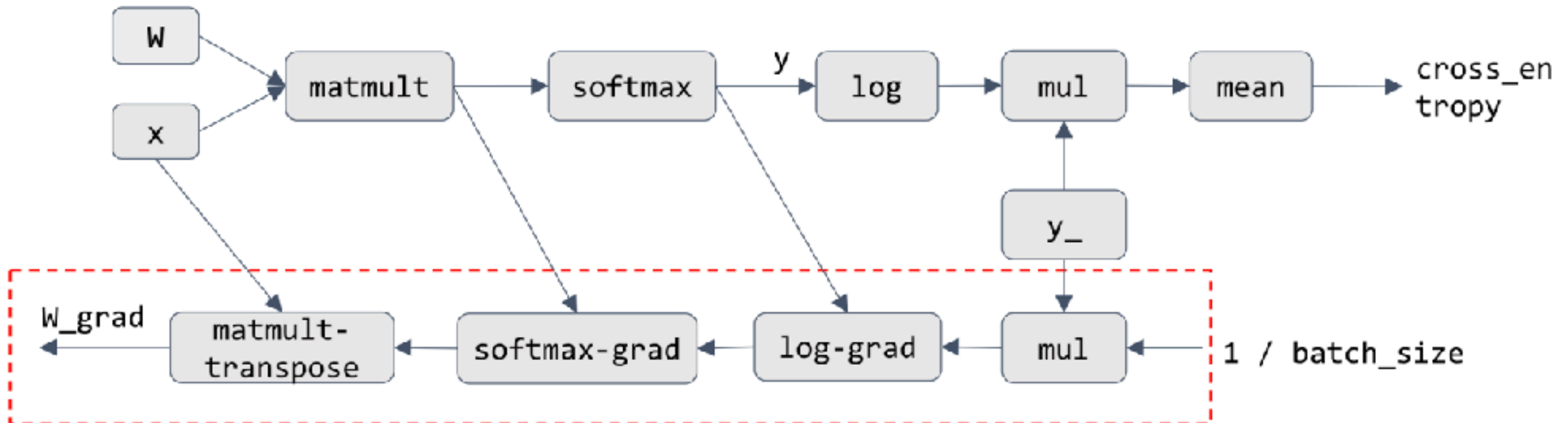


# Automatic Differentiation (autodiff)

- Create computation graph for gradient computation



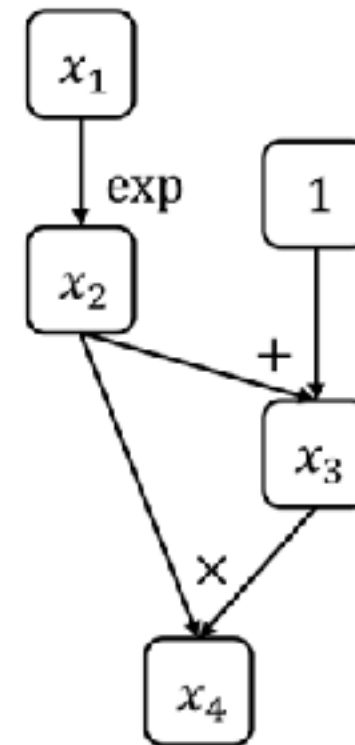
# AutoDiff Algorithm





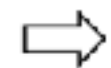
# AutoDiff Algorithm

```
⇒ def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```



# AutoDiff Algorithm

```
def gradient(out):
```



```
    node_to_grad[out] = 1
```

```
    nodes = get_node_list(out)
```

```
    for node in reverse_topo_order(nodes):
```

```
        grad ← sum partial adjoints from output edges
```

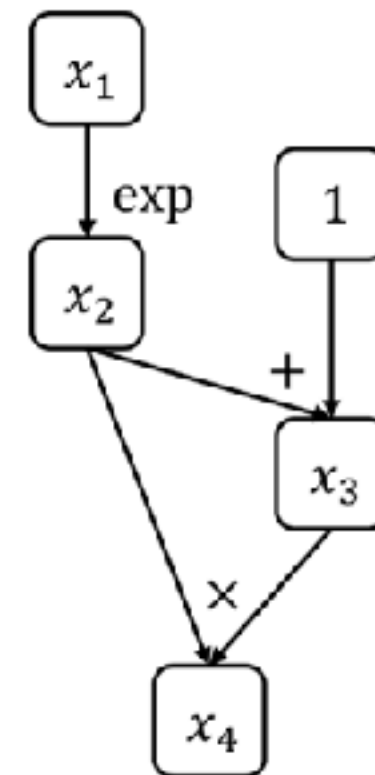
```
        input_grads ← node.op.gradient(input, grad) for  
input in node.inputs
```

```
        add input_grads to node_to_grad
```

```
    return node_to_grad
```

```
node_to_grad:
```

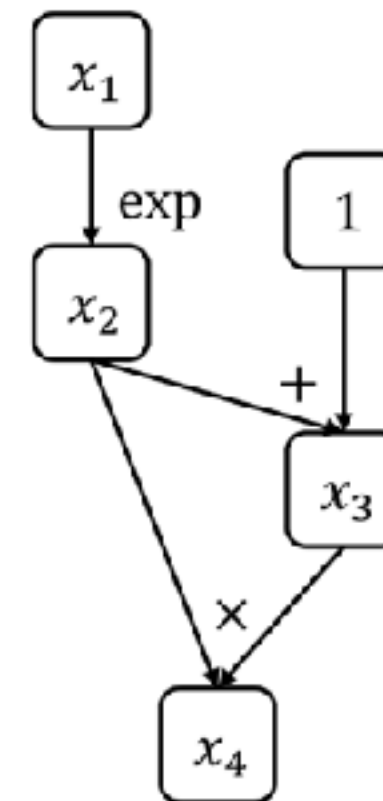
```
     $x_4$ :  $\overline{x_4}$ 
```



$\overline{x_4}$

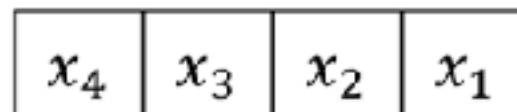
# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    ⇒ for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```



node\_to\_grad:

$x_4: \overline{x_4}$



$\overline{x_4}$

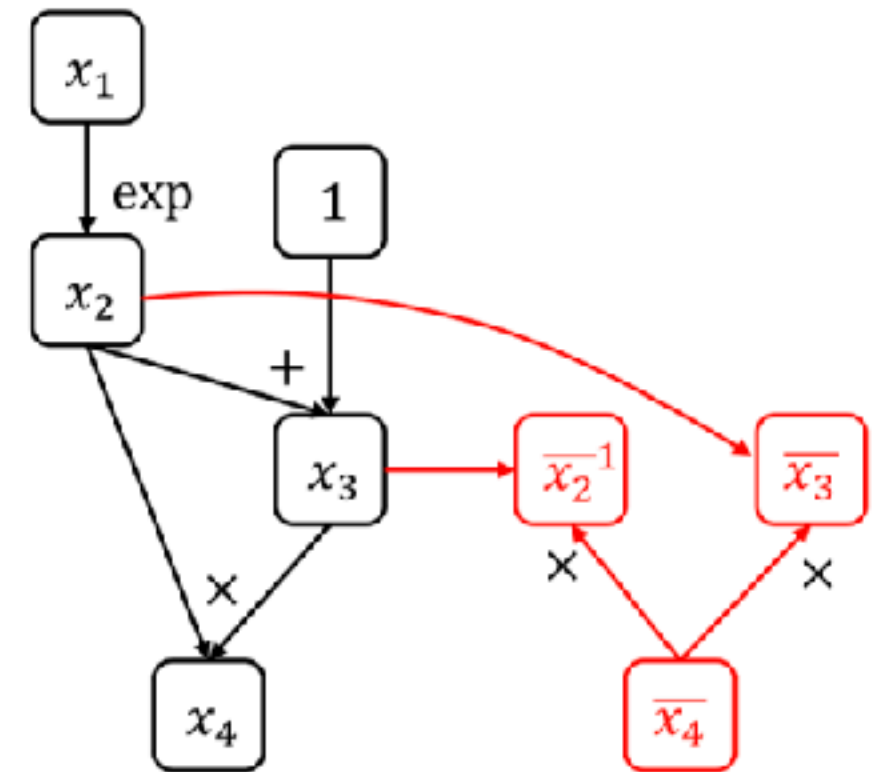
# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```

node\_to\_grad:

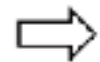
$x_4: \bar{x}_4$

$x_4$	$x_3$	$x_2$	$x_1$
-------	-------	-------	-------



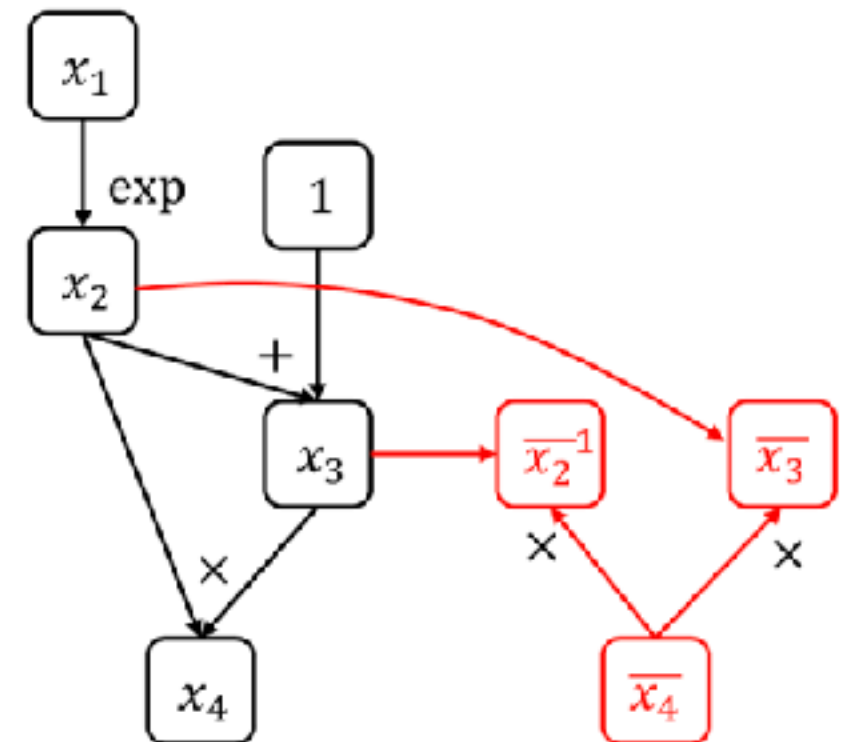
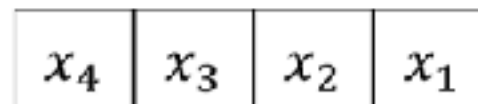
# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
    add input_grads to node_to_grad  
    return node_to_grad
```



node\_to\_grad:

$x_4: \overline{x_4}$   
 $x_3: \overline{x_3}$   
 $x_2: \overline{x_2}^{-1}$



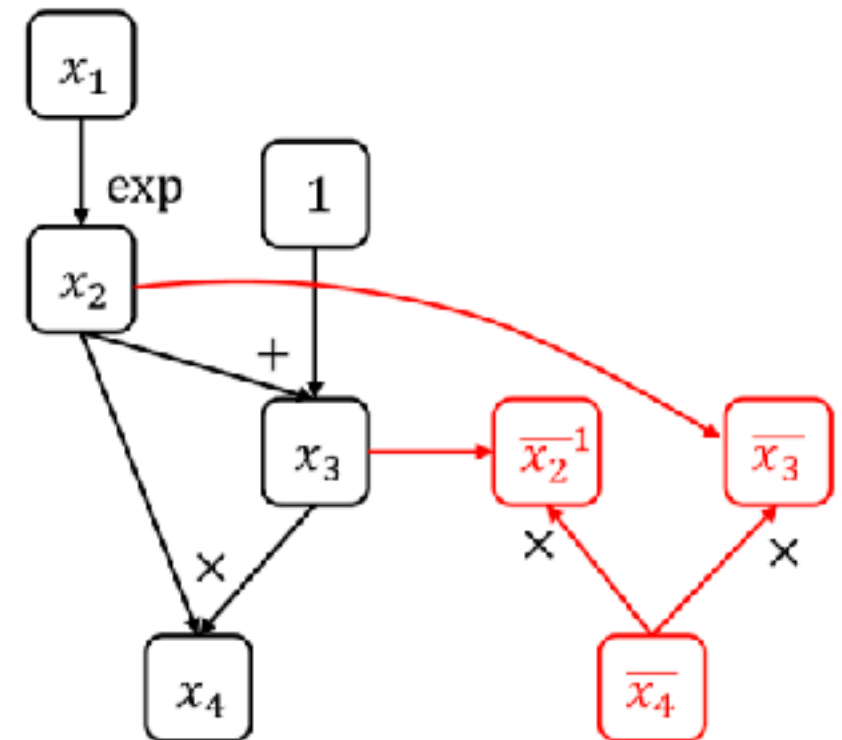
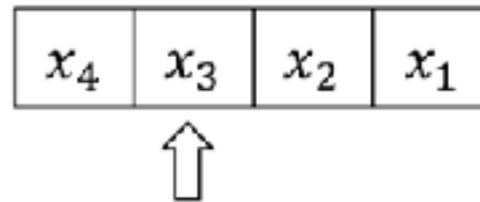
# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```



node\_to\_grad:

$x_4: \overline{x_4}$   
 $x_3: \overline{x_3}$   
 $x_2: \overline{x_2}^{-1}$



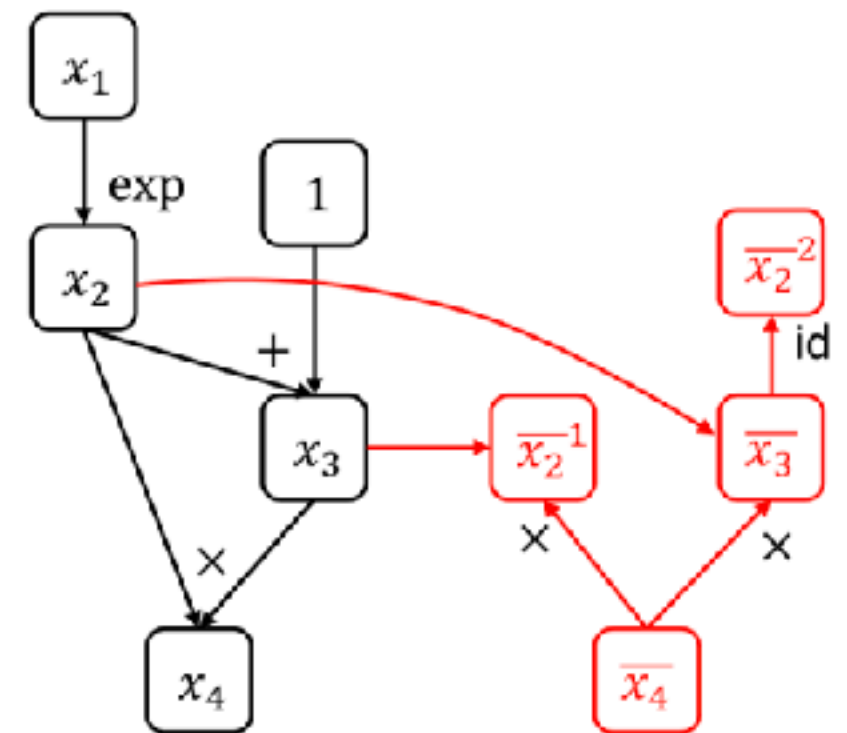
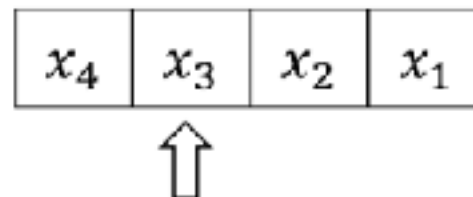
# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```



node\_to\_grad:

$x_4: \overline{x_4}$   
 $x_3: \overline{x_3}$   
 $x_2: \overline{x_2}^1$

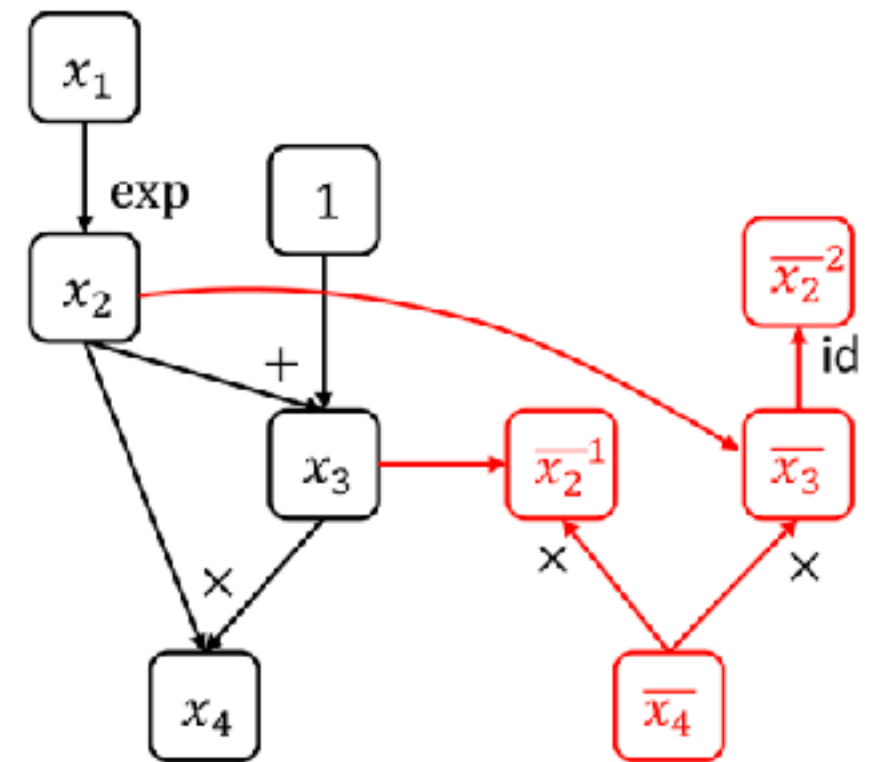
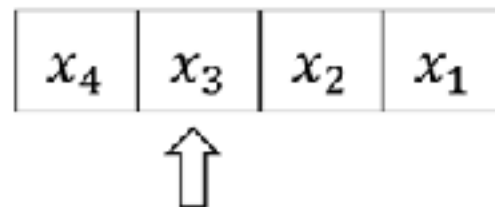


# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
    ⇒ add input_grads to node_to_grad  
    return node_to_grad
```

node\_to\_grad:

$x_4$ :  $\overline{x_4}$   
 $x_3$ :  $\overline{x_3}$   
 $x_2$ :  $\overline{x_2}^1, \overline{x_2}^2$



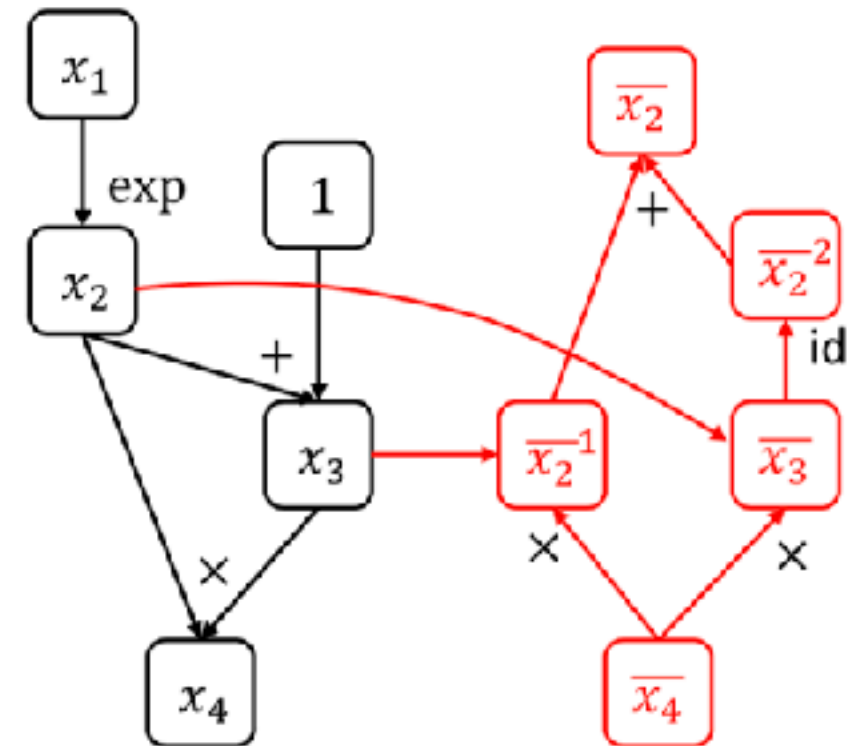
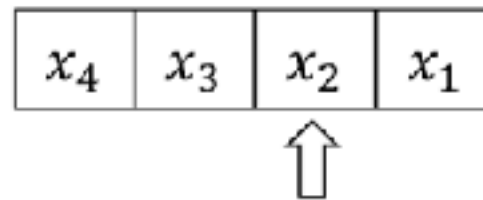


# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
input in node.inputs  
        add input_grads to node_to_grad  
    return node_to_grad
```

node\_to\_grad:

$x_4$ :  $\overline{x_4}$   
 $x_3$ :  $\overline{x_3}$   
 $x_2$ :  $\overline{x_2}^1$ ,  $\overline{x_2}^2$



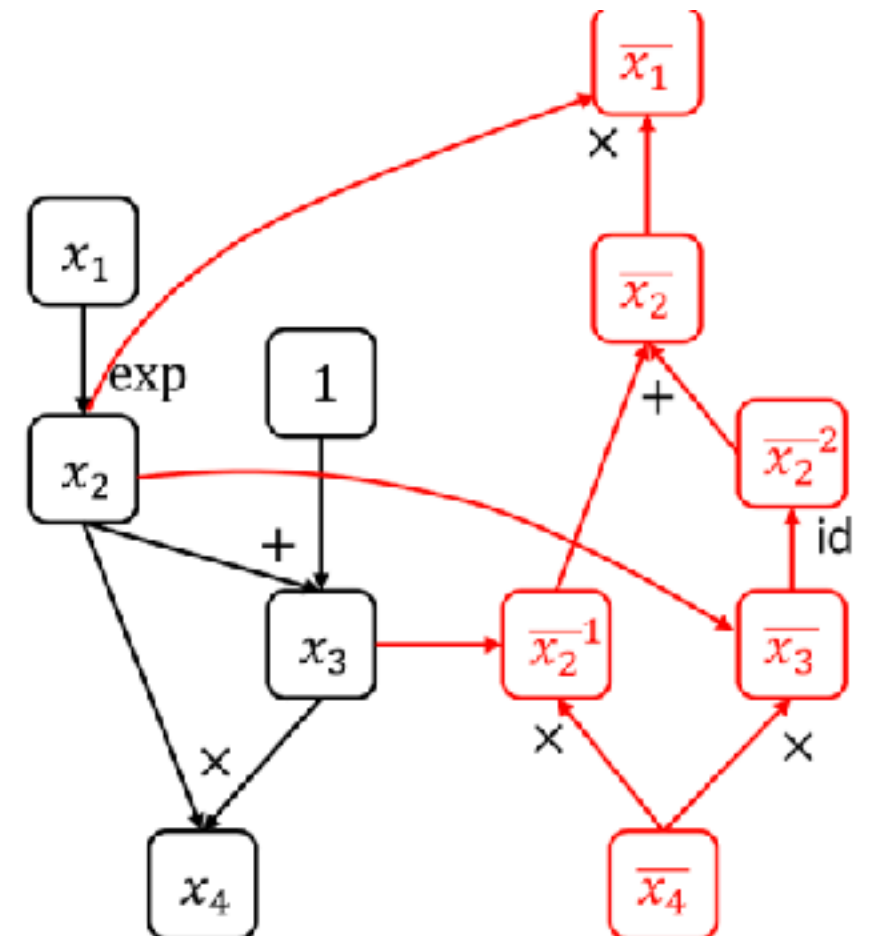
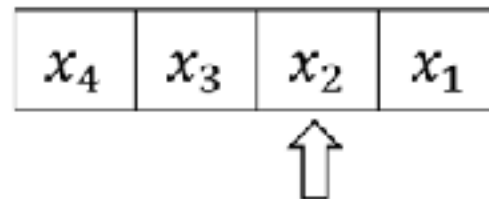
# AutoDiff Algorithm

```
def gradient(out):
    node_to_grad[out] = 1
    nodes = get_node_list(out)
    for node in reverse_topo_order(nodes):
        grad ← sum partial adjoints from output edges
        input_grads ← node.op.gradient(input, grad)
        for input in node.inputs:
            add input_grads to node_to_grad
    return node_to_grad
```



node\_to\_grad:

$x_4: \bar{x}_4$   
 $x_3: \bar{x}_3$   
 $x_2: \bar{x}_2^{-1}, \bar{x}_2^{-2}$

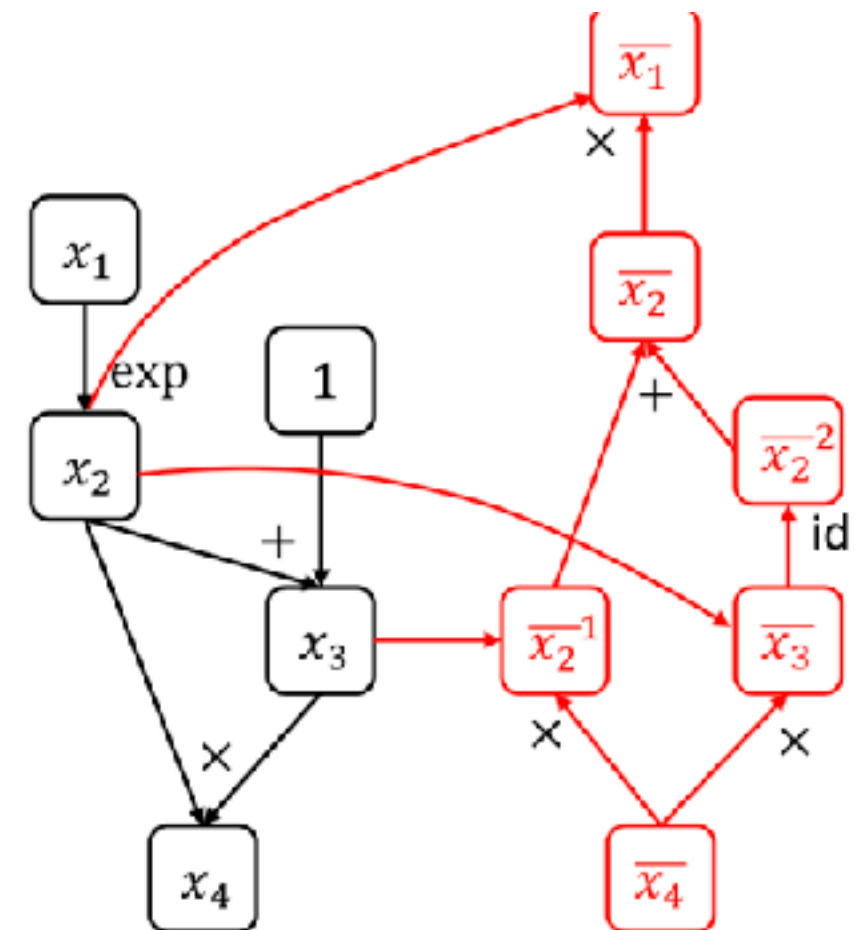
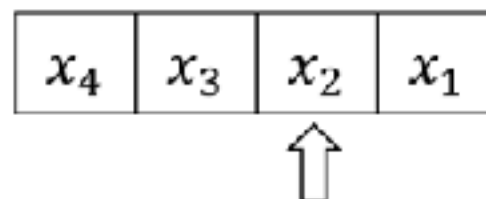


# AutoDiff Algorithm

```
def gradient(out):  
    node_to_grad[out] = 1  
    nodes = get_node_list(out)  
    for node in reverse_topo_order(nodes):  
        grad ← sum partial adjoints from output edges  
        input_grads ← node.op.gradient(input, grad) for  
        input in node.inputs  
    ⇒ add input_grads to node_to_grad  
    return node_to_grad
```

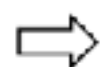
node\_to\_grad:

$x_4$ :  $\overline{x_4}$   
 $x_3$ :  $\overline{x_3}$   
 $x_2$ :  $\overline{x_2}^1, \overline{x_2}^2$   
 $x_1$ :  $\overline{x_1}$



# AutoDiff Algorithm

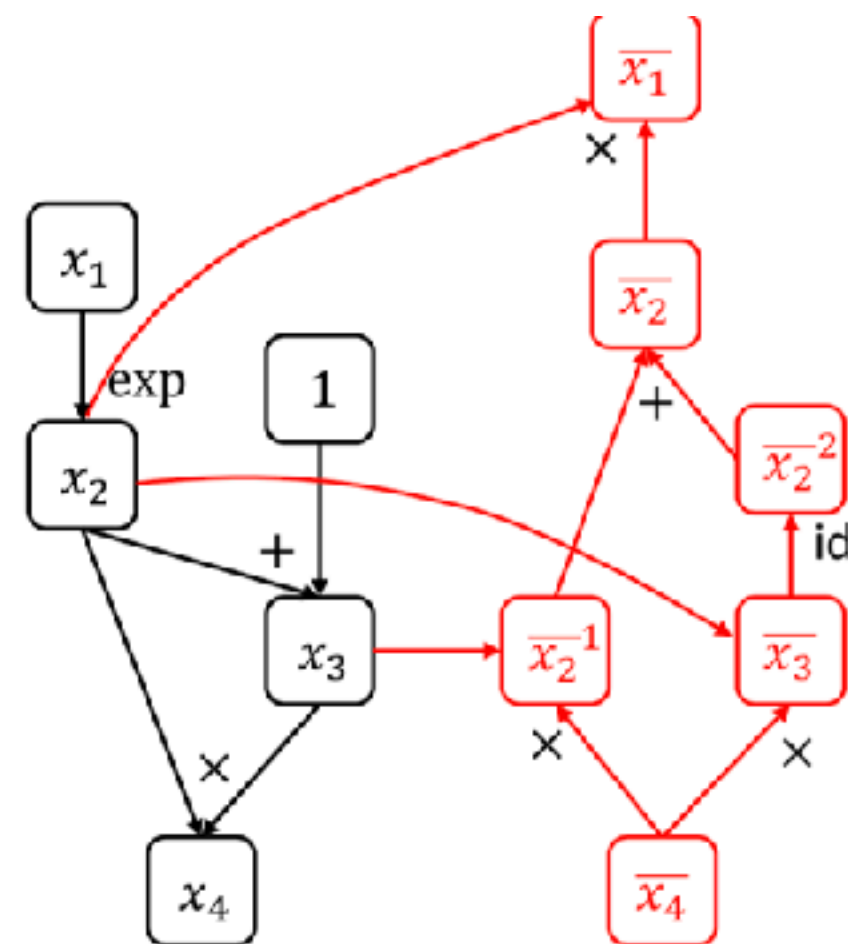
```
def gradient(out):
    node_to_grad[out] = 1
    nodes = get_node_list(out)
    for node in reverse_topo_order(nodes):
        grad ← sum partial adjoints from output edges
        input_grads ← node.op.gradient(input, grad) for
        input in node.inputs
        add input_grads to node_to_grad
    return node_to_grad
```



node\_to\_grad:

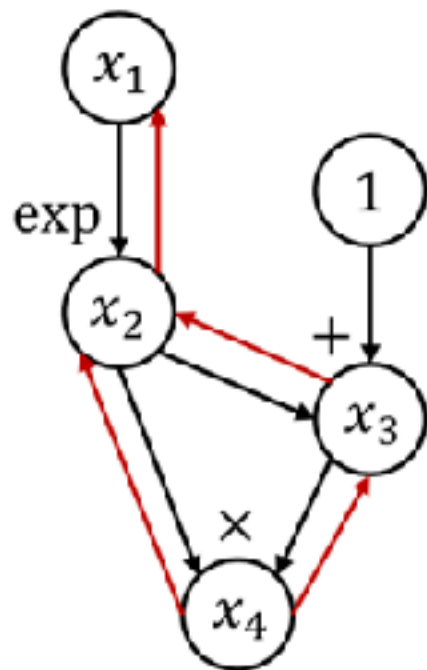
$x_4: \bar{x}_4$   
 $x_3: \bar{x}_3$   
 $x_2: \bar{x}_2^1, \bar{x}_2^2$   
 $x_1: \bar{x}_1$

$x_4$	$x_3$	$x_2$	$x_1$
-------	-------	-------	-------

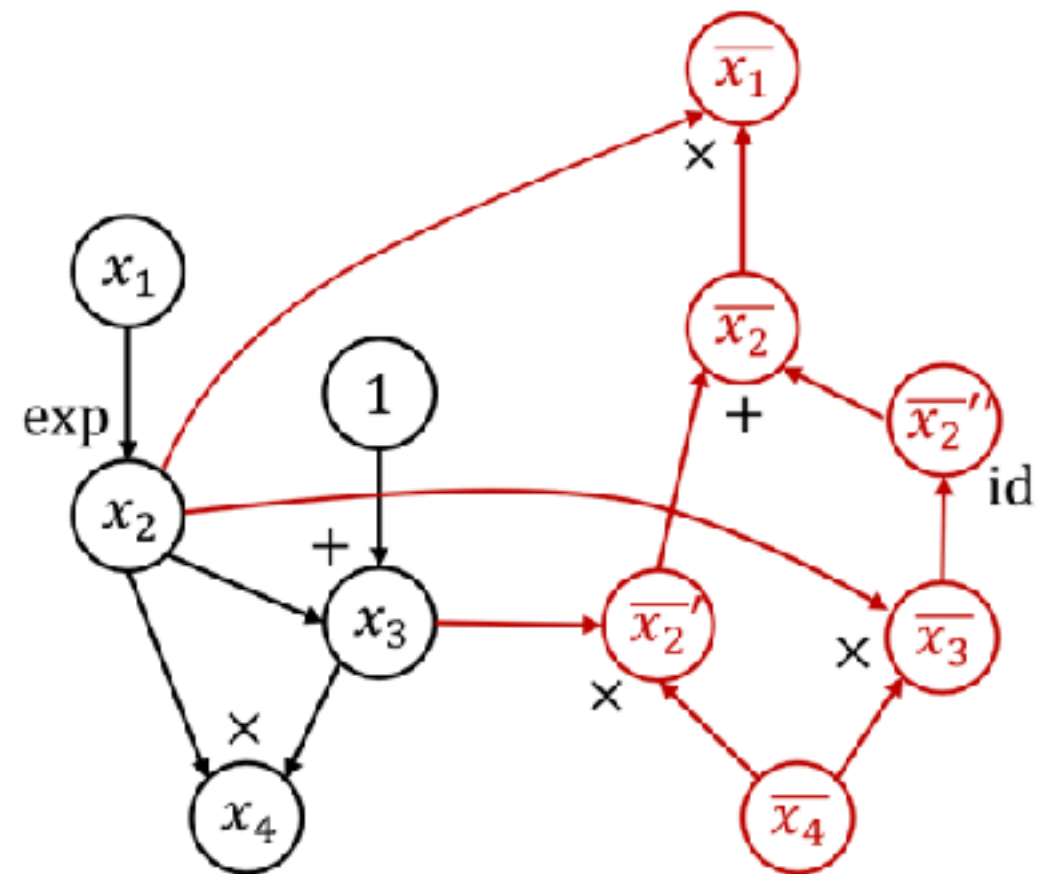


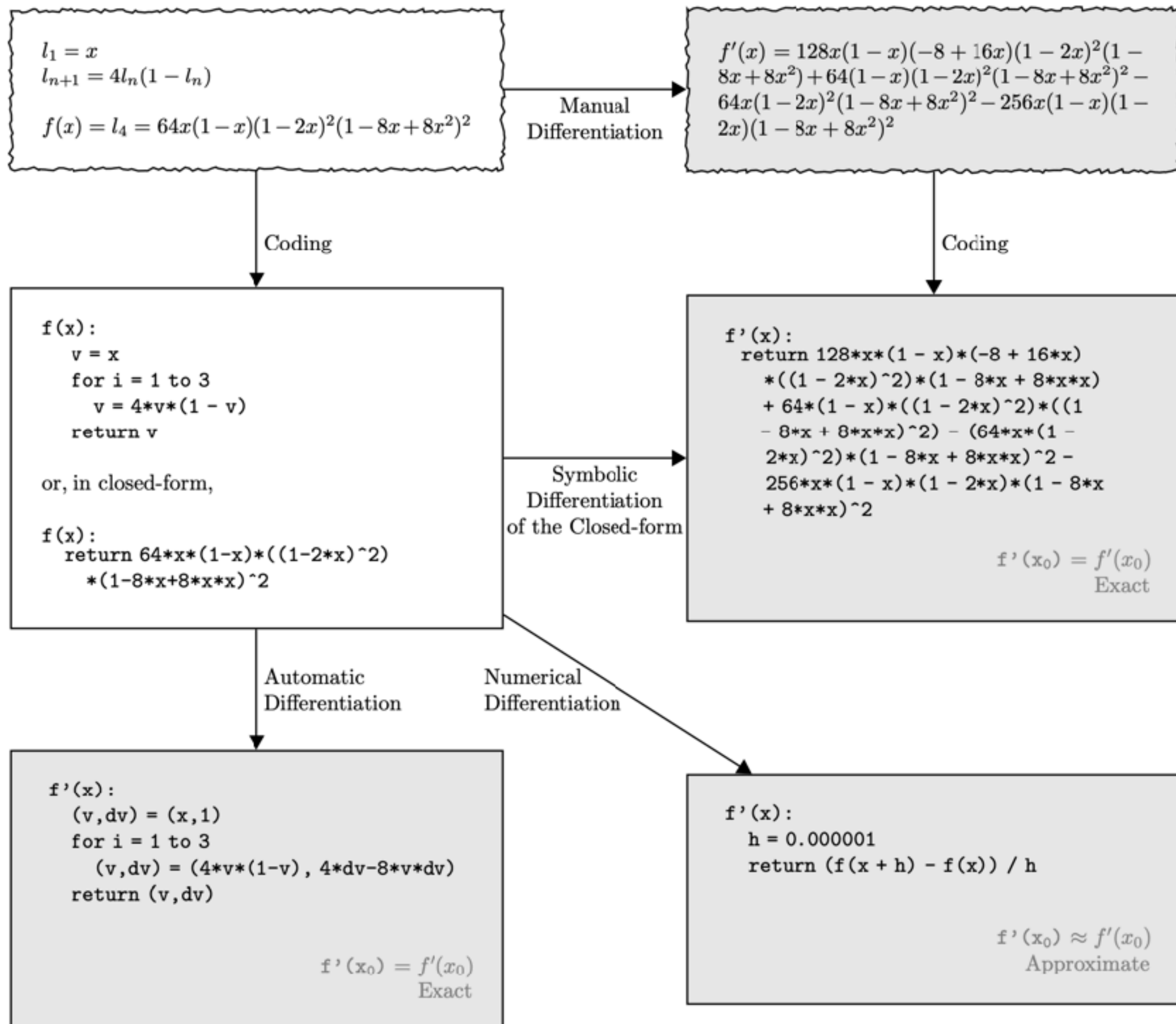
# Backpropagation vs AutoDiff

Backpropagation



AutoDiff





# Recap

- Numerical differentiation
  - Tool to check the correctness of implementation
- Backpropagation
  - Easy to understand and implement
  - Bad for memory use and schedule optimization
- Automatic differentiation
  - Generate gradient computation to entire computation graph
  - Better for system optimization

# References

- Automatic differentiation in machine learning: a survey <https://arxiv.org/abs/1502.05767>
- CS231n backpropagation: <http://cs231n.github.io/optimization-2/>