

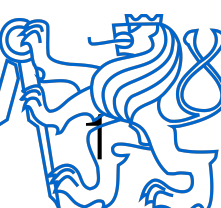
Fully connected networks

Neurons, fully connected networks, computational graphs,
Jacobians and vector-Jacobian product

Karel Zimmermann

Czech Technical University in Prague

Faculty of Electrical Engineering, Department of Cybernetics



Linear classifier and neuron

Labels


RGB images

+1

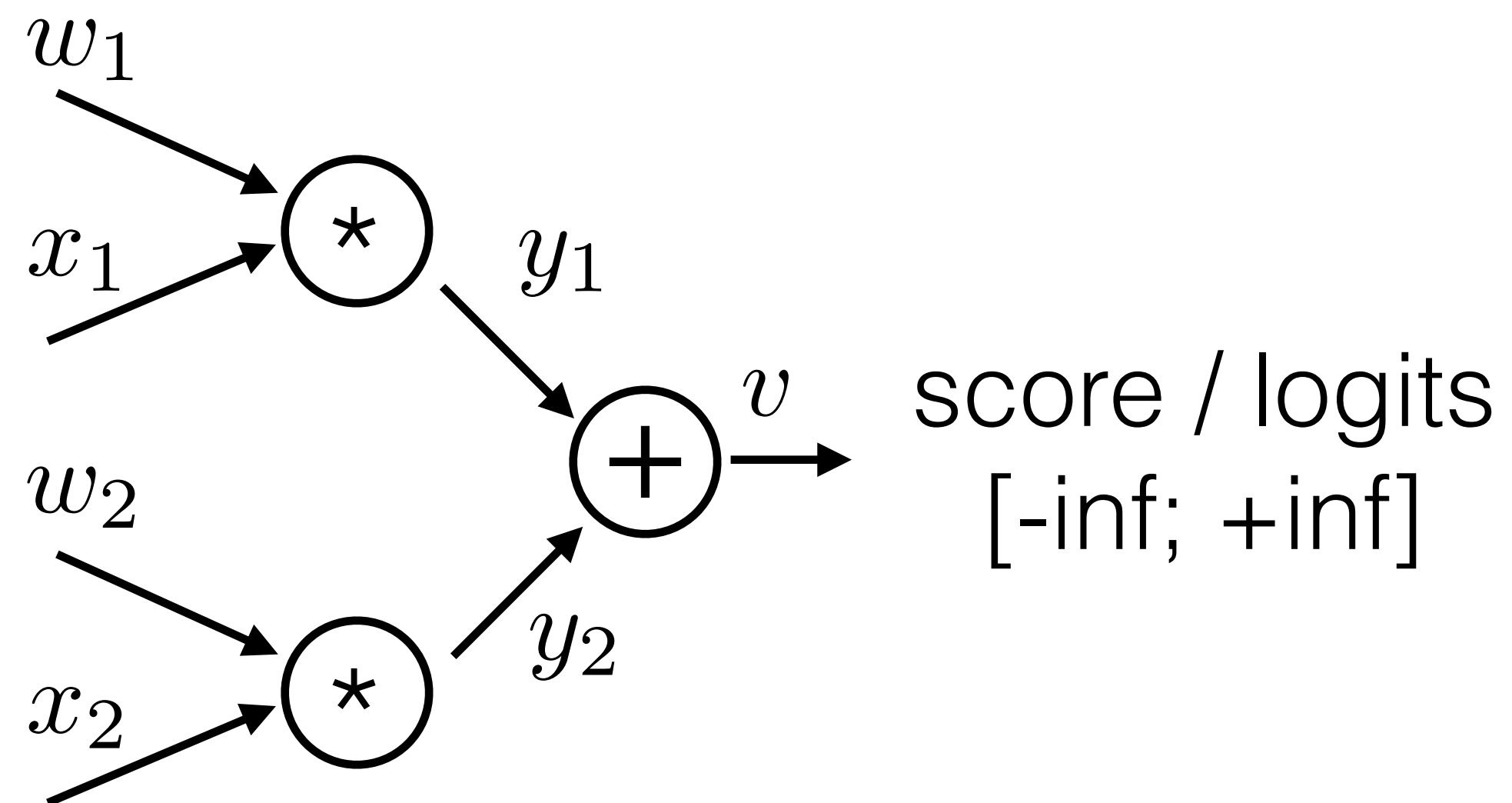


-1



```
def classify(  
    return  $\mathbf{w}^T \mathbf{x}$ 
```

Computational graph of linear classifier



Linear classifier and neuron

Labels


RGB images

+1



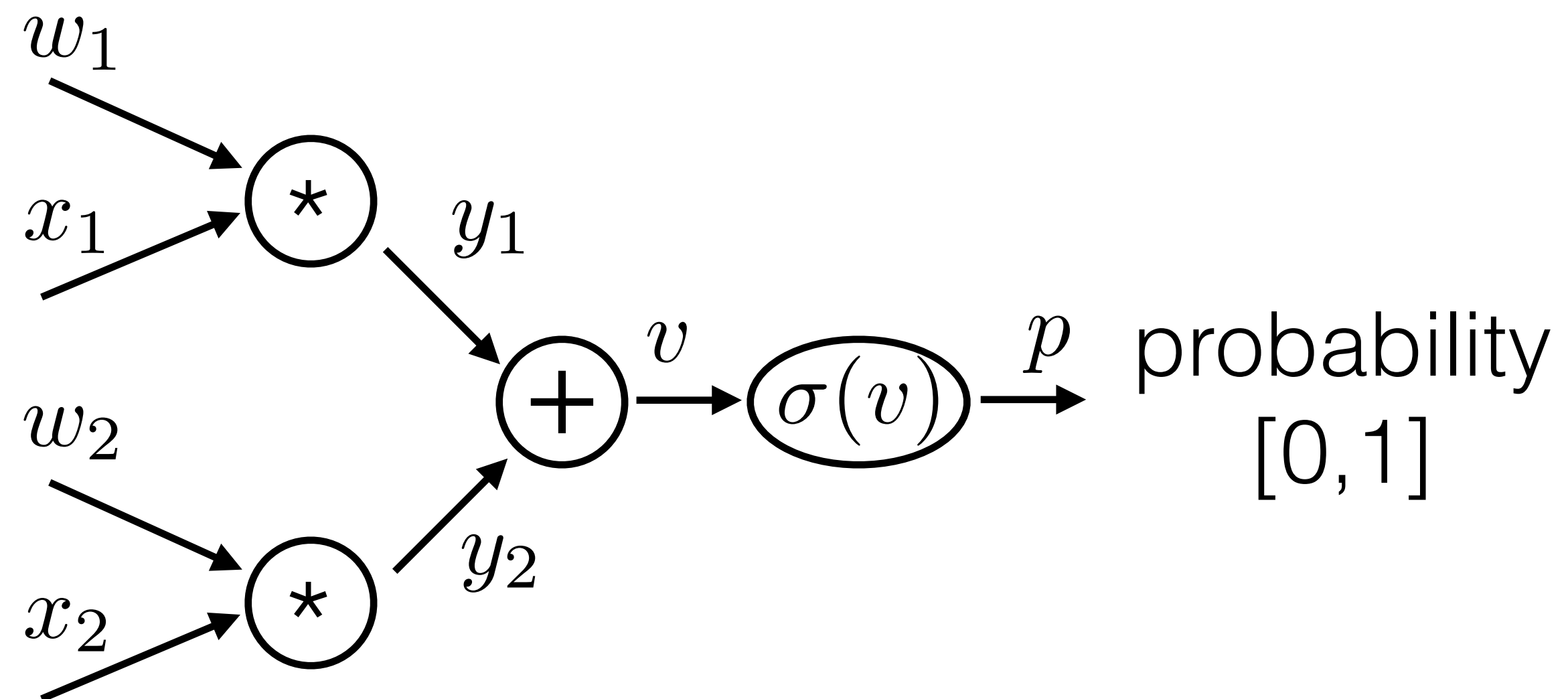
-1



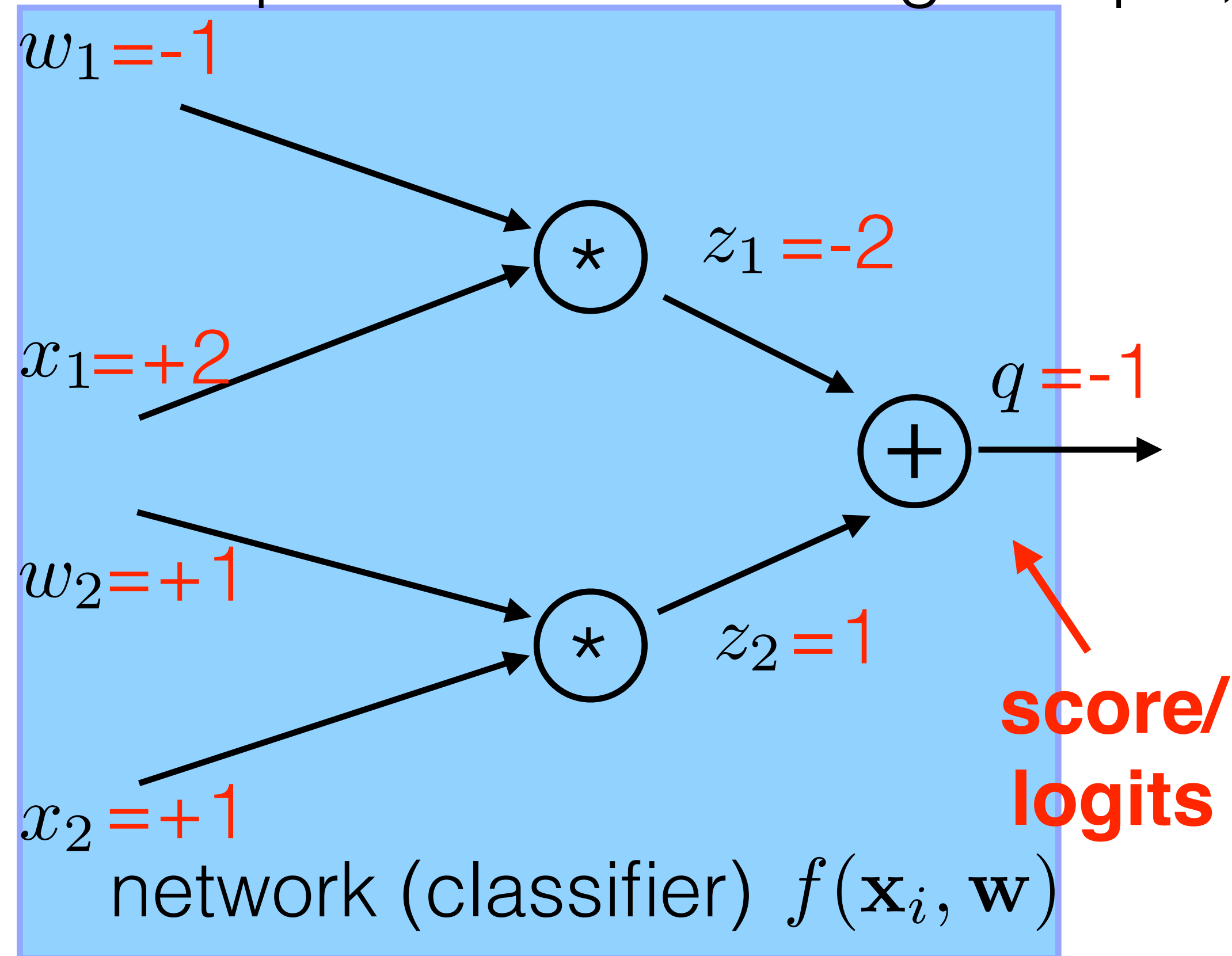
```
def classify(

Computational graph of neuron

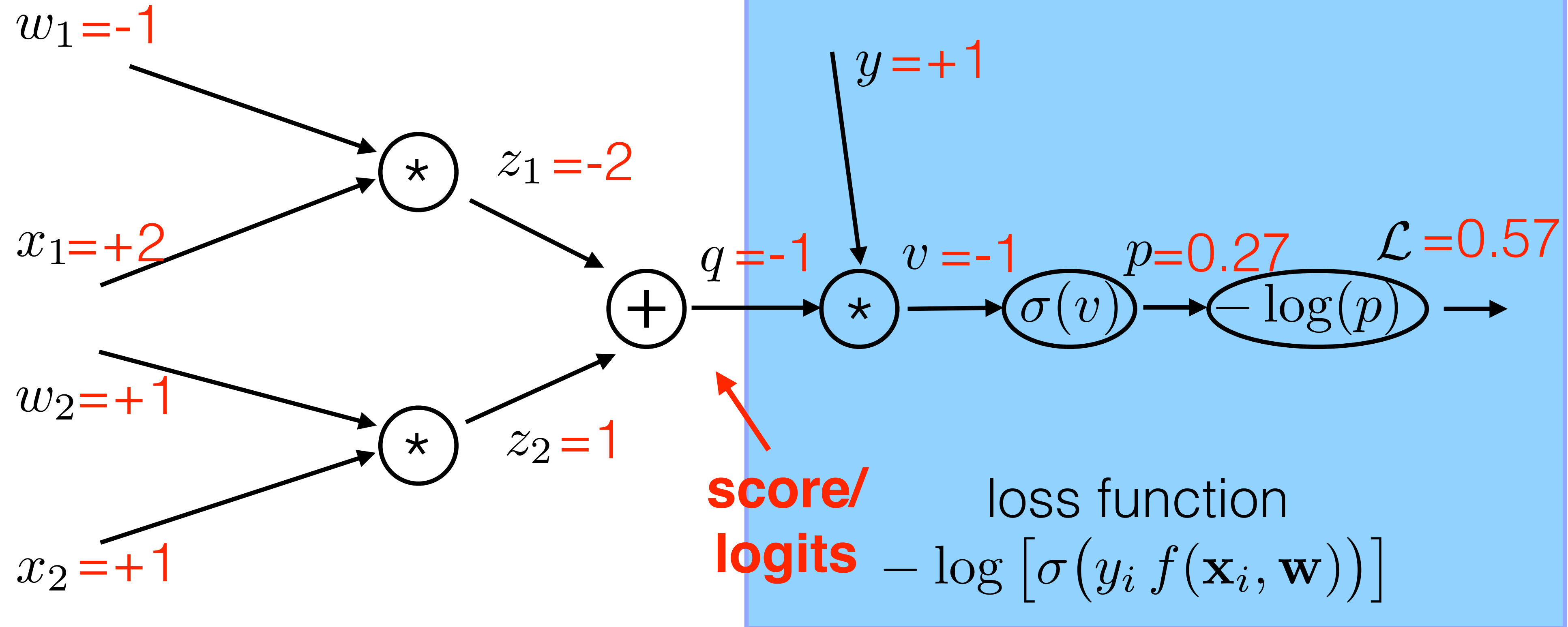

```



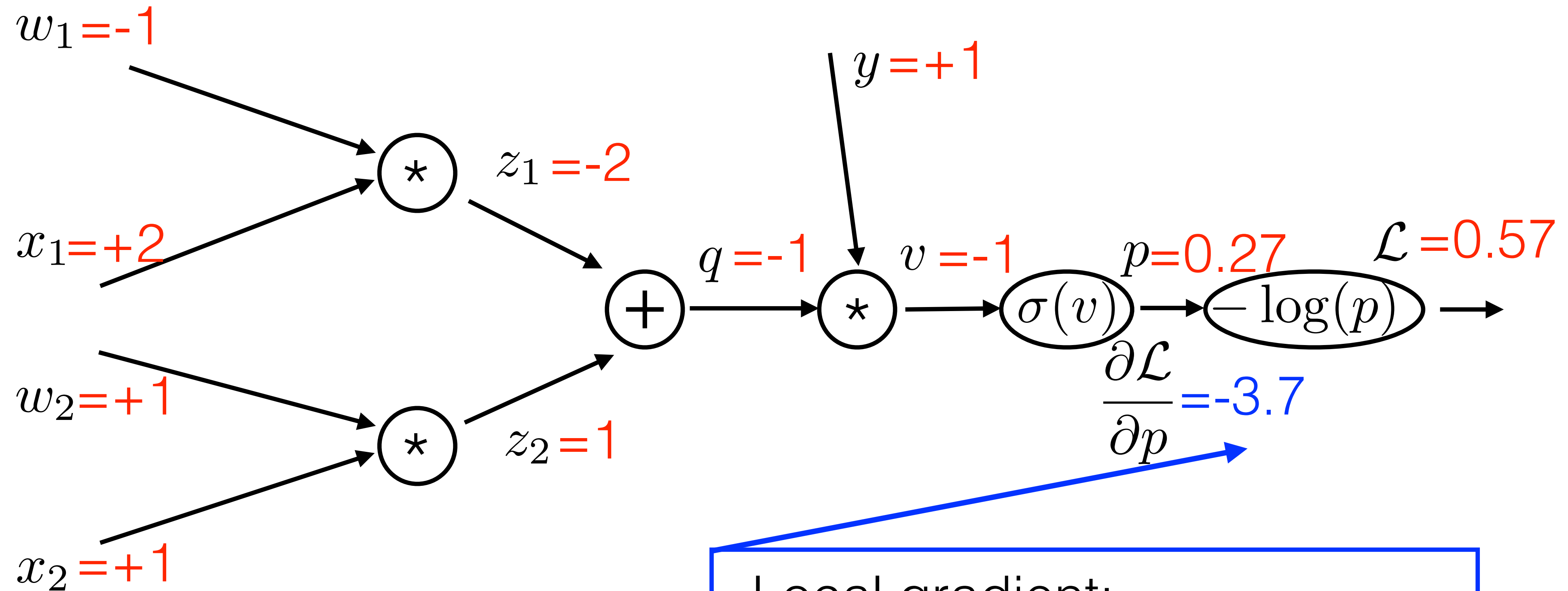
Example II: Given training sample, how do you learn weights?



Example II: Given training sample, how do you learn weights?



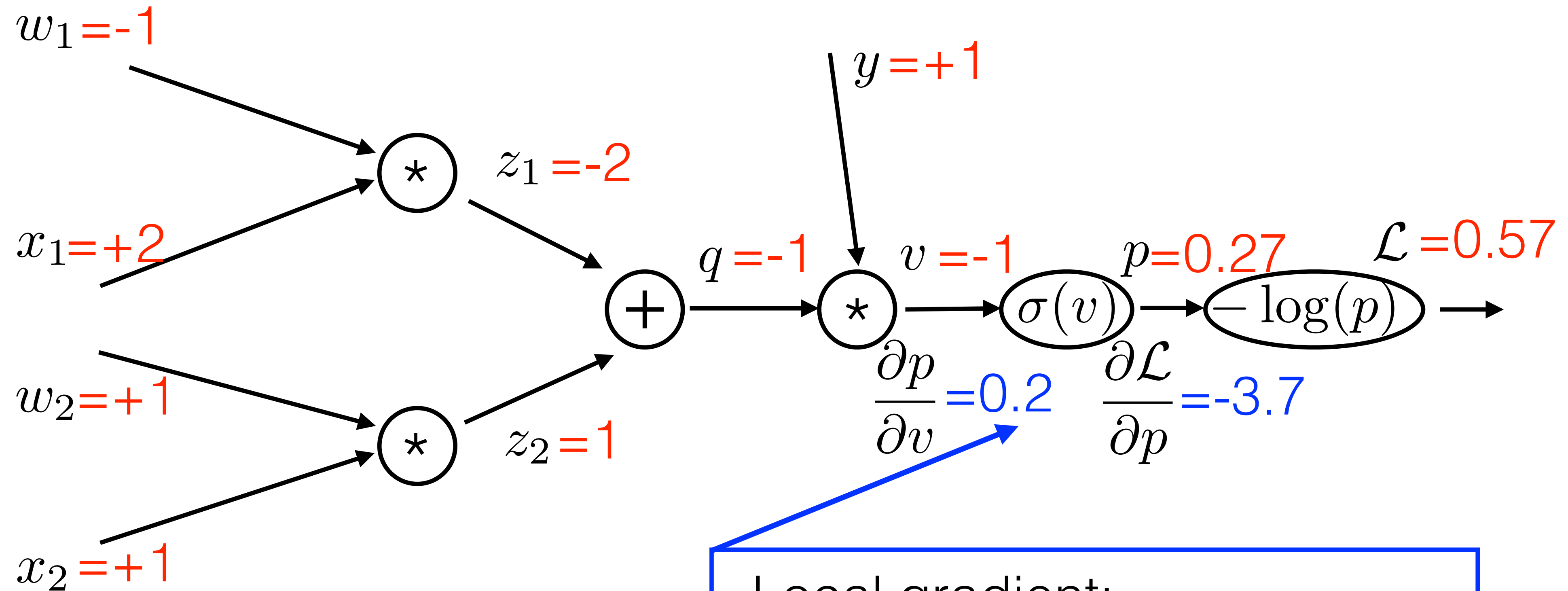
Example II: Given training sample, how do you learn weights?



Local gradient:

$$\frac{\partial \mathcal{L}}{\partial p} = \frac{\partial(-\log(p))}{\partial p} = -\frac{1}{p}$$

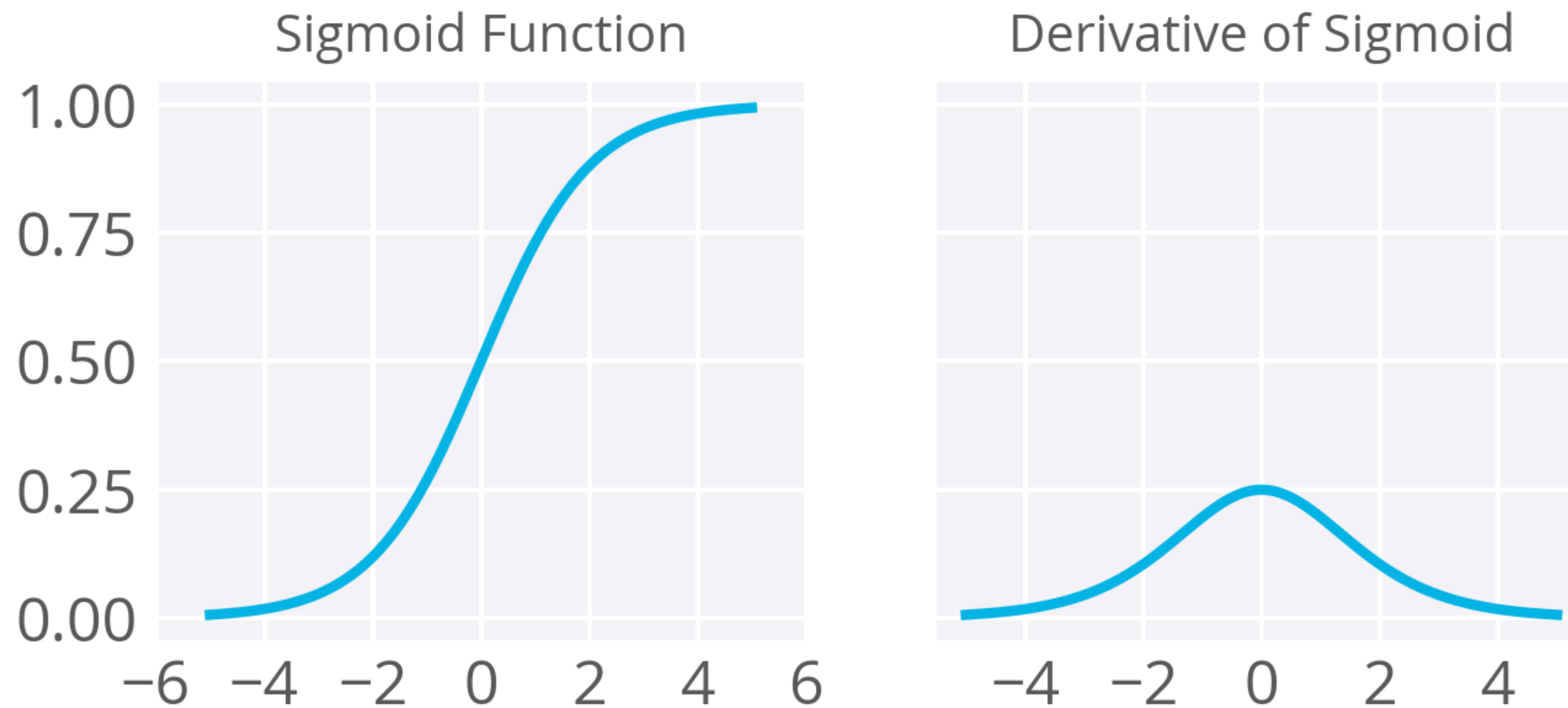
Example II: Given training sample, how do you learn weights?



Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

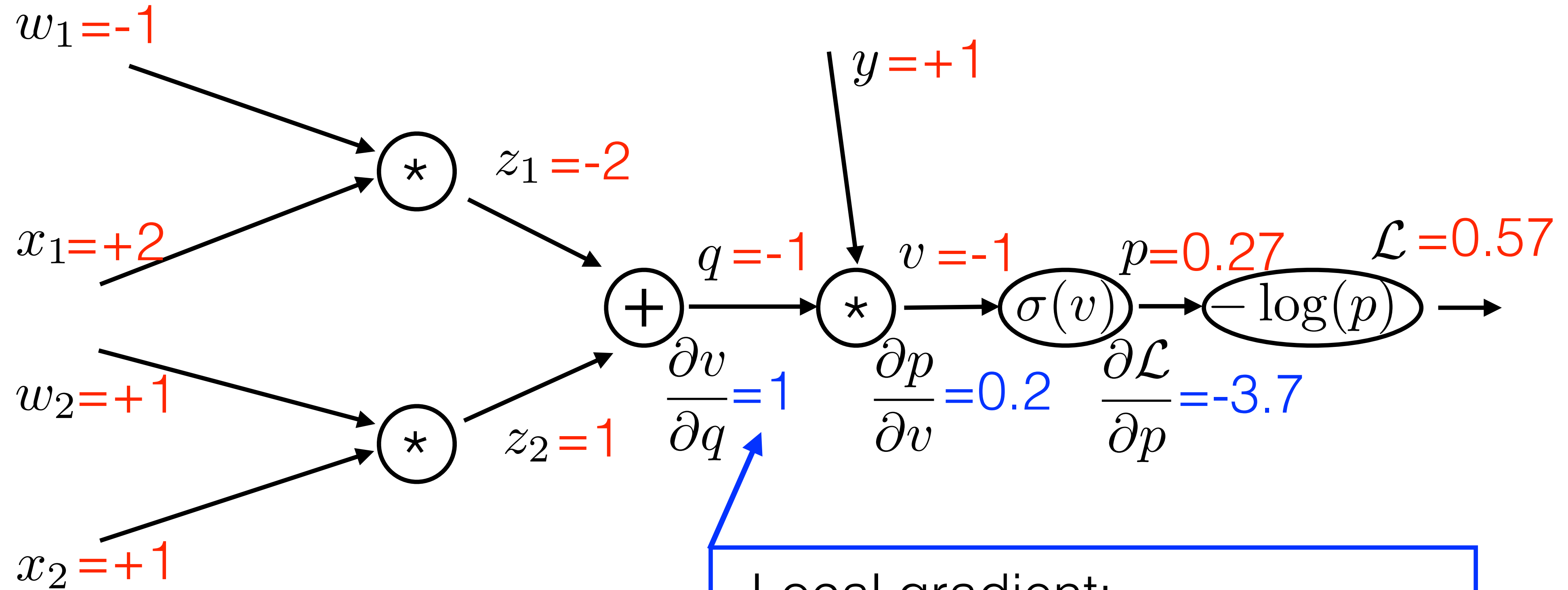
Example II: Given training sample, how do you learn weights?



Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

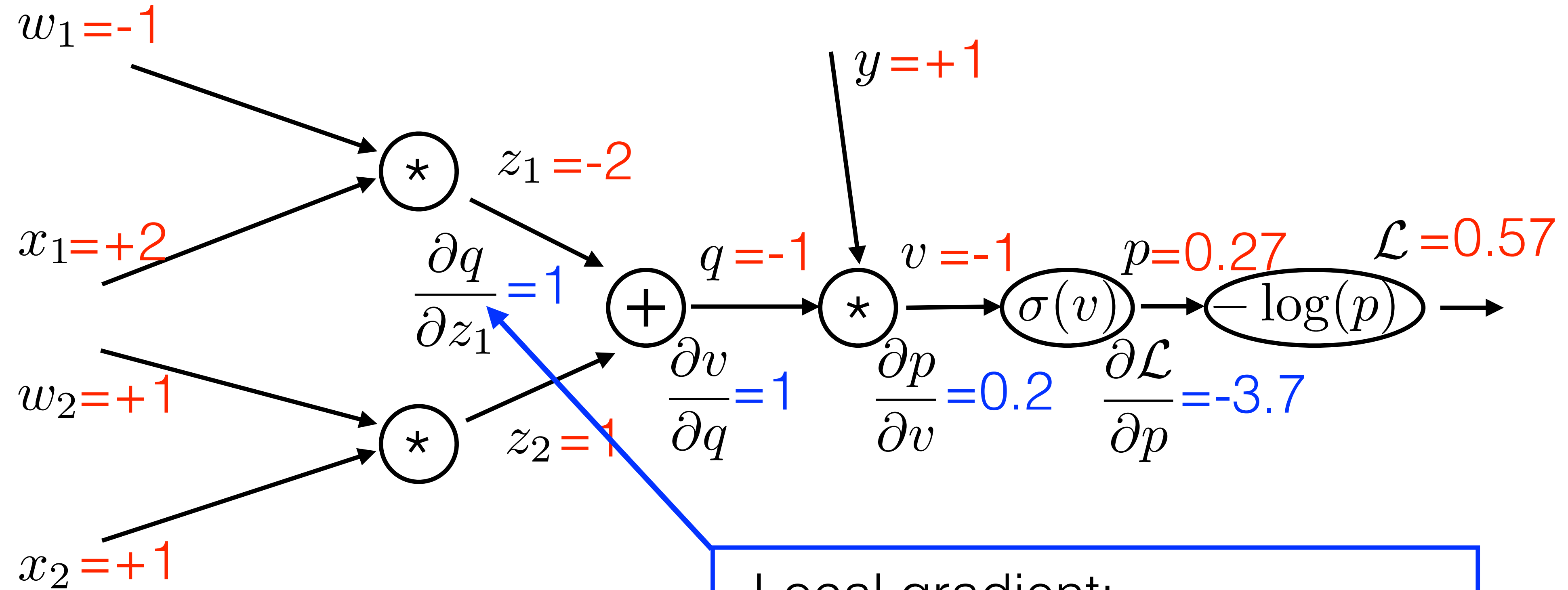
Example II: Given training sample, how do you learn weights?



Local gradient:

$$\frac{\partial v}{\partial q} = \frac{\partial(yq)}{\partial q} = y$$

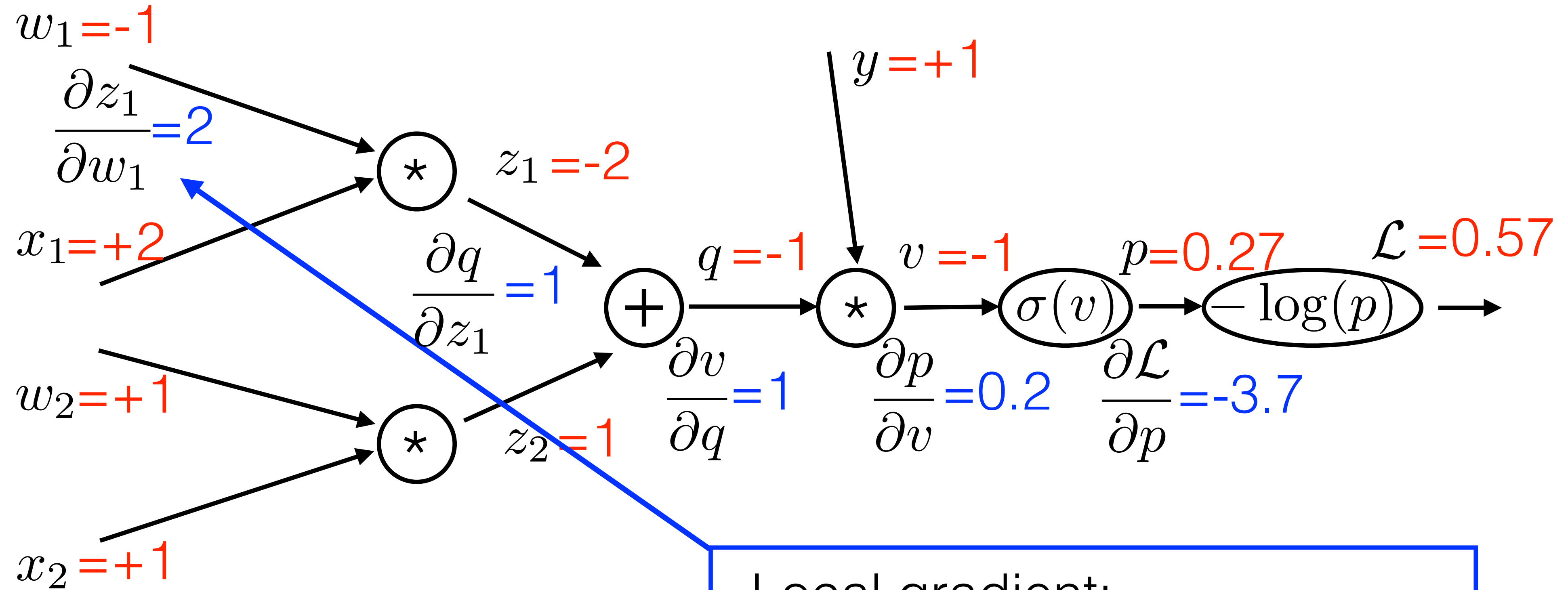
Example II: Given training sample, how do you learn weights?



Local gradient:

$$\frac{\partial q}{\partial z_1} = \frac{\partial(z_1 + z_2)}{\partial z_1} = 1$$

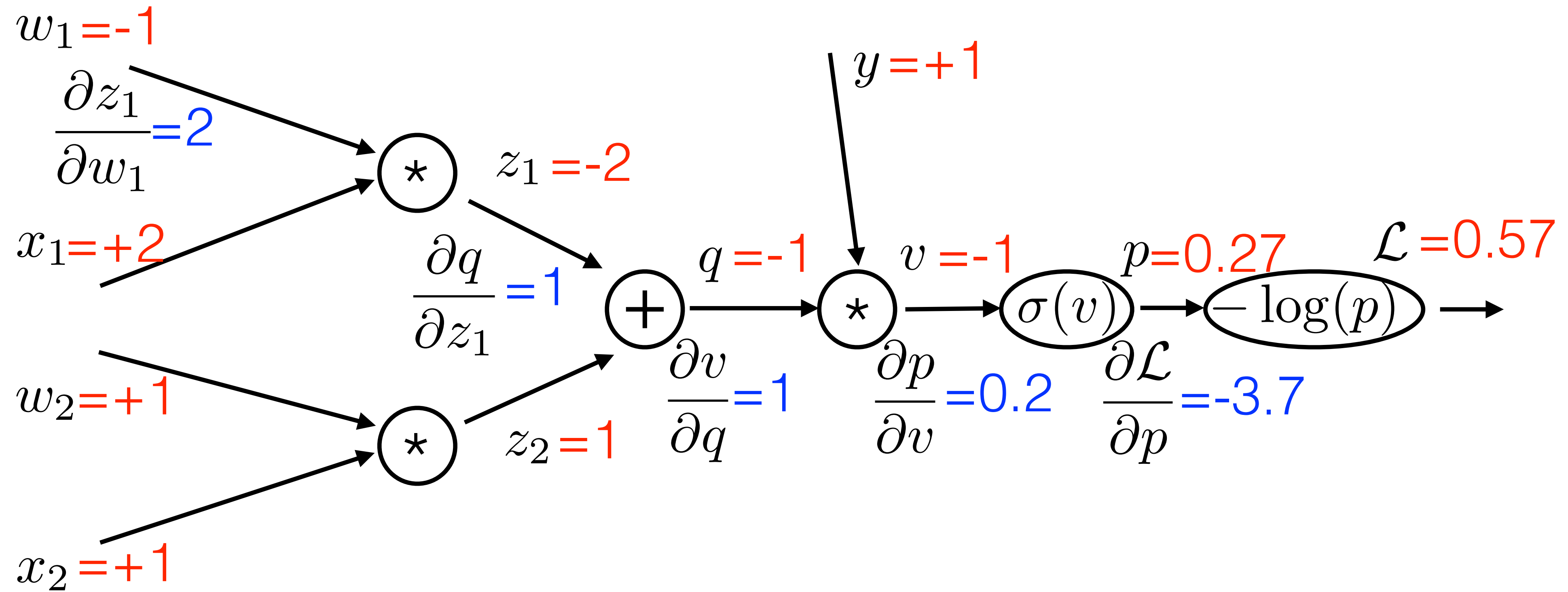
Example II: Given training sample, how do you learn weights?



Local gradient:

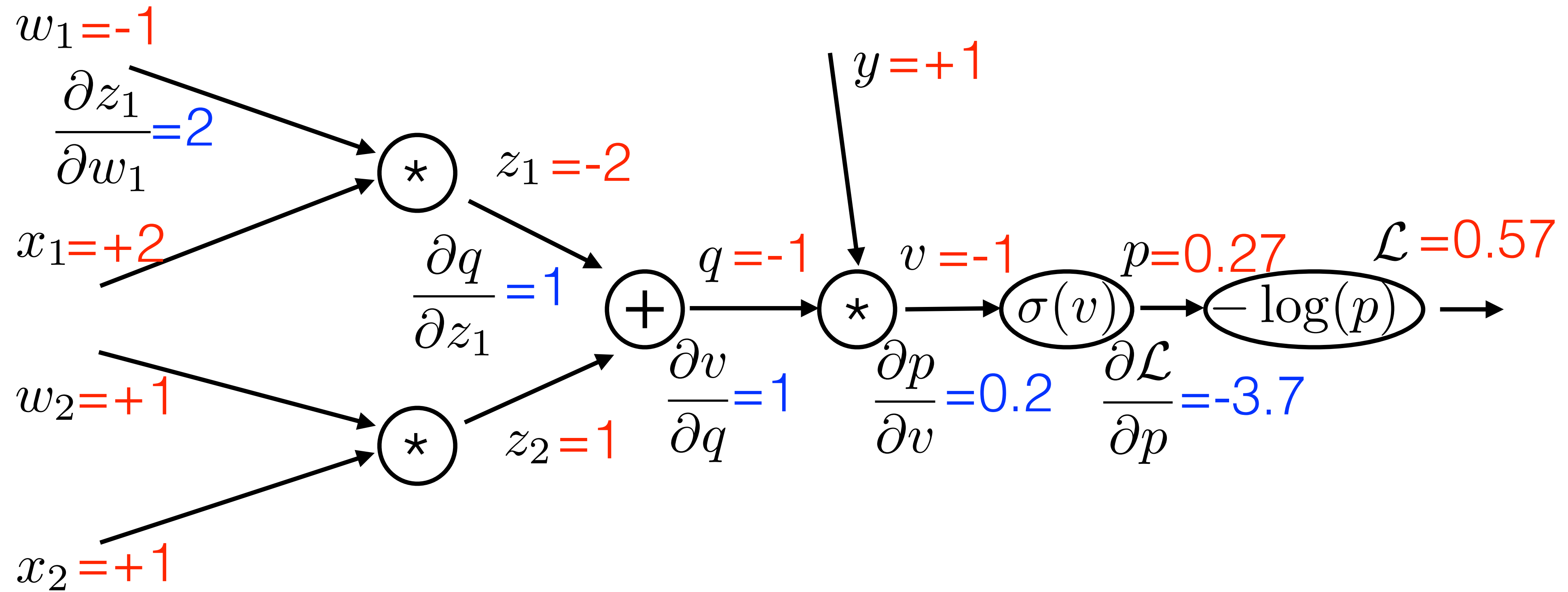
$$\frac{\partial z_1}{\partial w_1} = \frac{\partial(w_1 x_1)}{\partial w_1} = x_1$$

Example II: Given training sample, how do you learn weights?



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

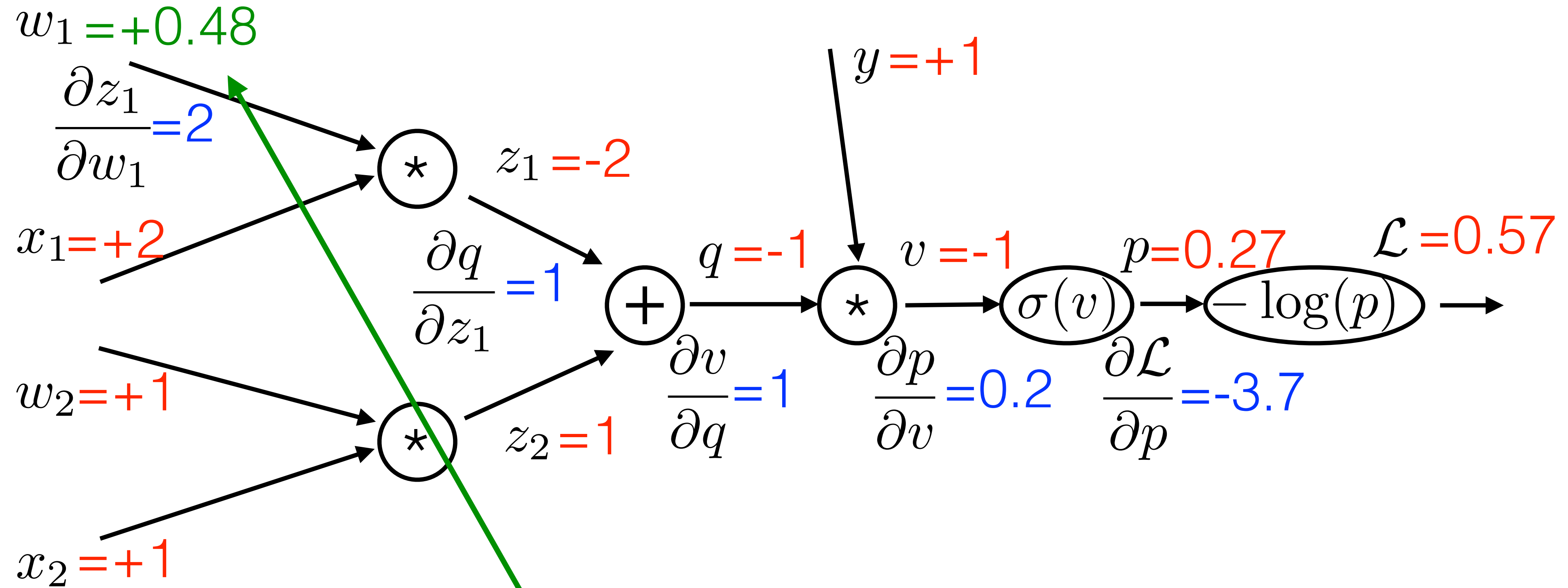
Example II: Given training sample, how do you learn weights?



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$

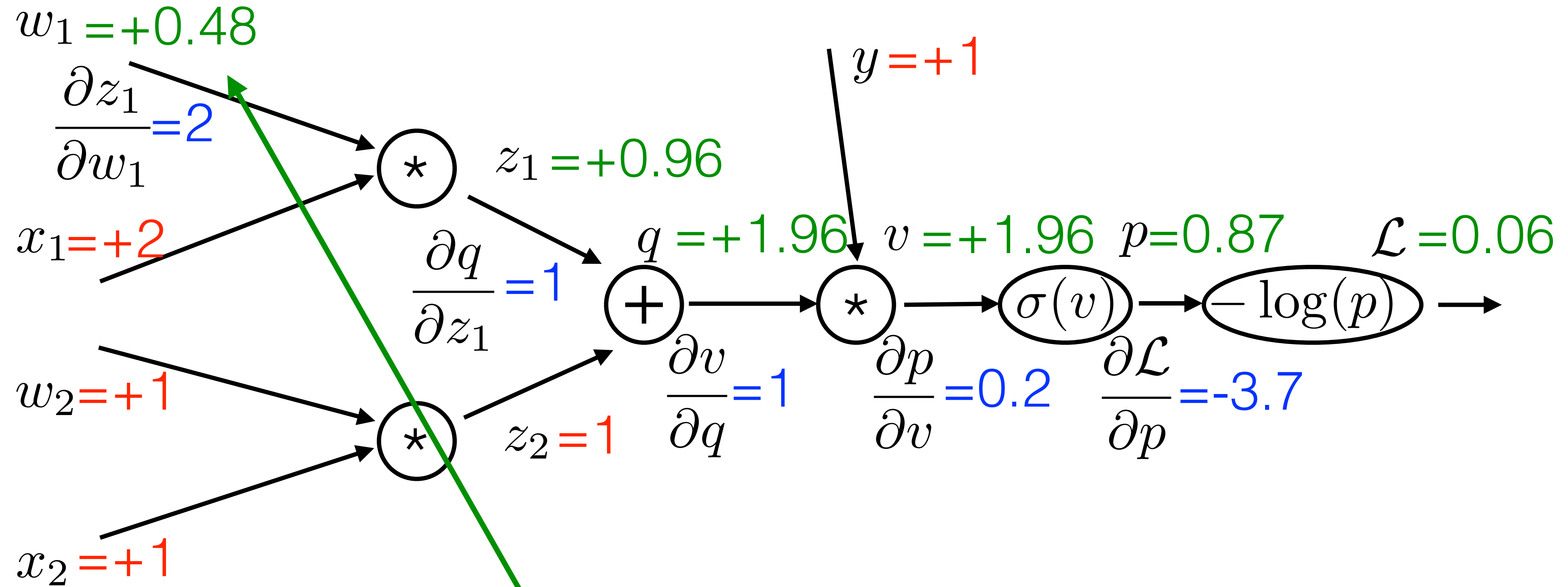
Example II: Given training sample, how do you learn weights?



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$

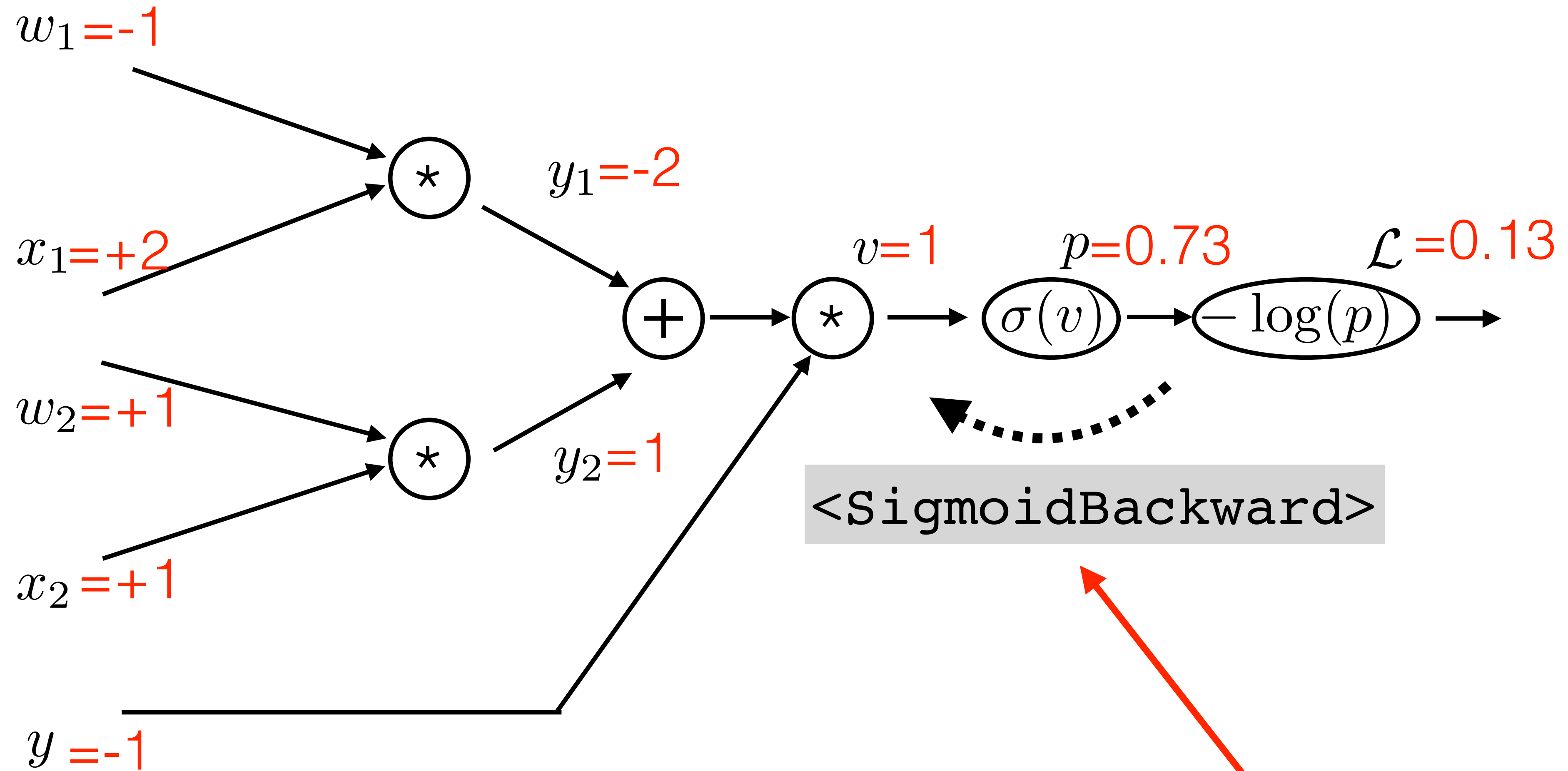
Example II: Given training sample, how do you learn weights?



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

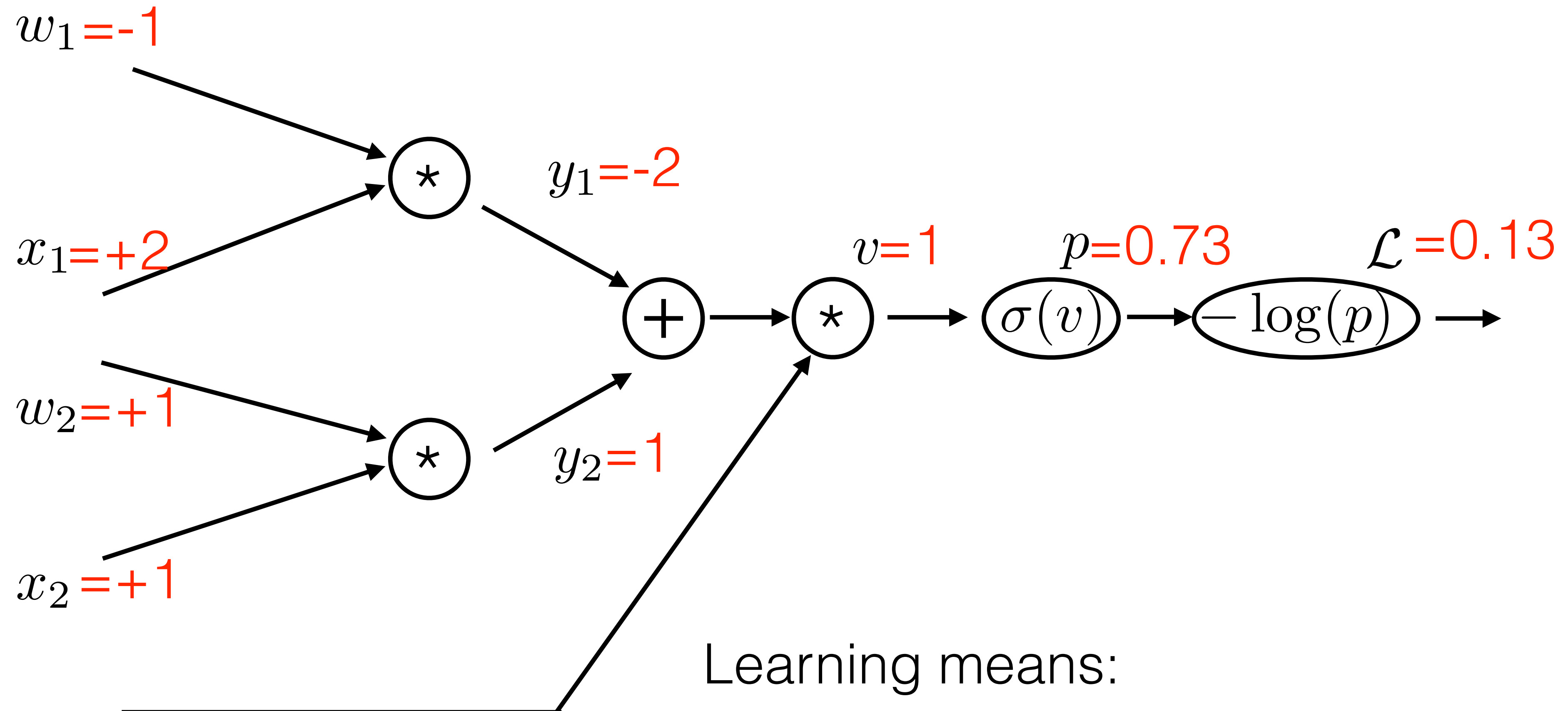
$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$

PyTorch representation



```
p = torch.sigmoid(v)
tensor(0.7311, dtype=torch.float64, grad_fn=<SigmoidBackward>)
```


Example III: vector representation

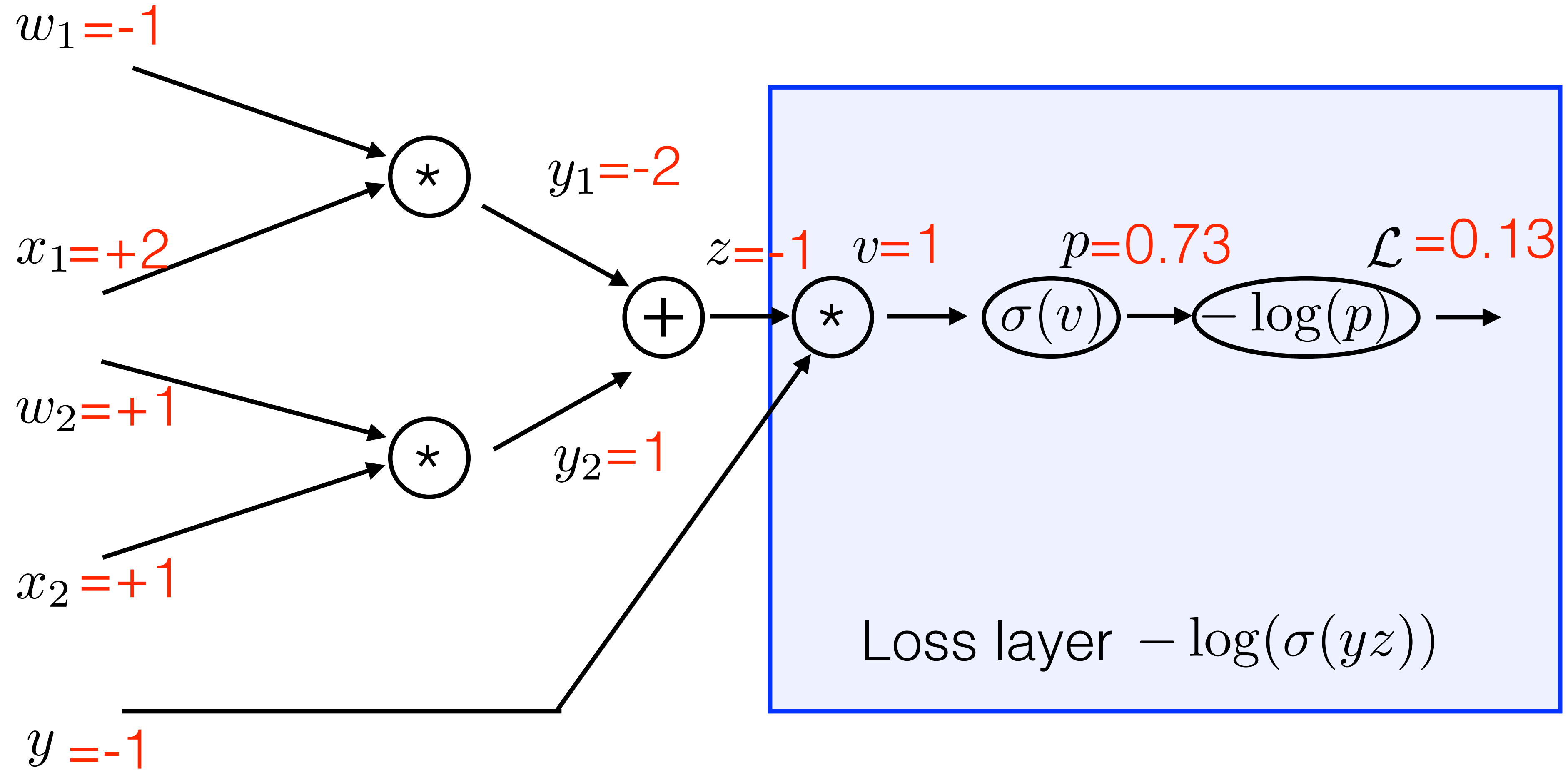


Learning means:

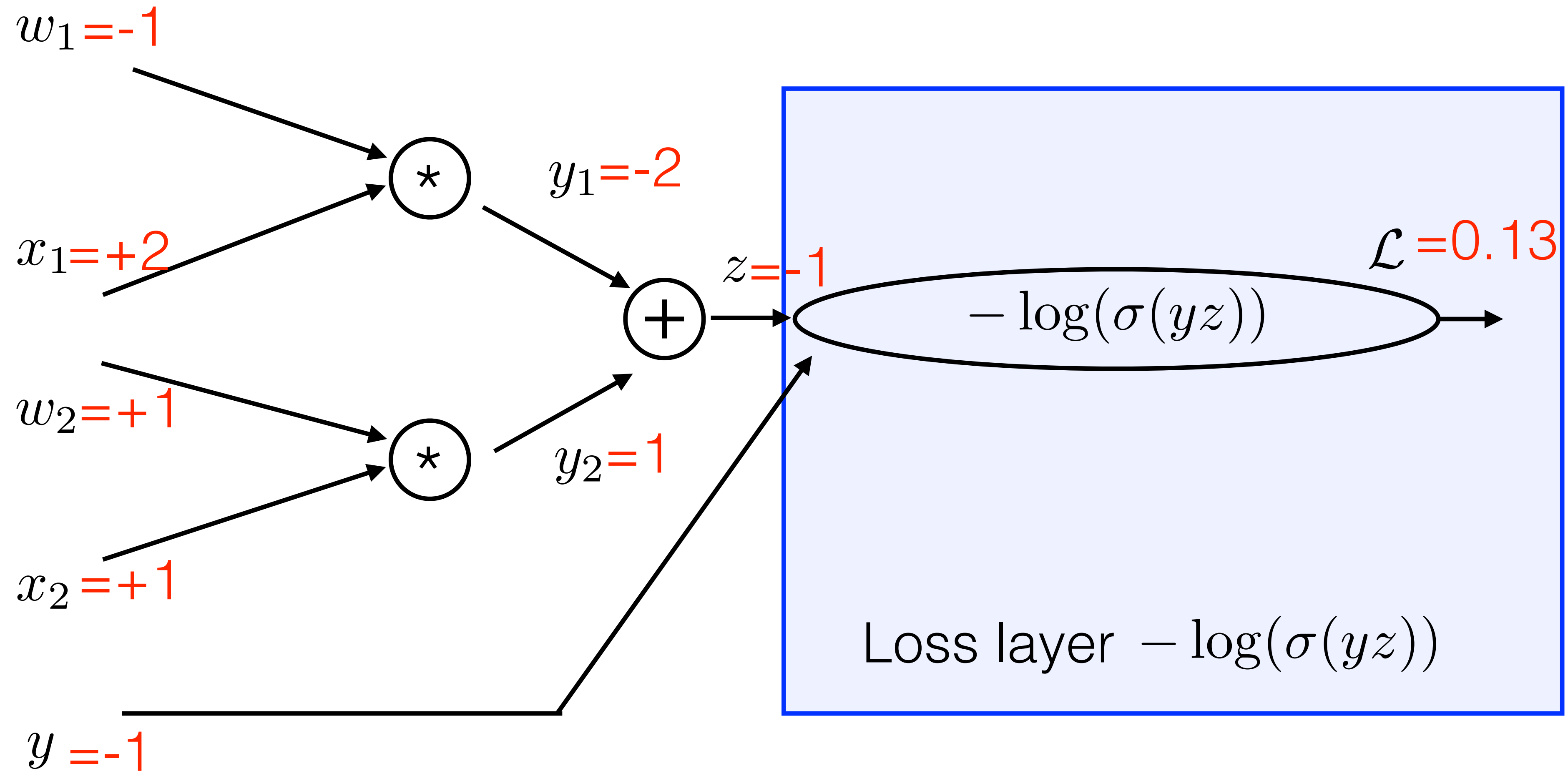
Iteratively change all weights \mathbf{w} to minimize \mathcal{L}

$$\mathbf{w} = \mathbf{w} - \alpha \left[\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^\top \quad \text{where} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \left[\frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots \right]$$

Computational graph of the learning



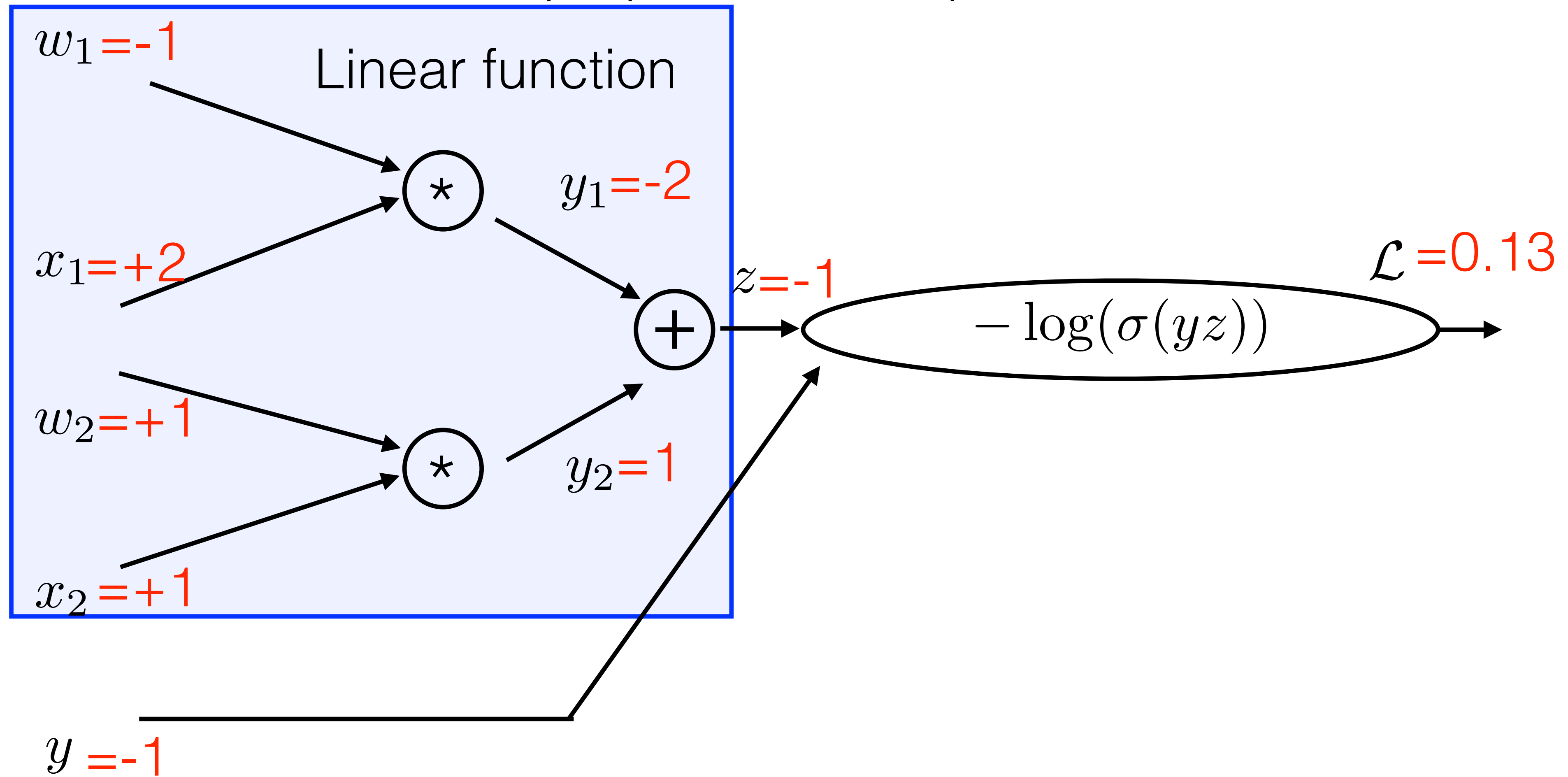
Backprop in vector representation



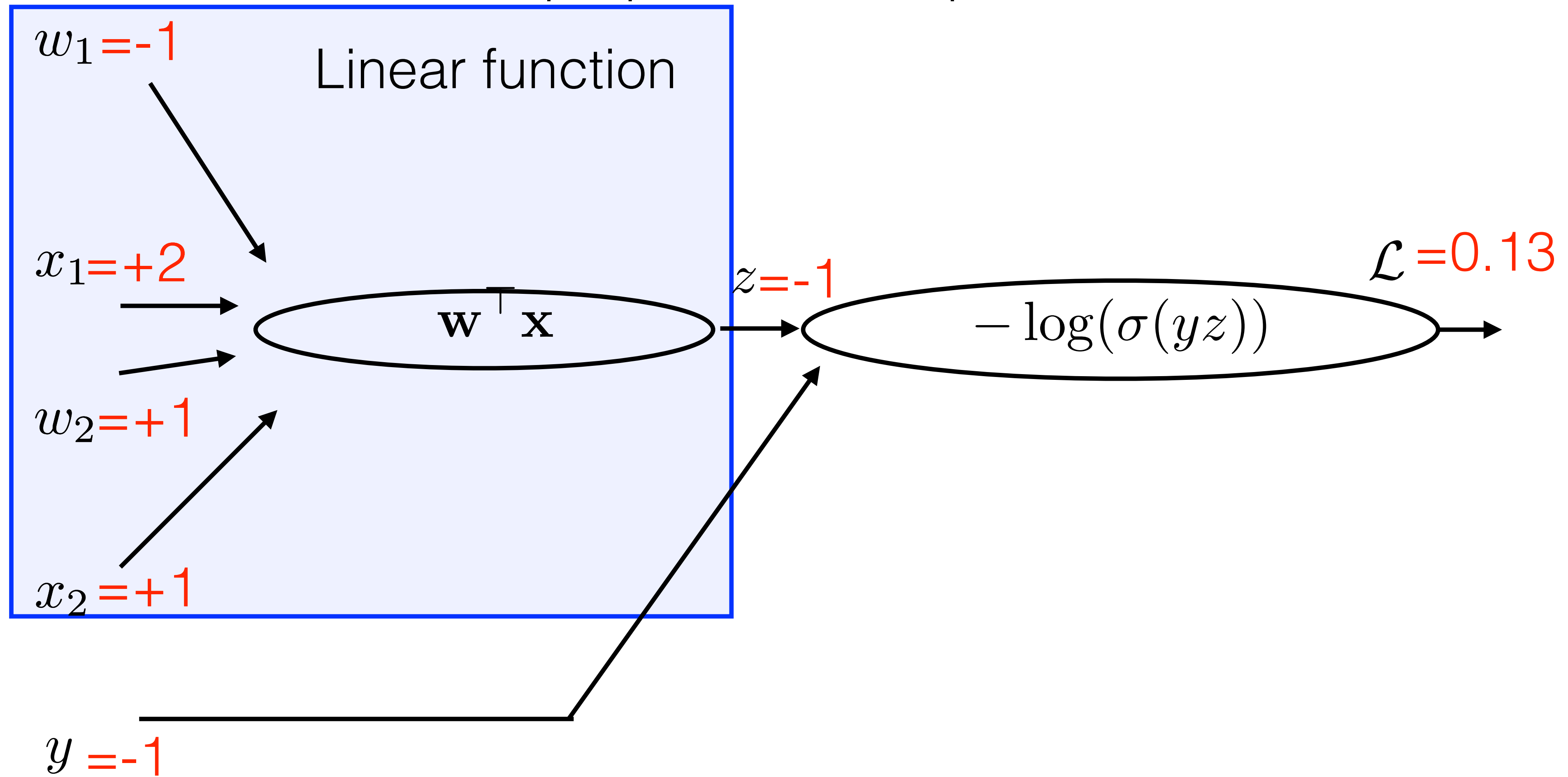
This is the logistic loss!

$$\mathcal{L}(y, z) = -\log(\sigma(yz)) = \log(1 + \exp(-yz))$$

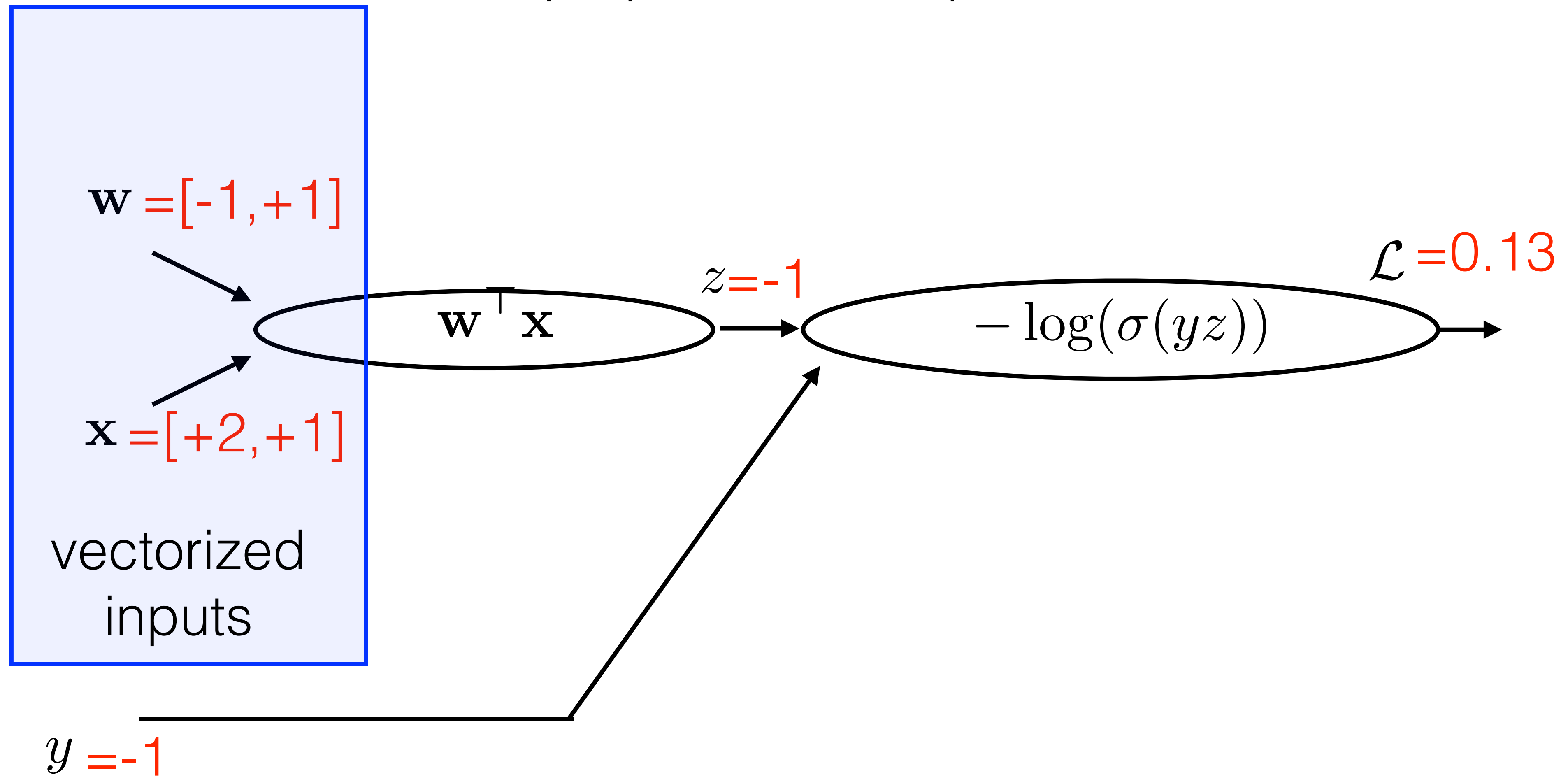
Backprop in vector representation



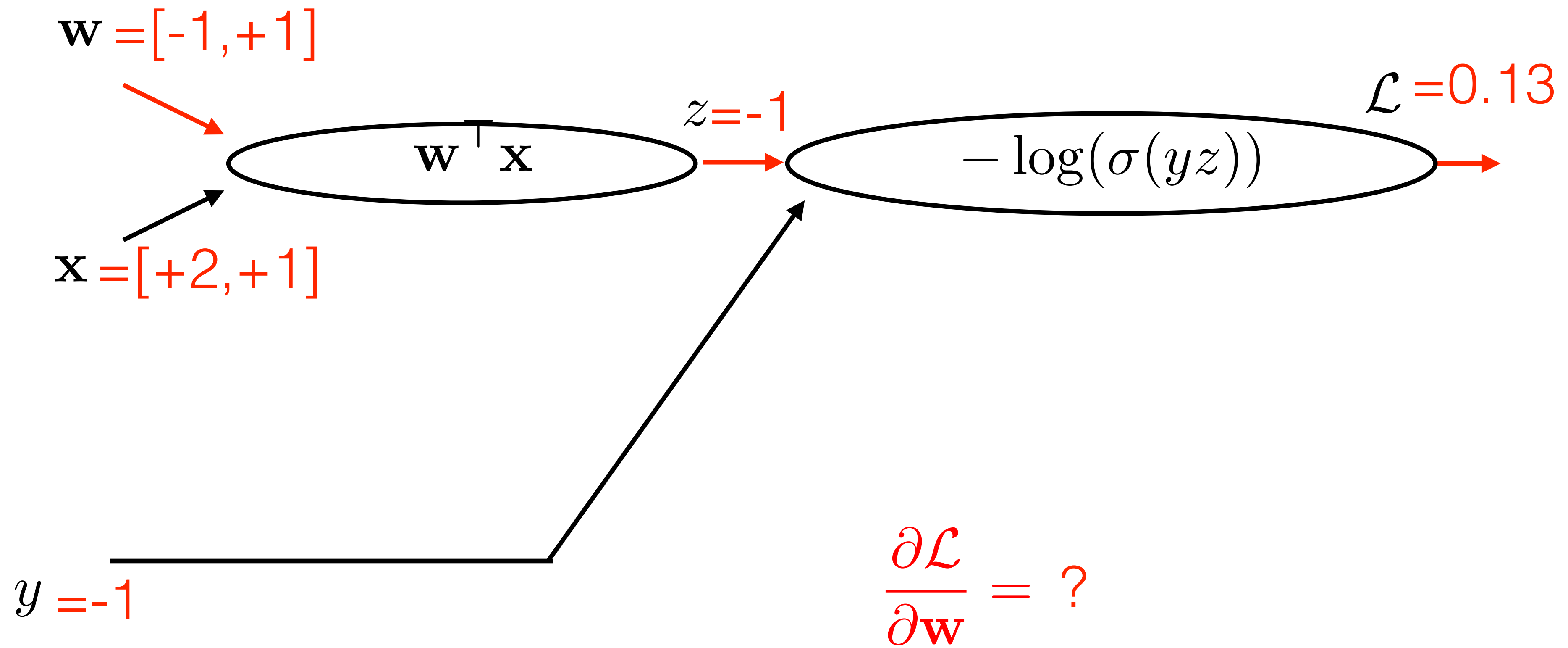
Backprop in vector representation



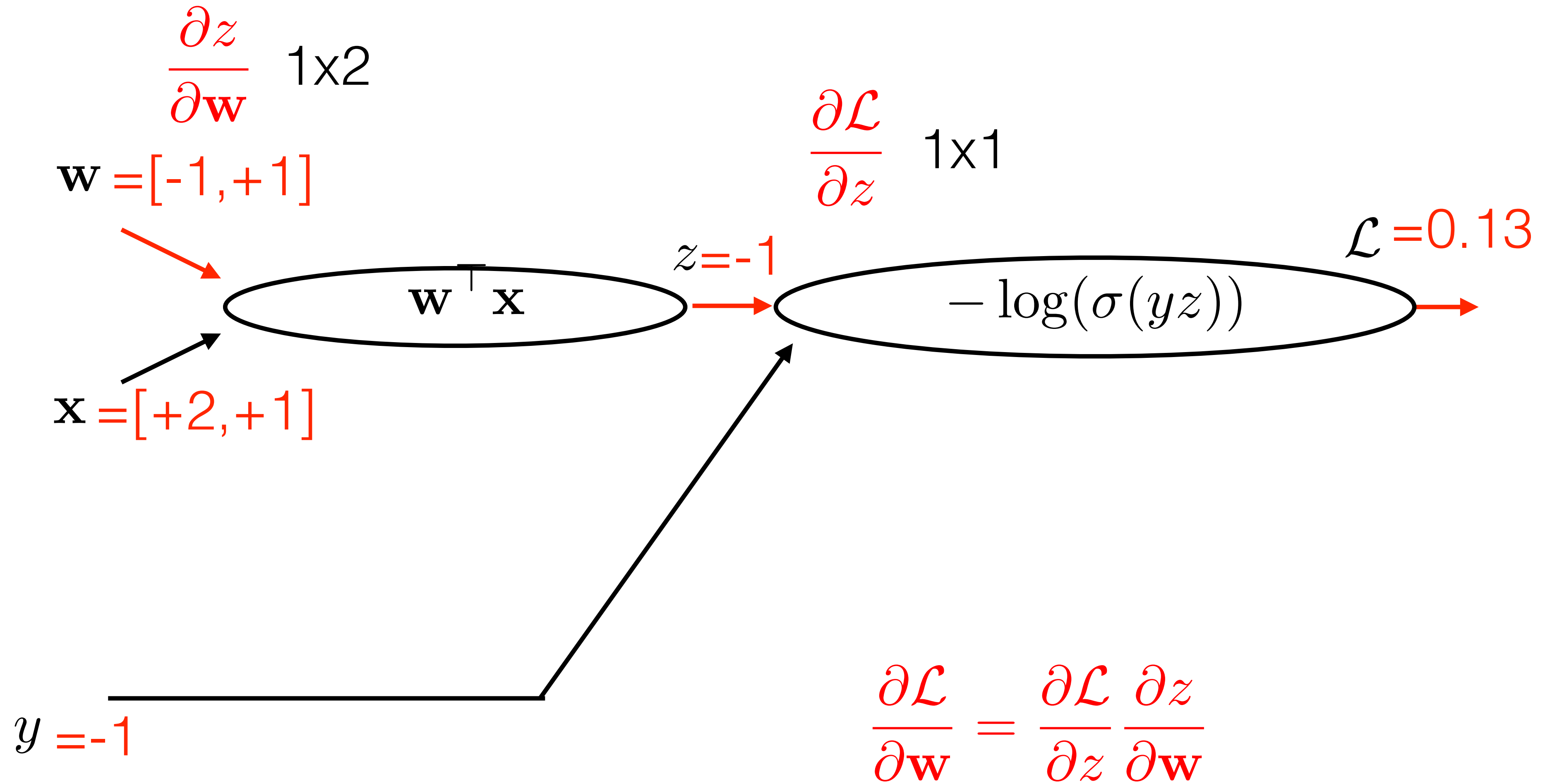
Backprop in vector representation



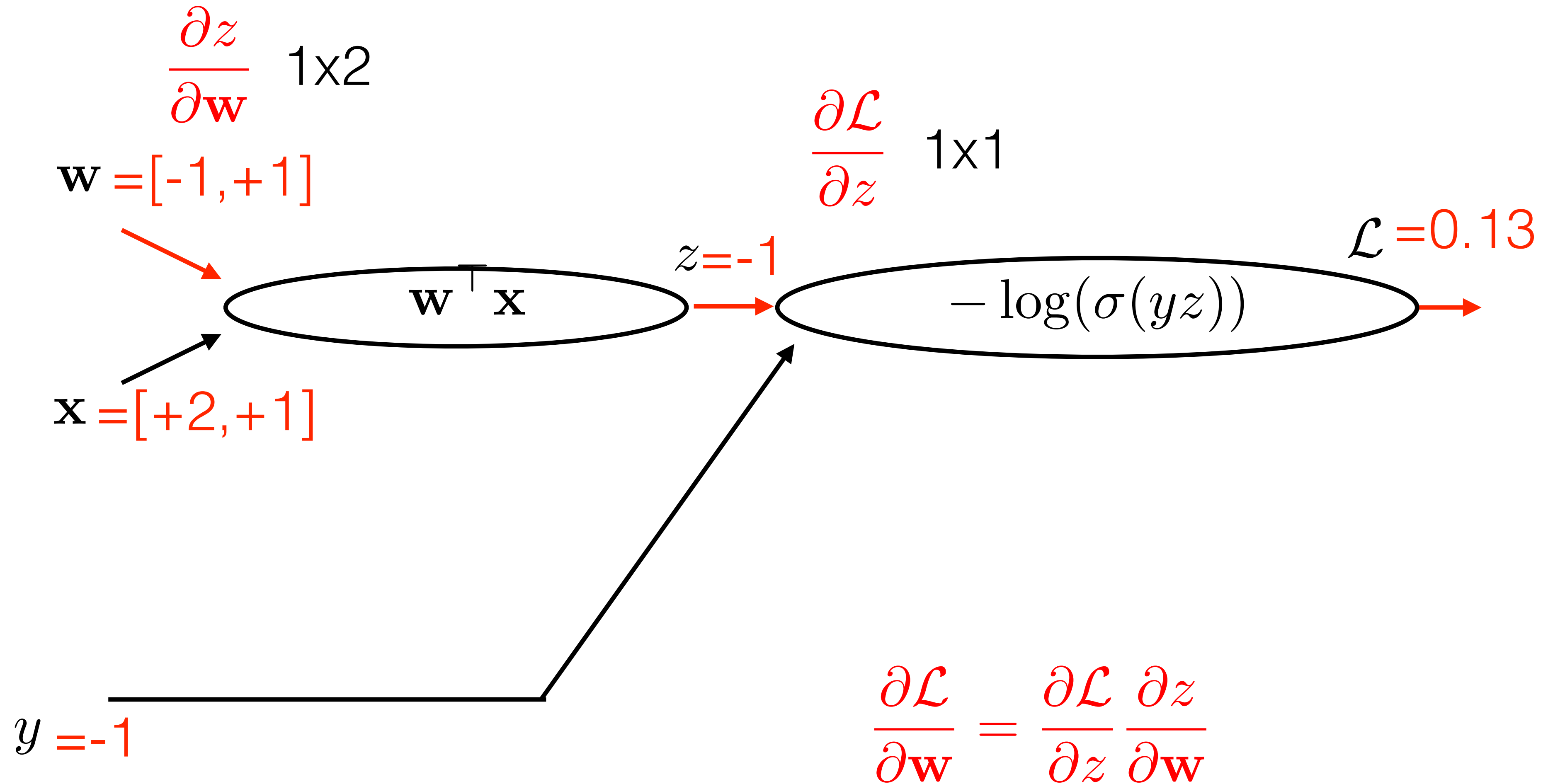
Backprop in vector representation



Backprop in vector representation

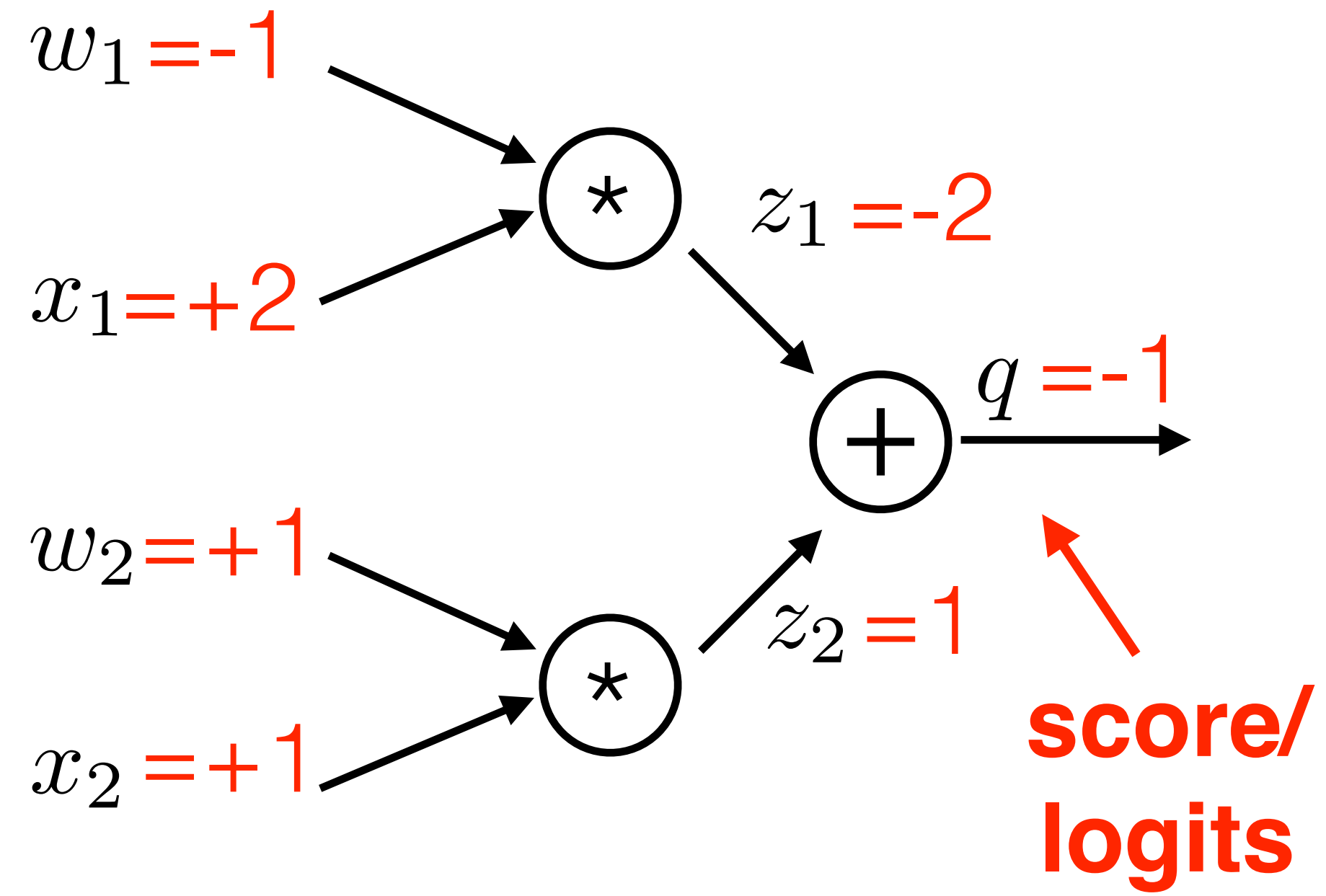


Backprop in vector representation

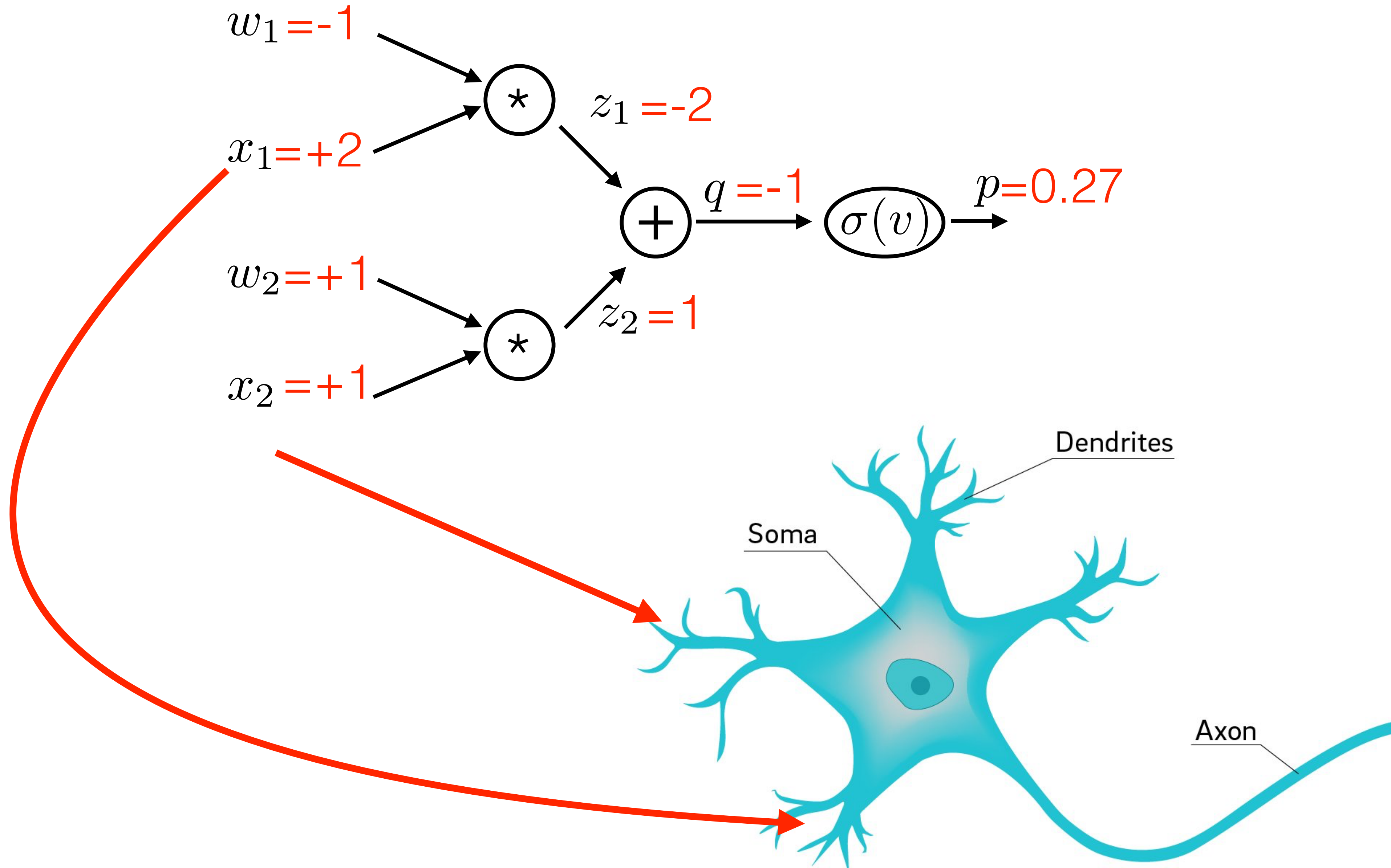


Learning from multiple training samples means summing up the partial derivative over all samples

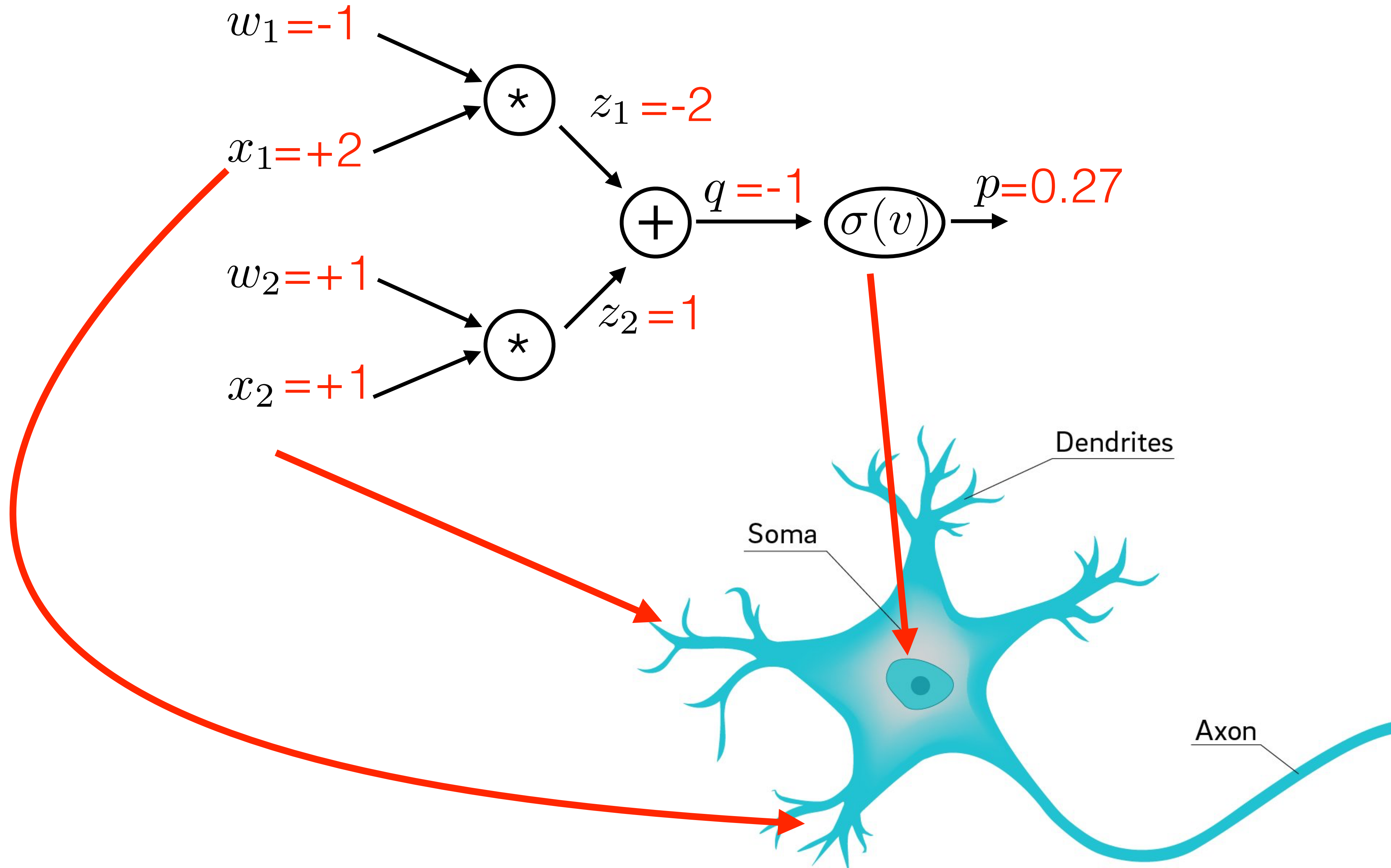
Linear classifier + sigmoid = neuron



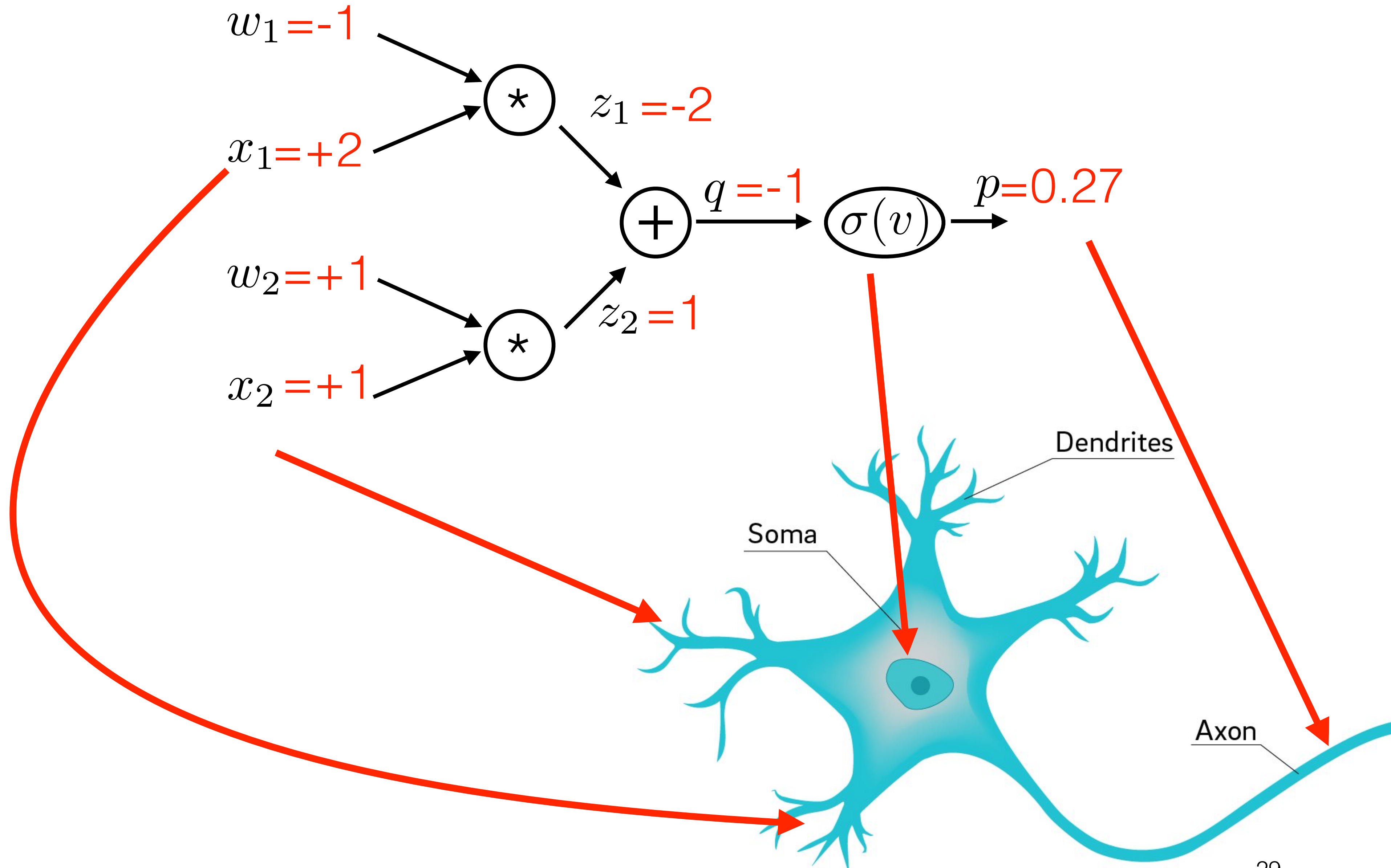
Relation to biological neuron



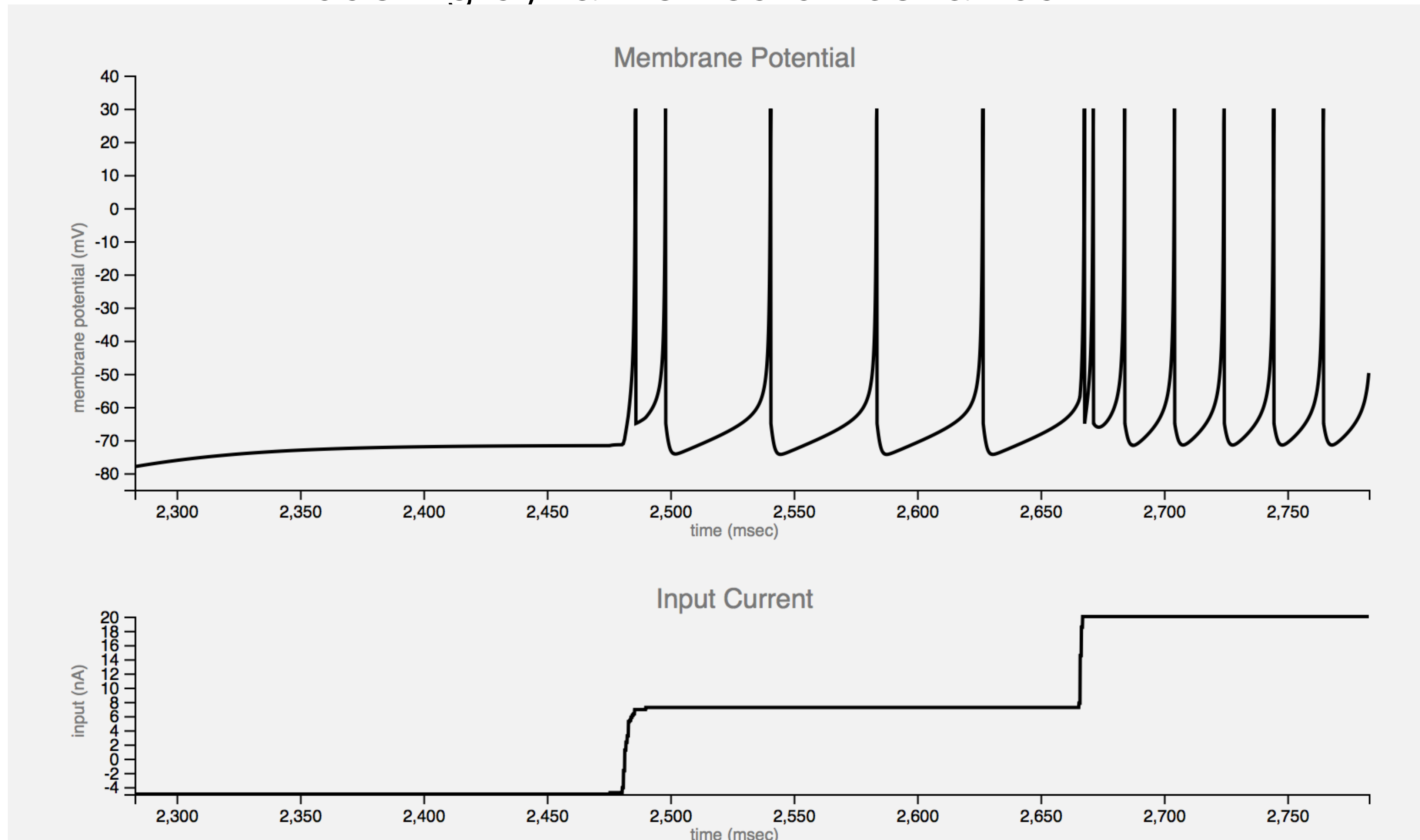
Relation to biological neuron



Relation to biological neuron

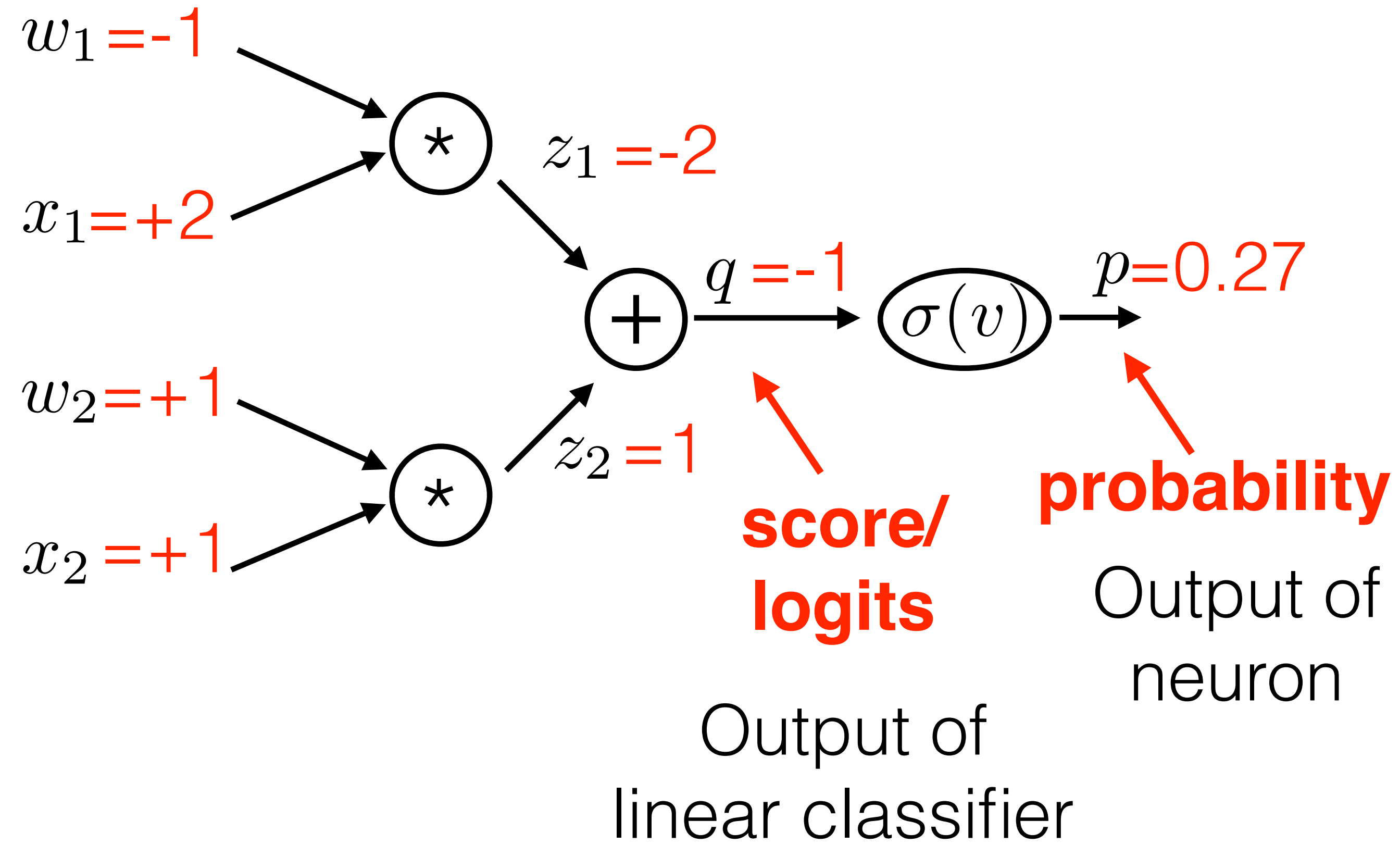


Modeling dynamic neuron behaviour

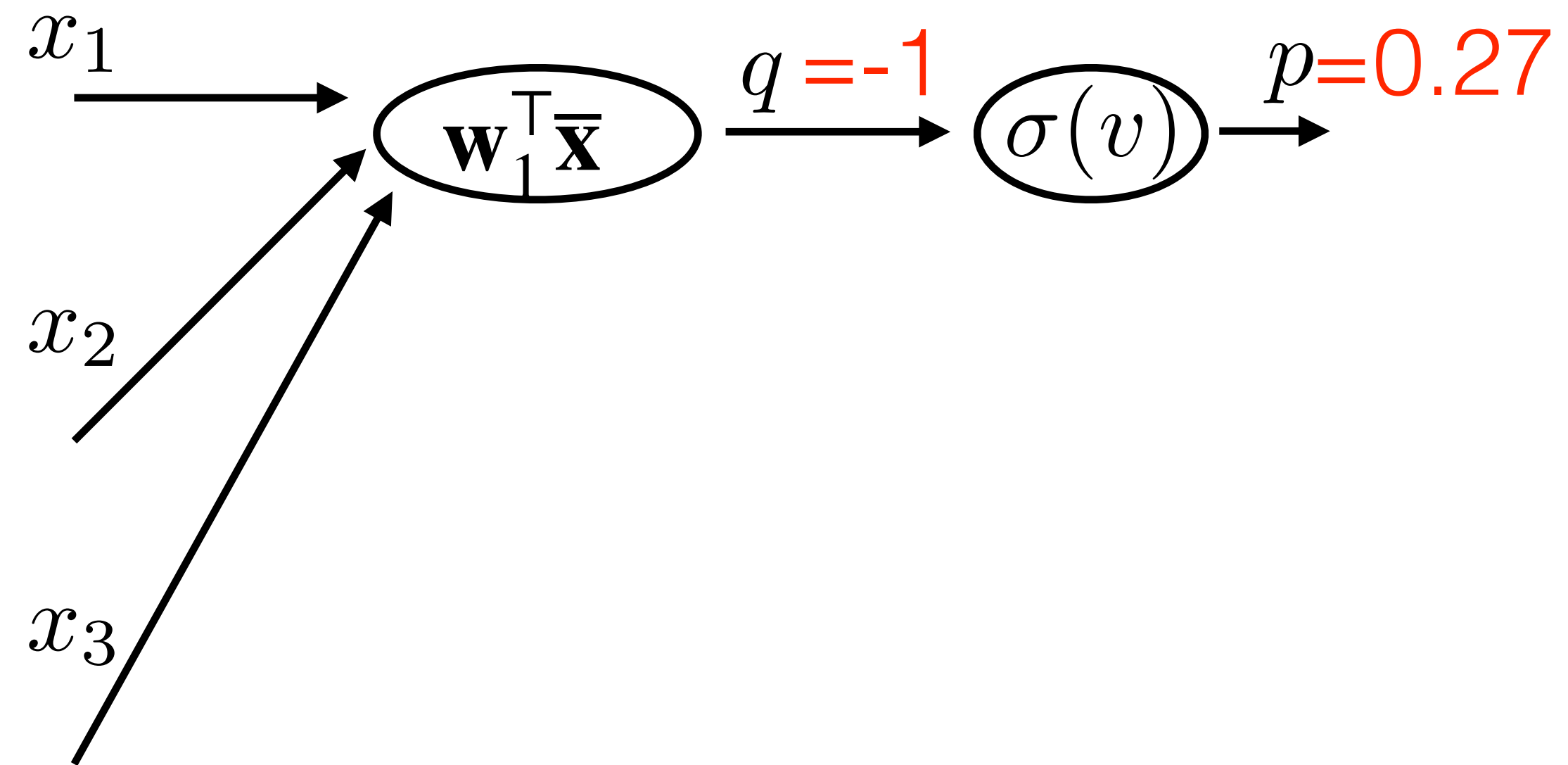


<http://jackterwilliger.com/biological-neural-networks-part-i-spiking-neurons/>

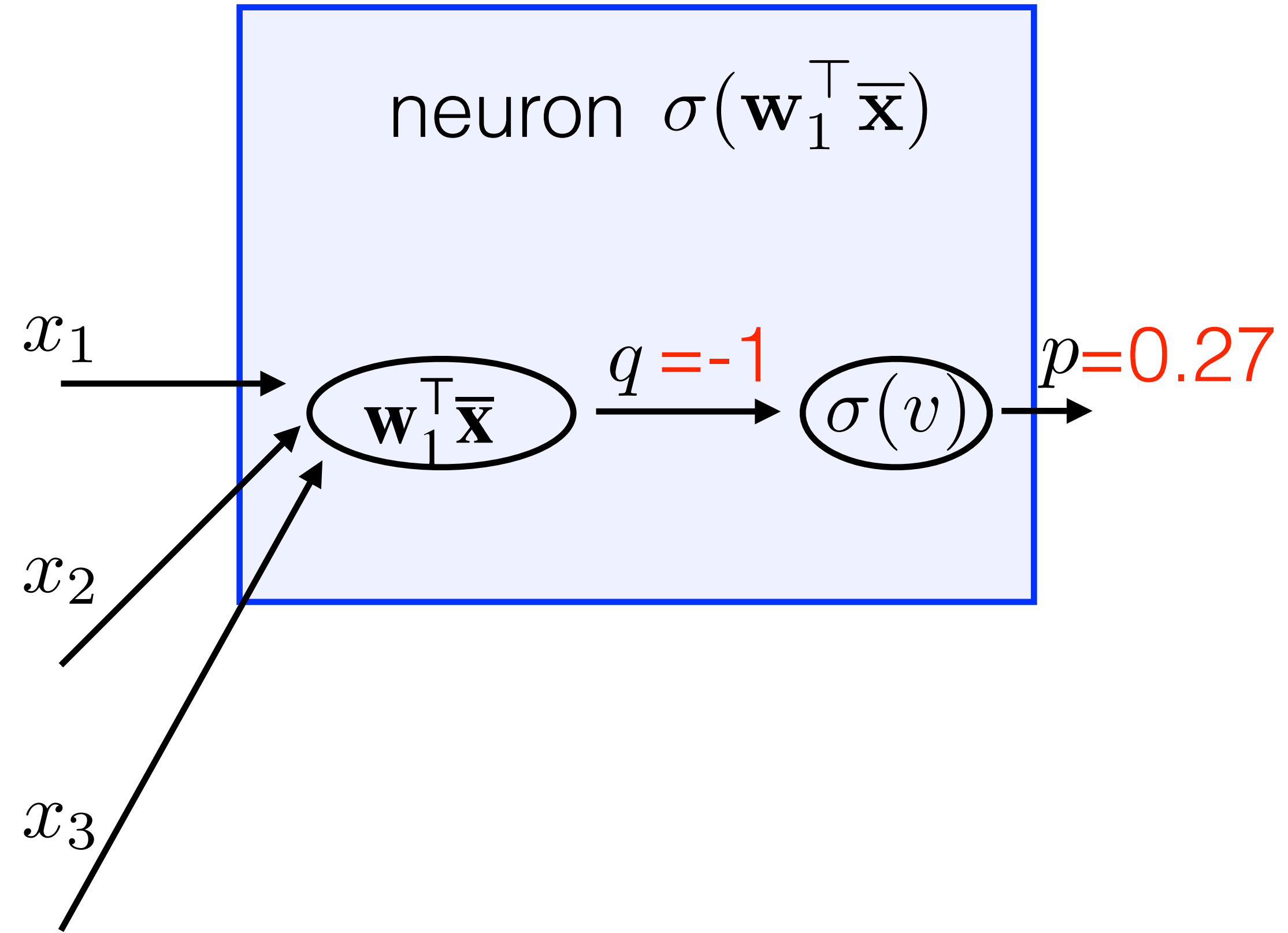
Neuron



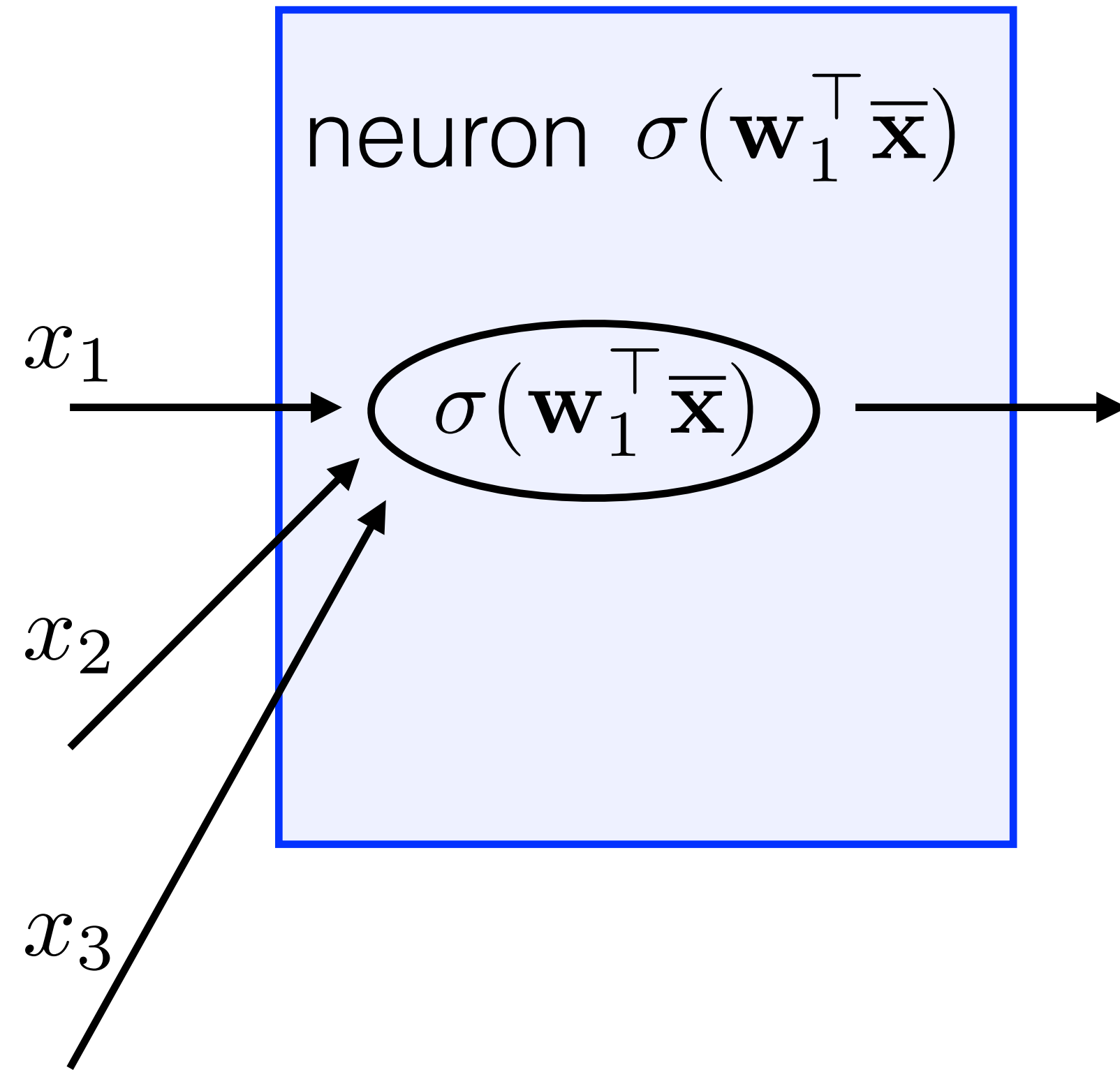
Neuron



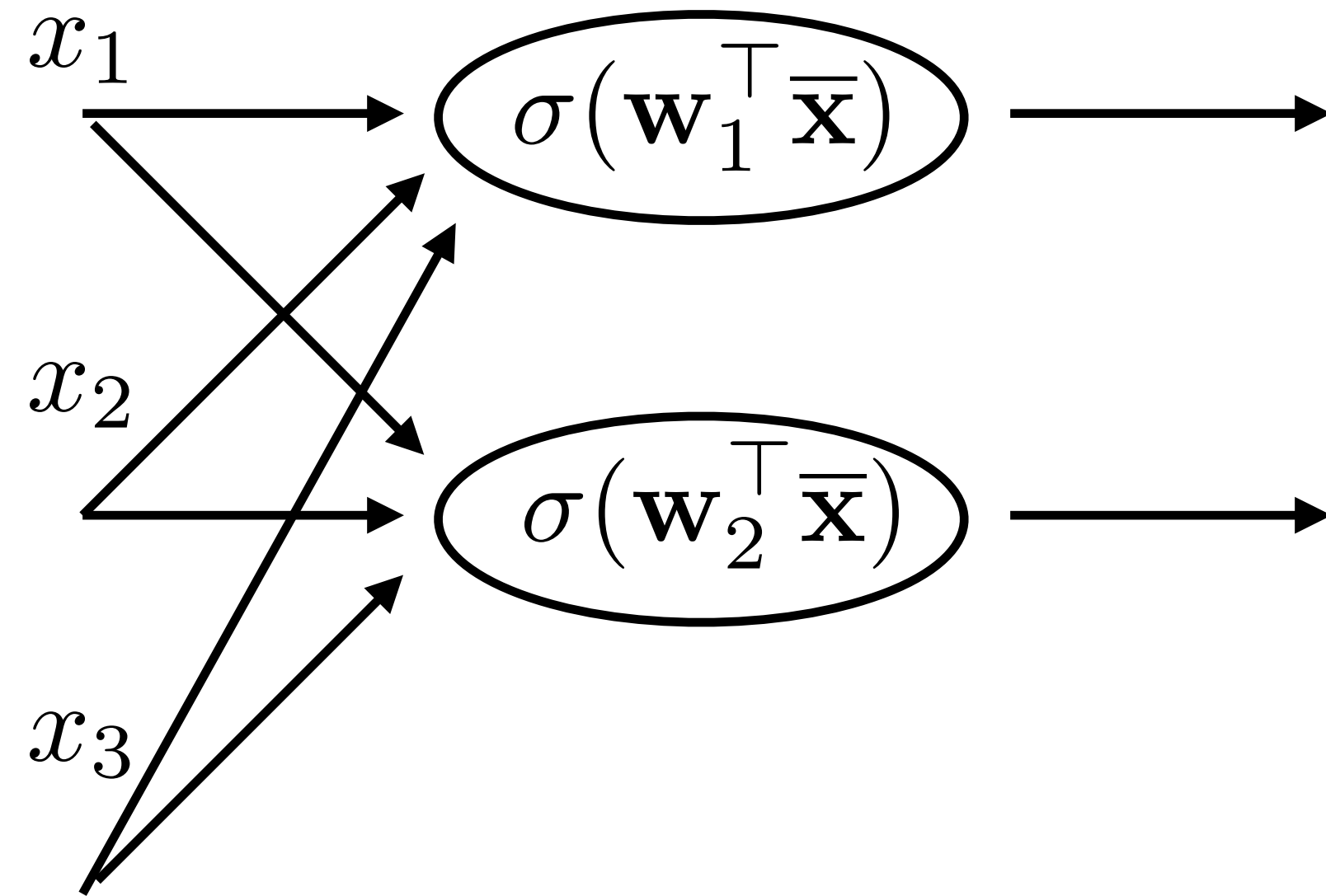
Neuron



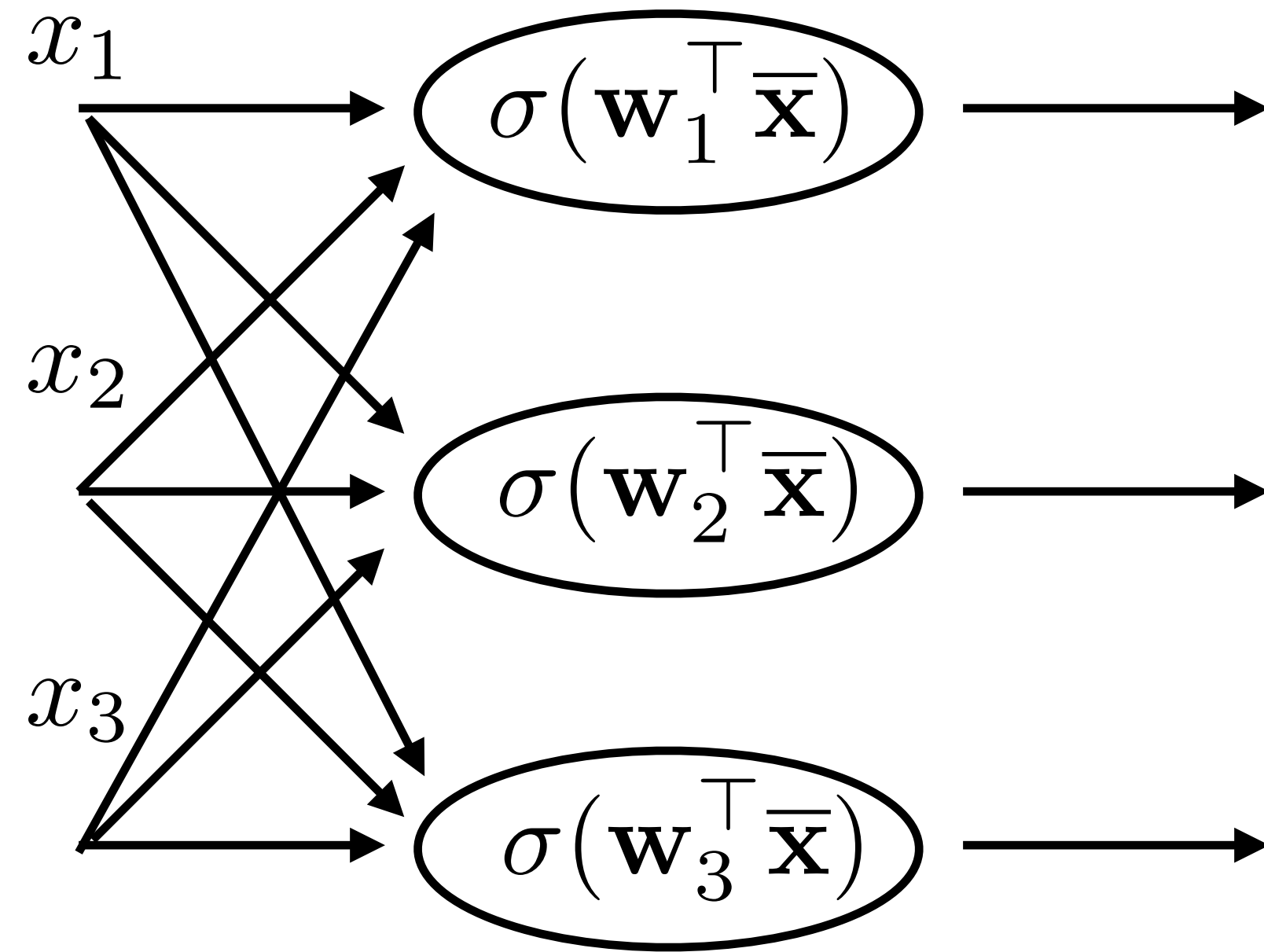
Fully-connected neural network



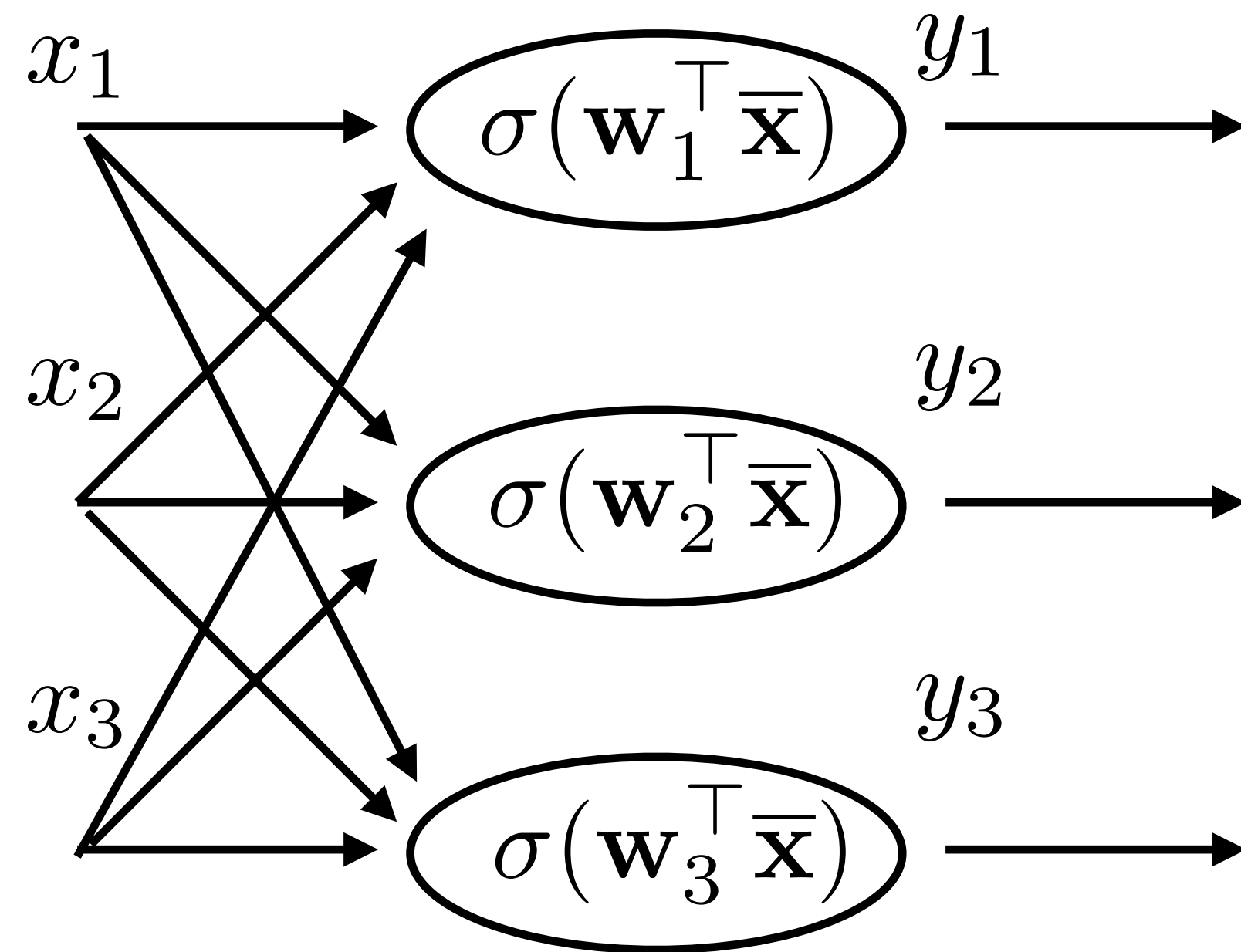
Fully-connected neural network



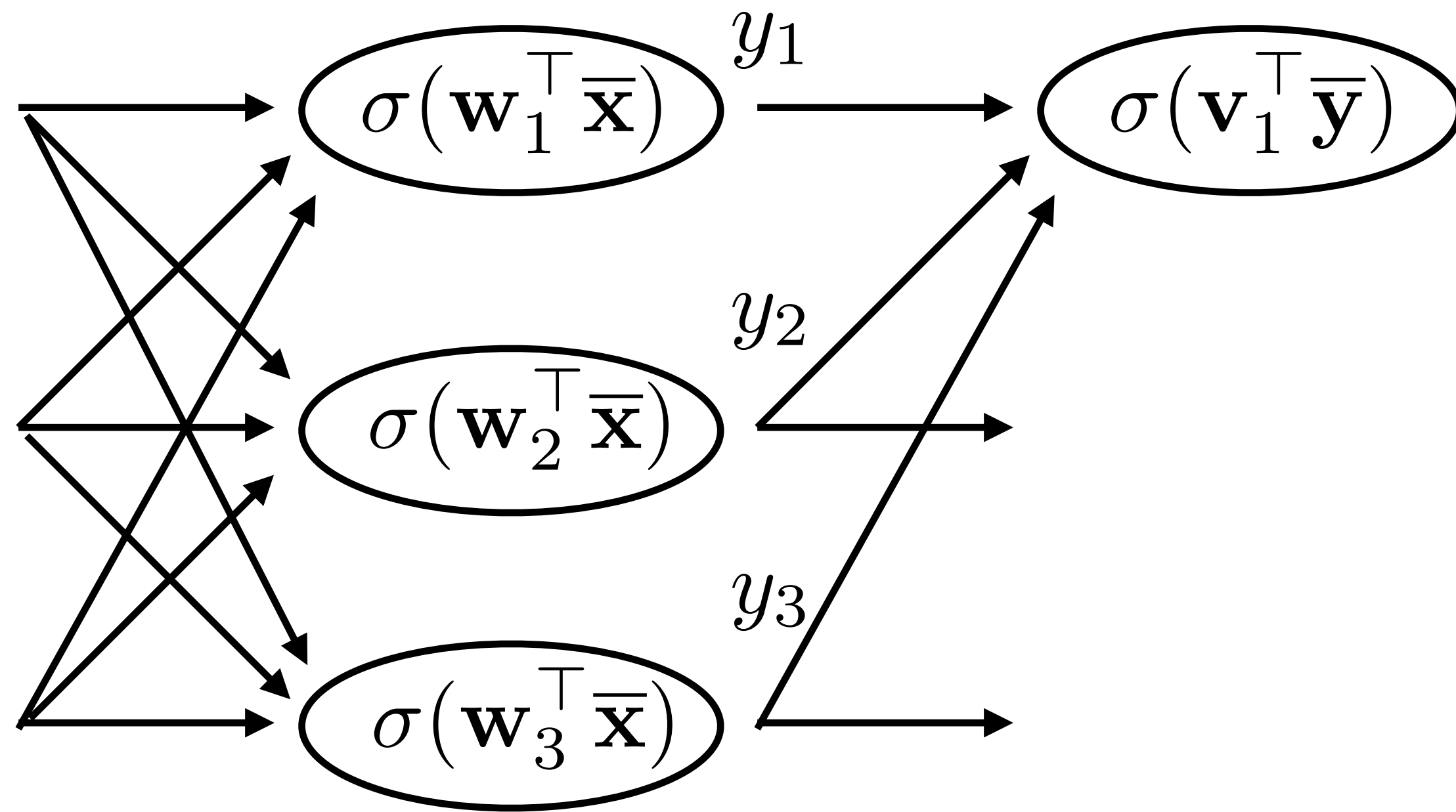
Fully-connected neural network



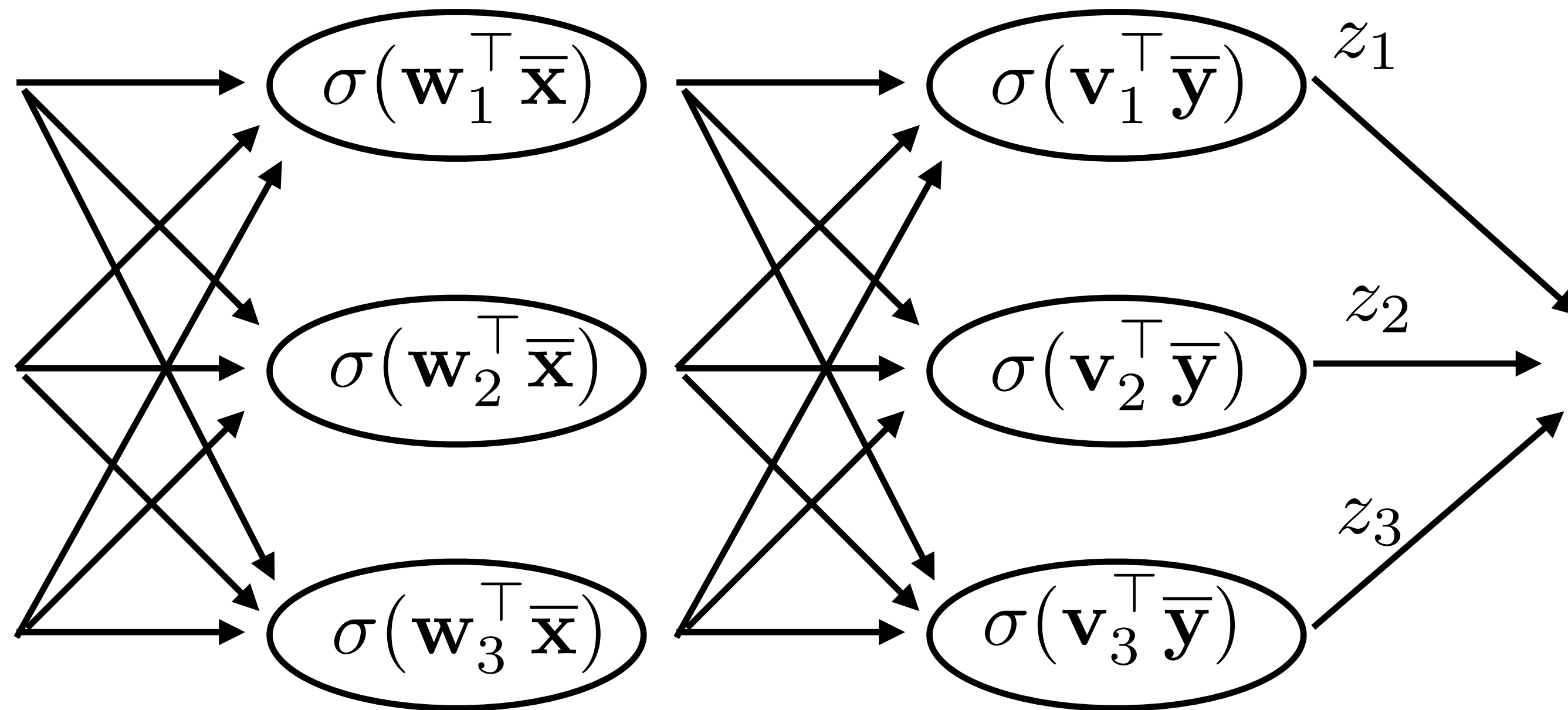
Fully-connected neural network



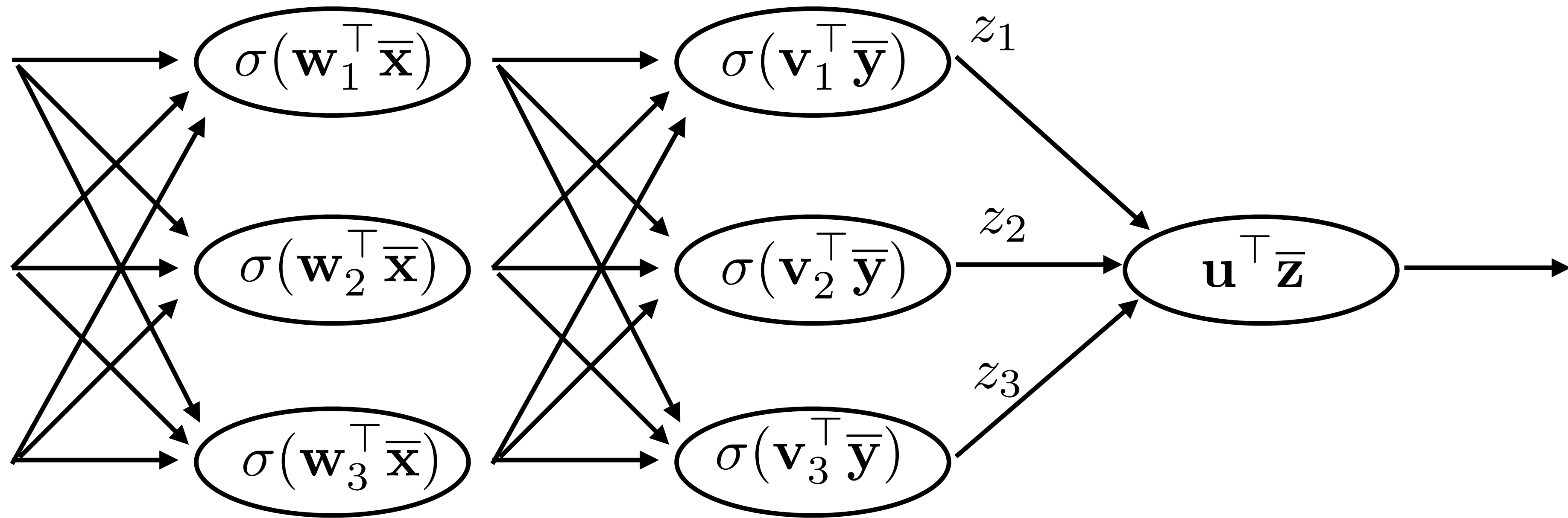
Fully-connected neural network



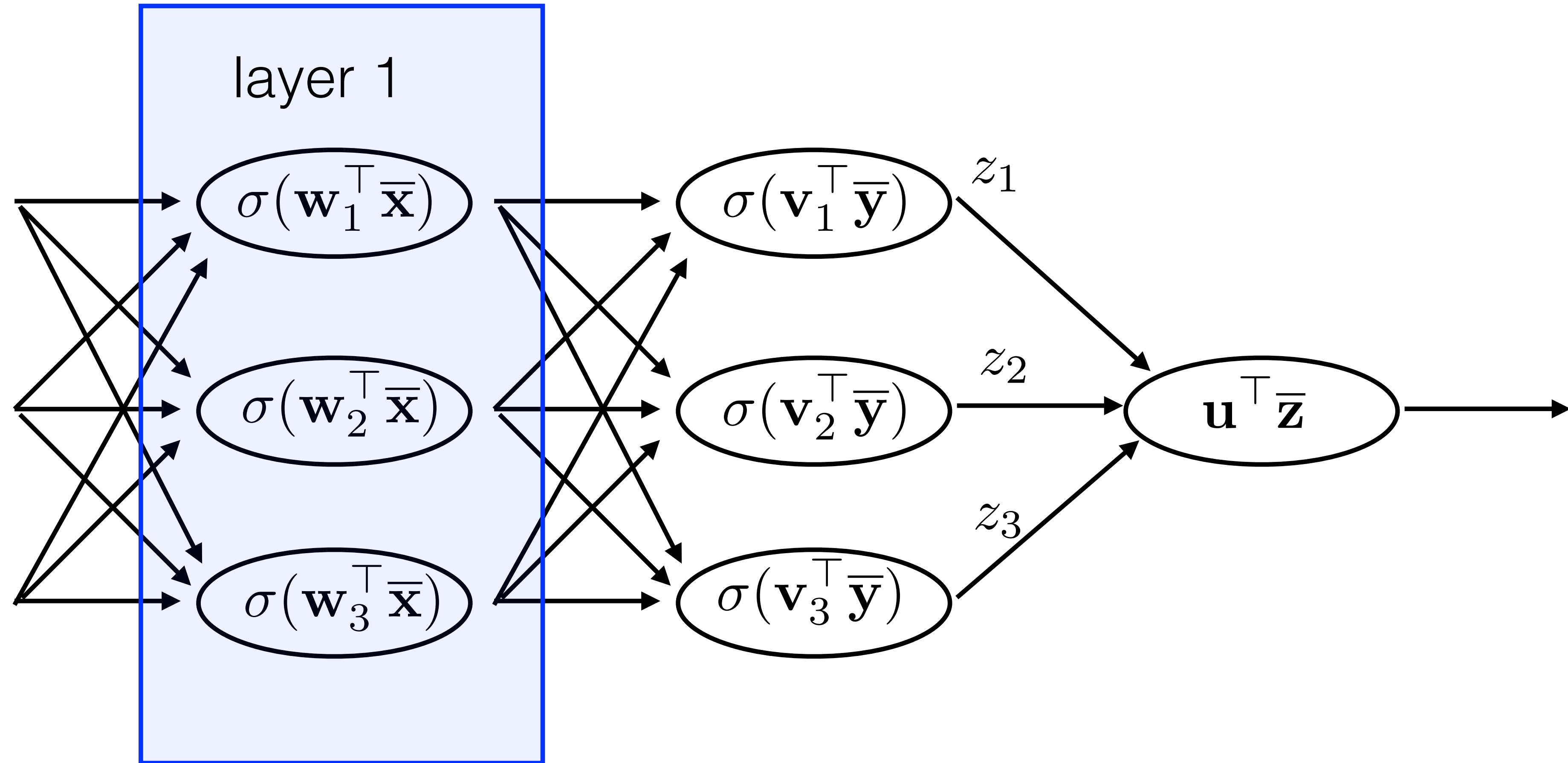
Fully-connected neural network



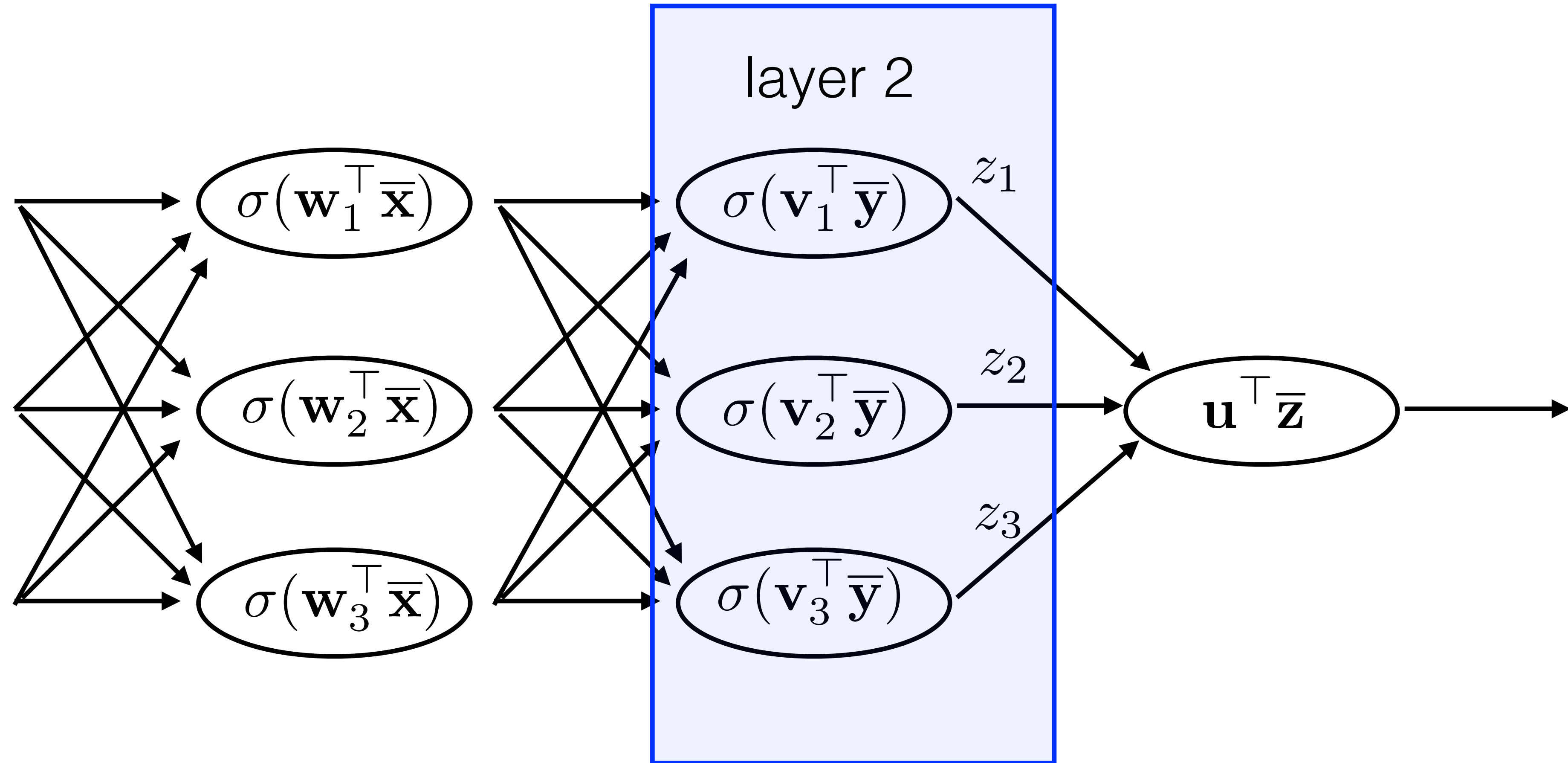
Fully-connected neural network



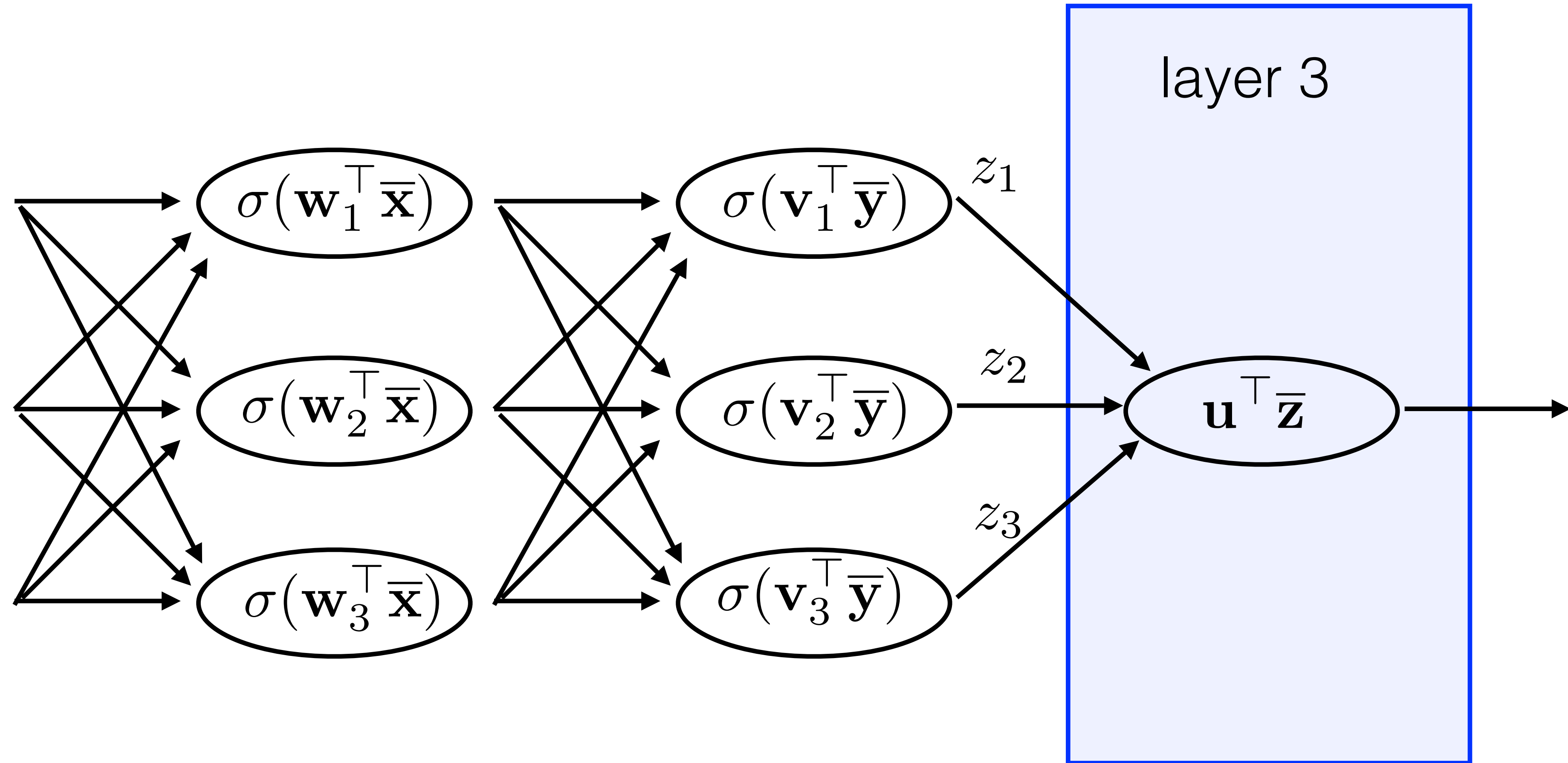
Fully-connected neural network



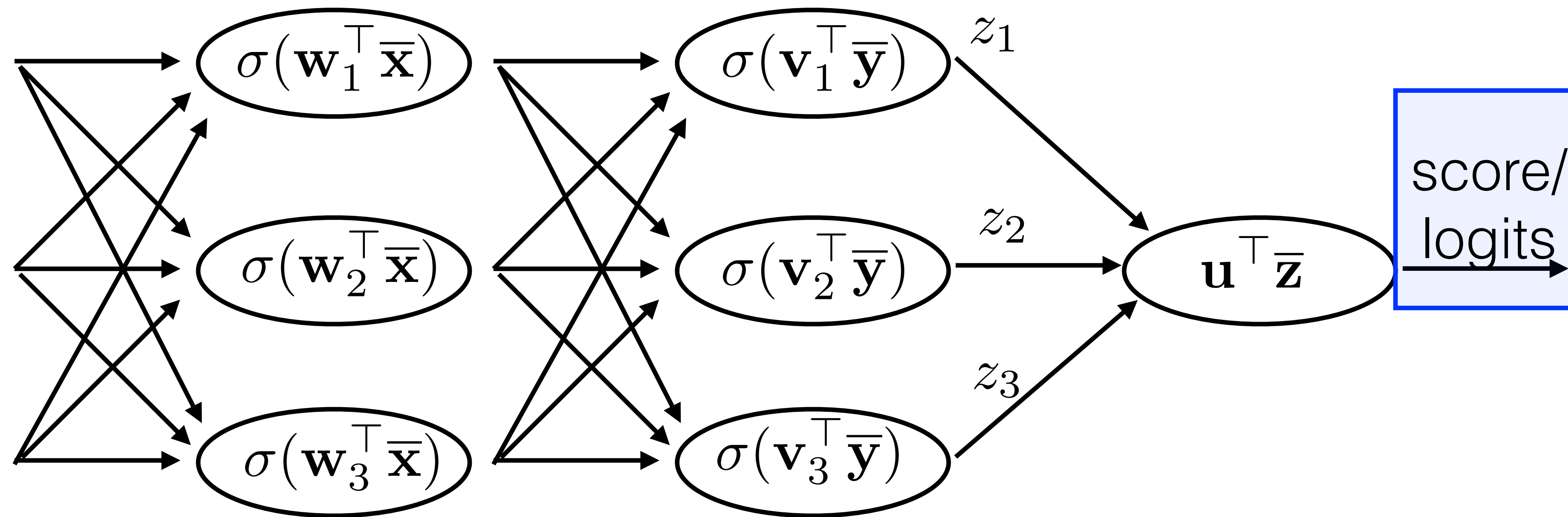
Fully-connected neural network



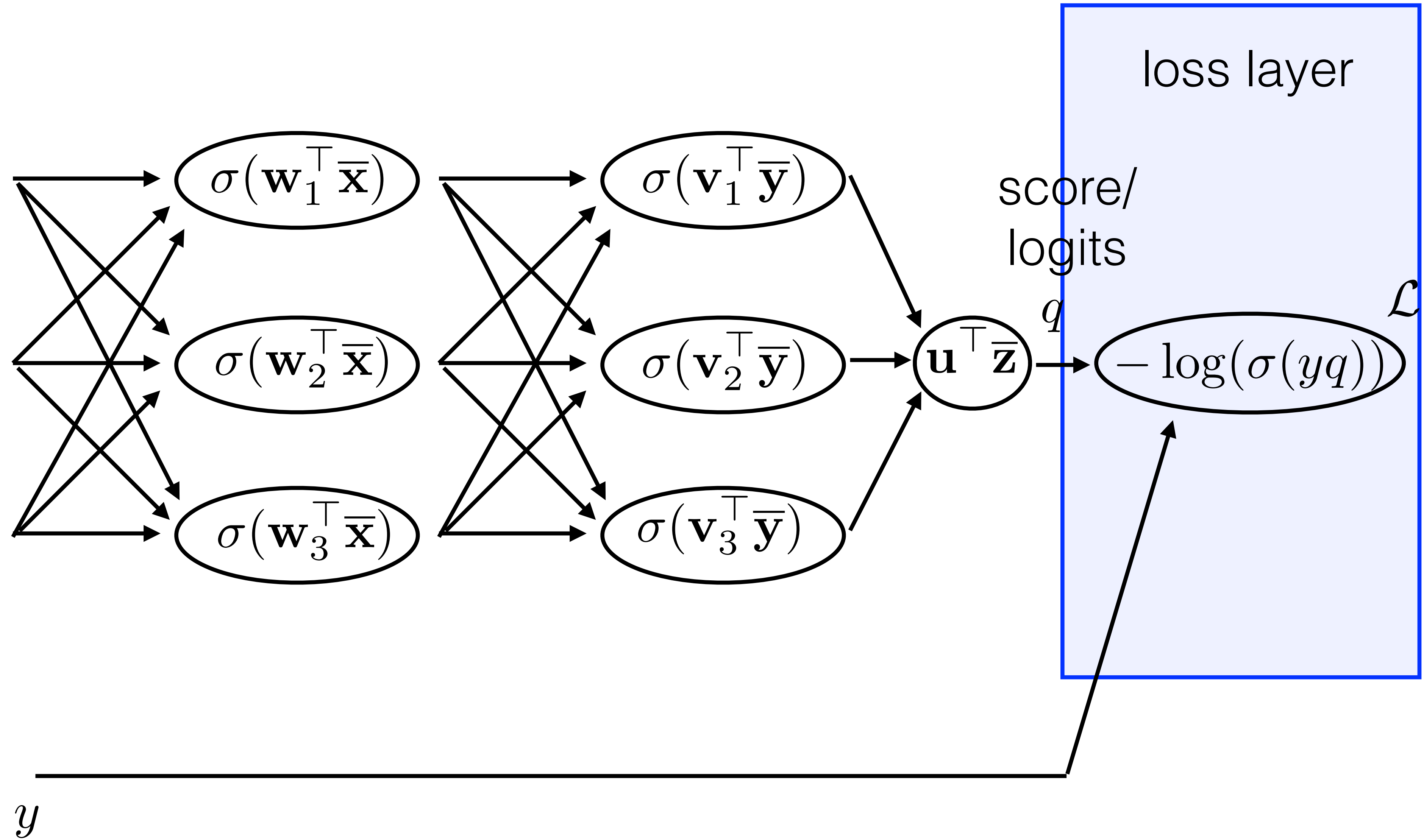
Fully-connected neural network



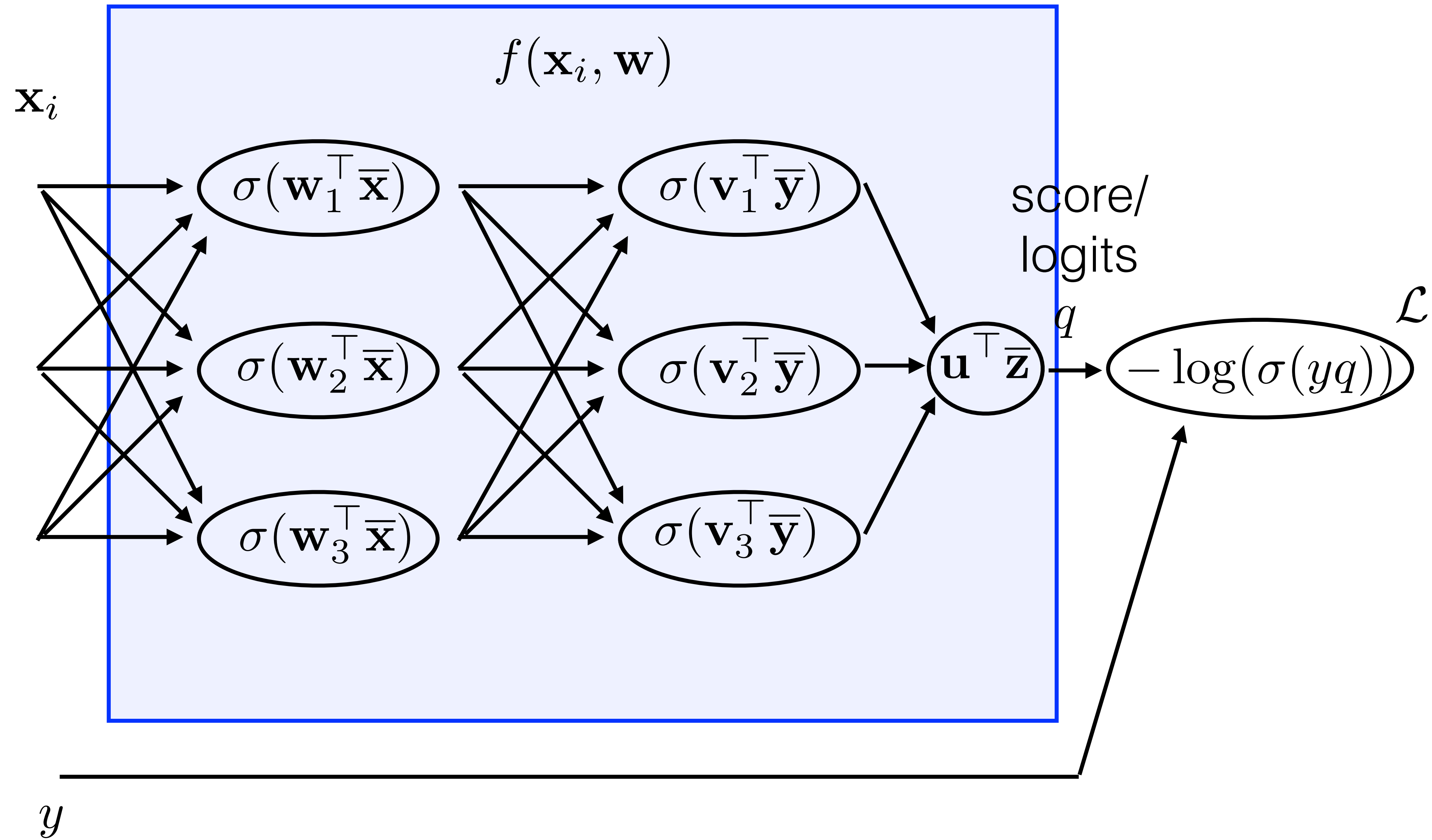
Fully-connected neural network



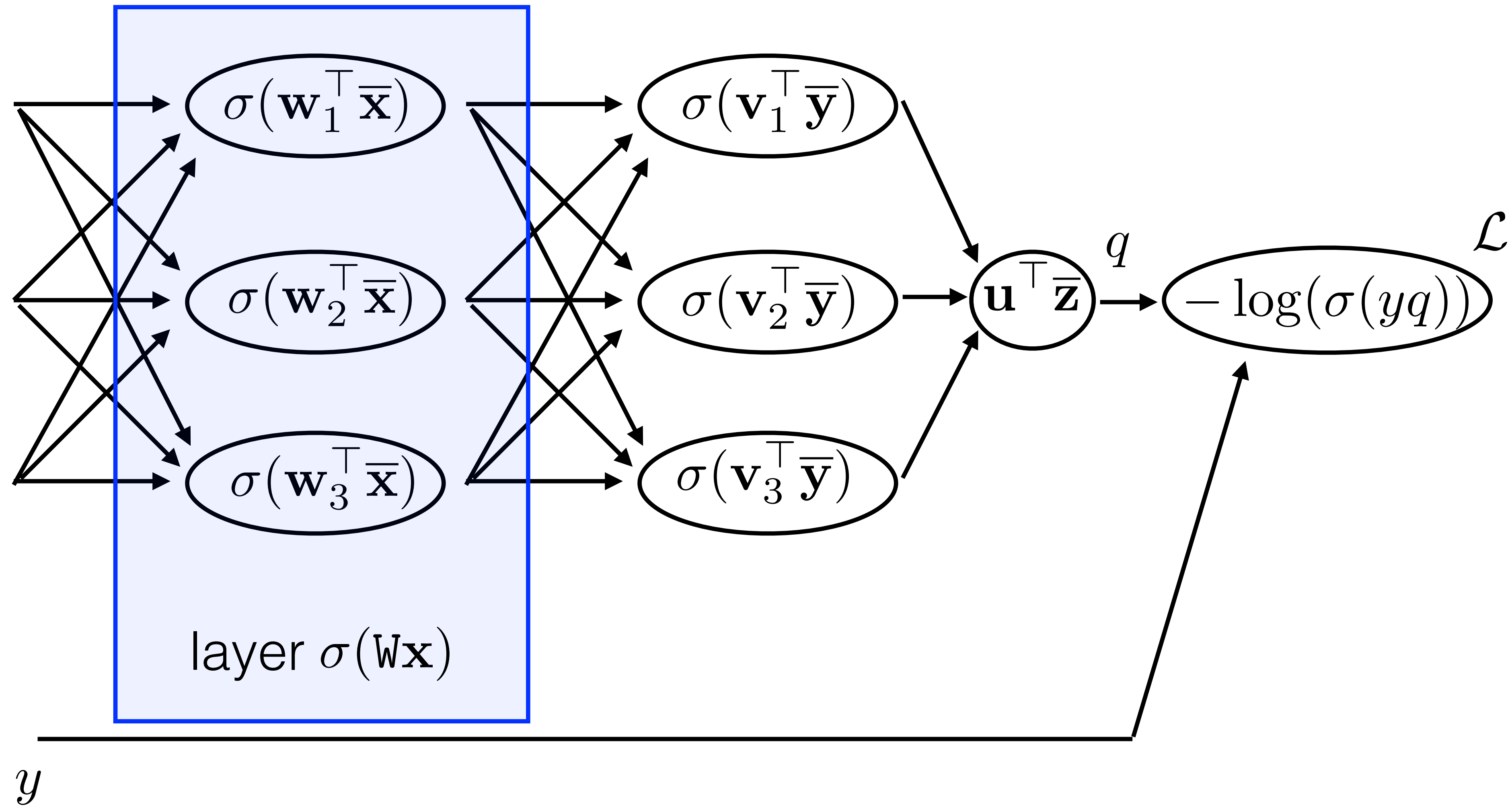
Fully-connected neural network



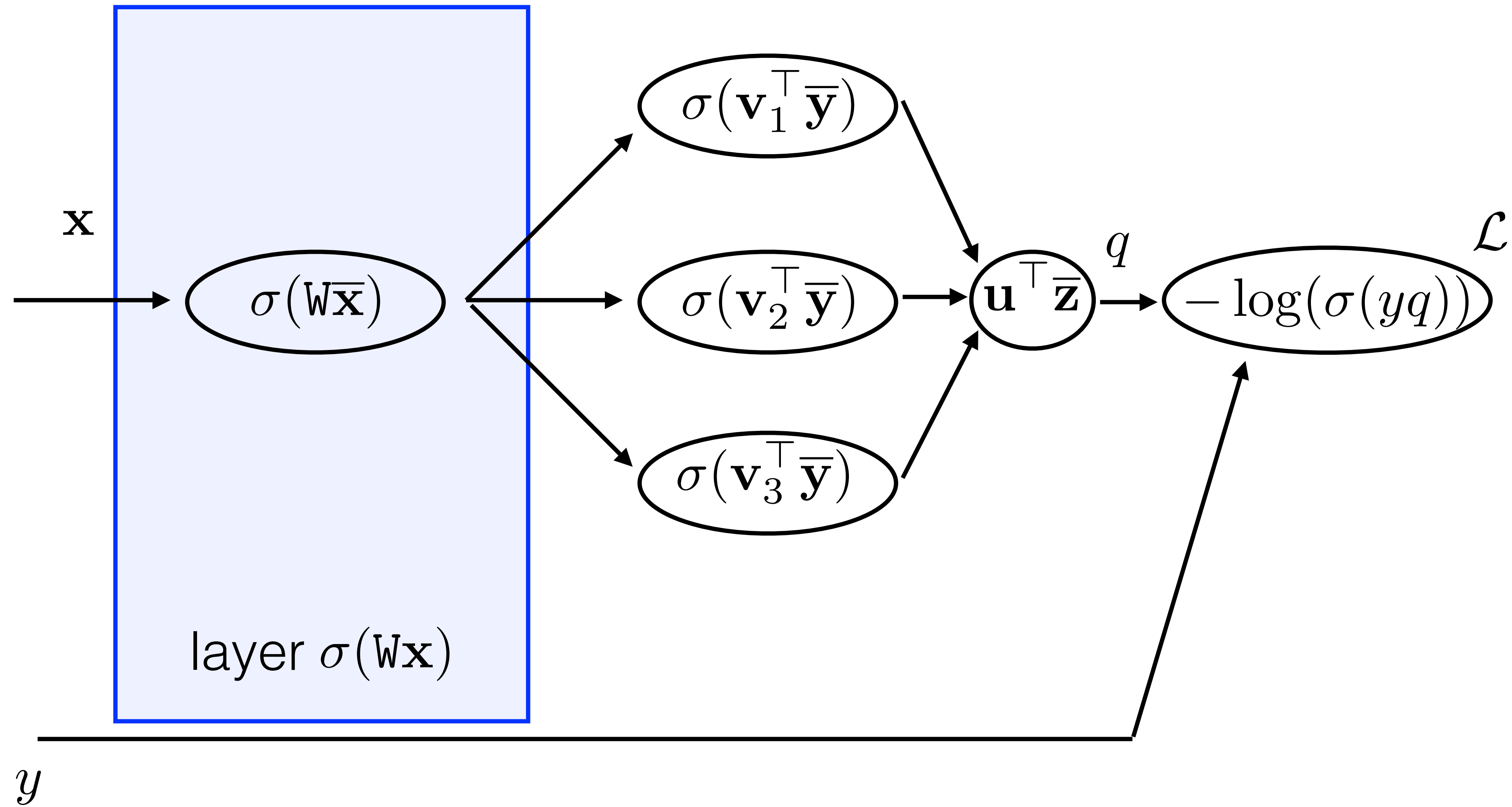
Fully-connected neural network



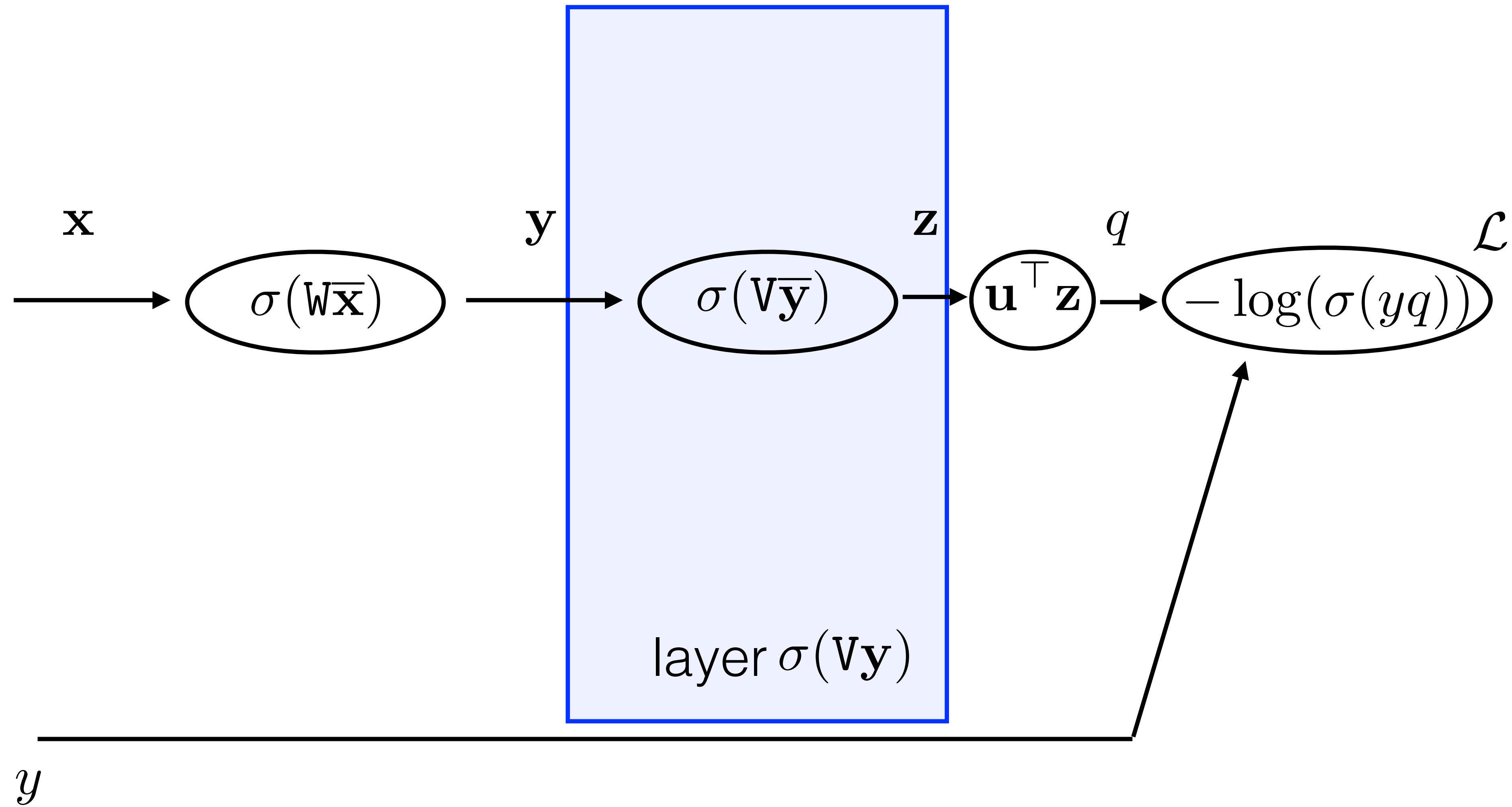
Fully-connected neural network



Fully-connected neural network



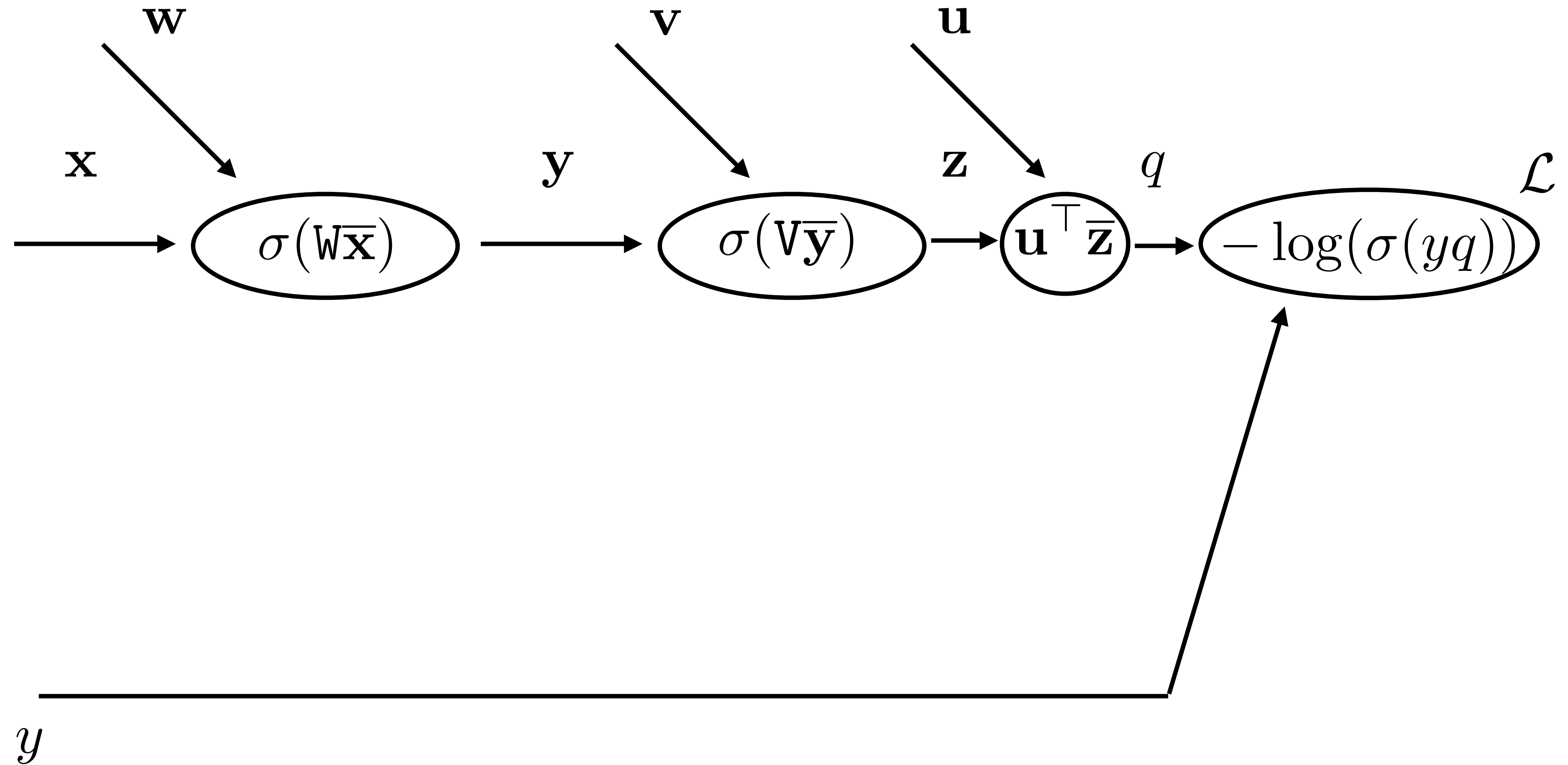
Fully-connected neural network



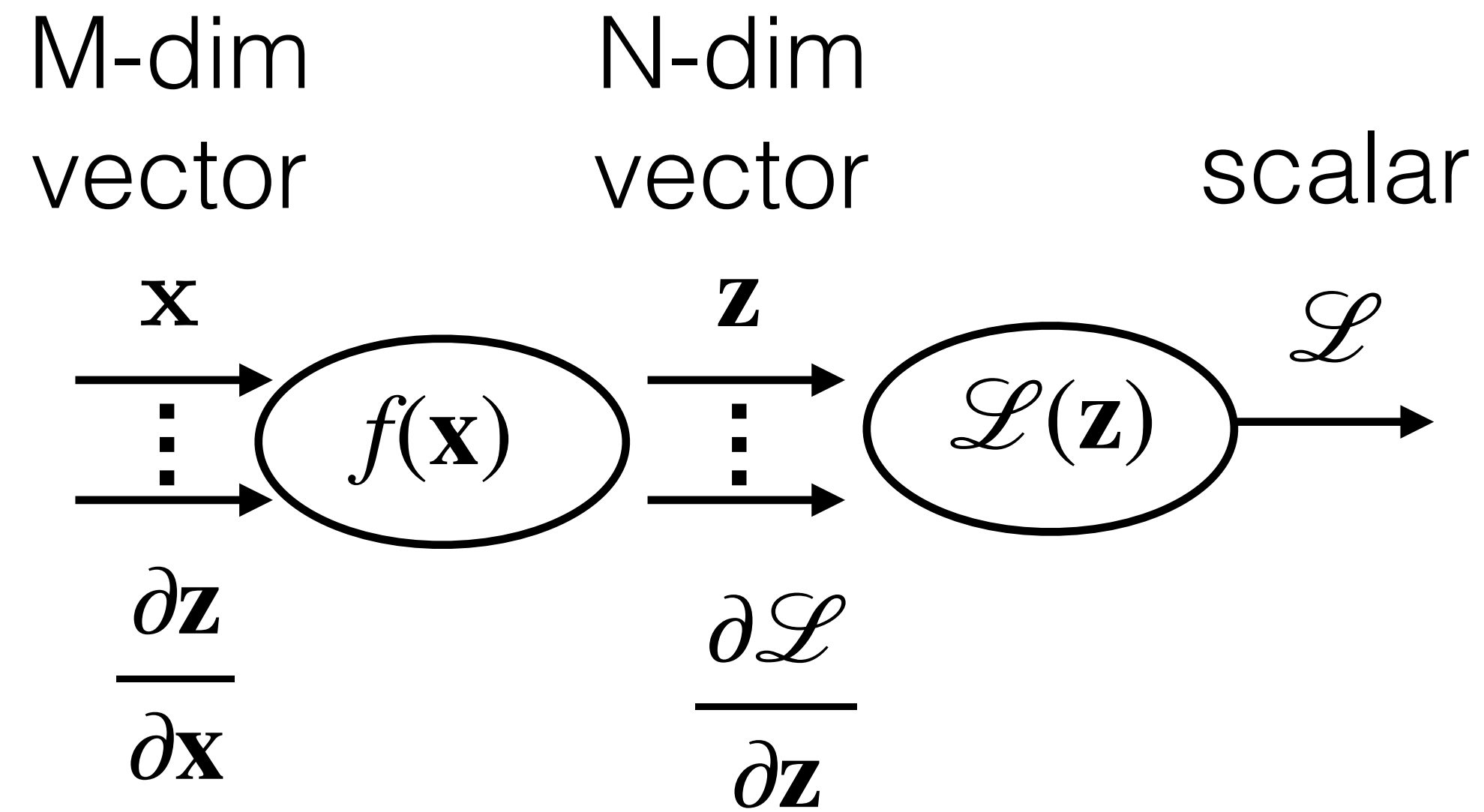
Fully-connected neural network

$$\mathbf{w} = \text{vec}(W)$$

$$\mathbf{v} = \text{vec}(V)$$



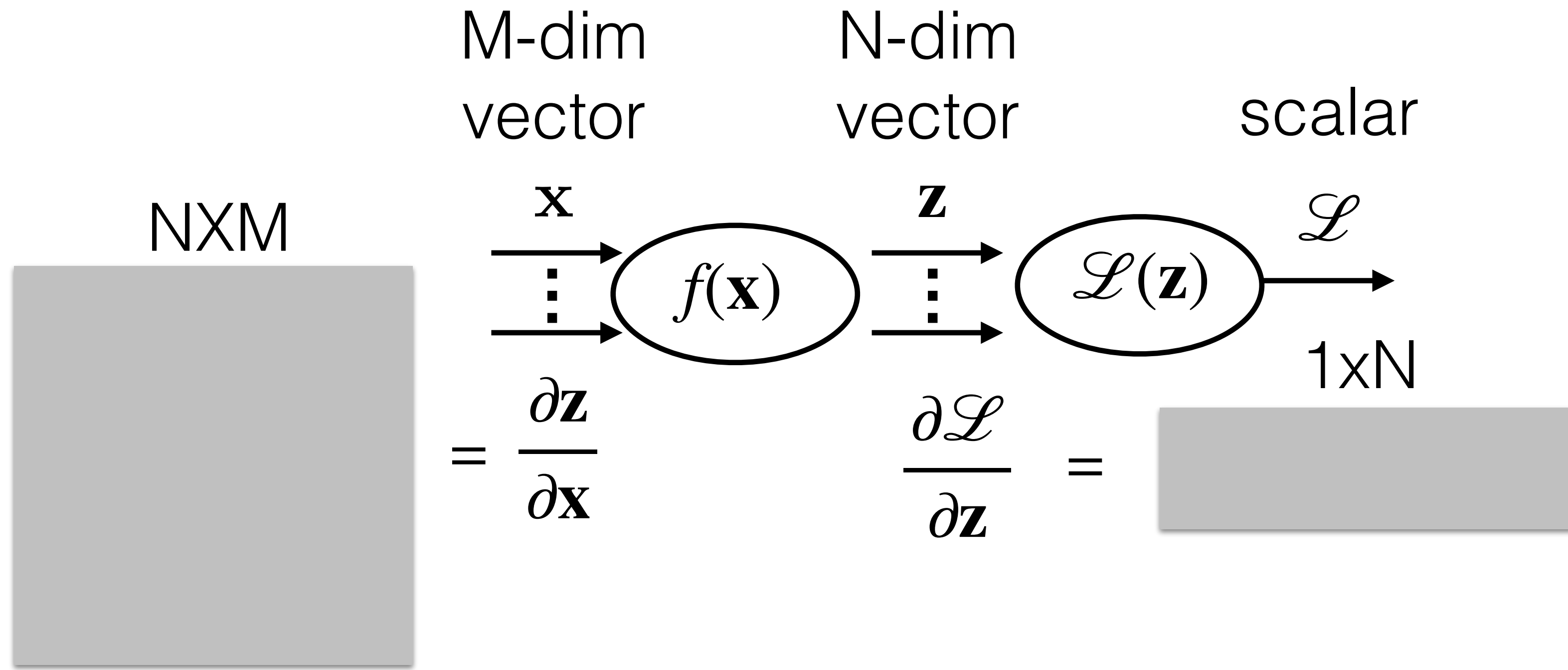
Chainrule and Jacobian



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$

Jacobian $1 \times M$ Jacobian $1 \times N$ Jacobian $N \times N$

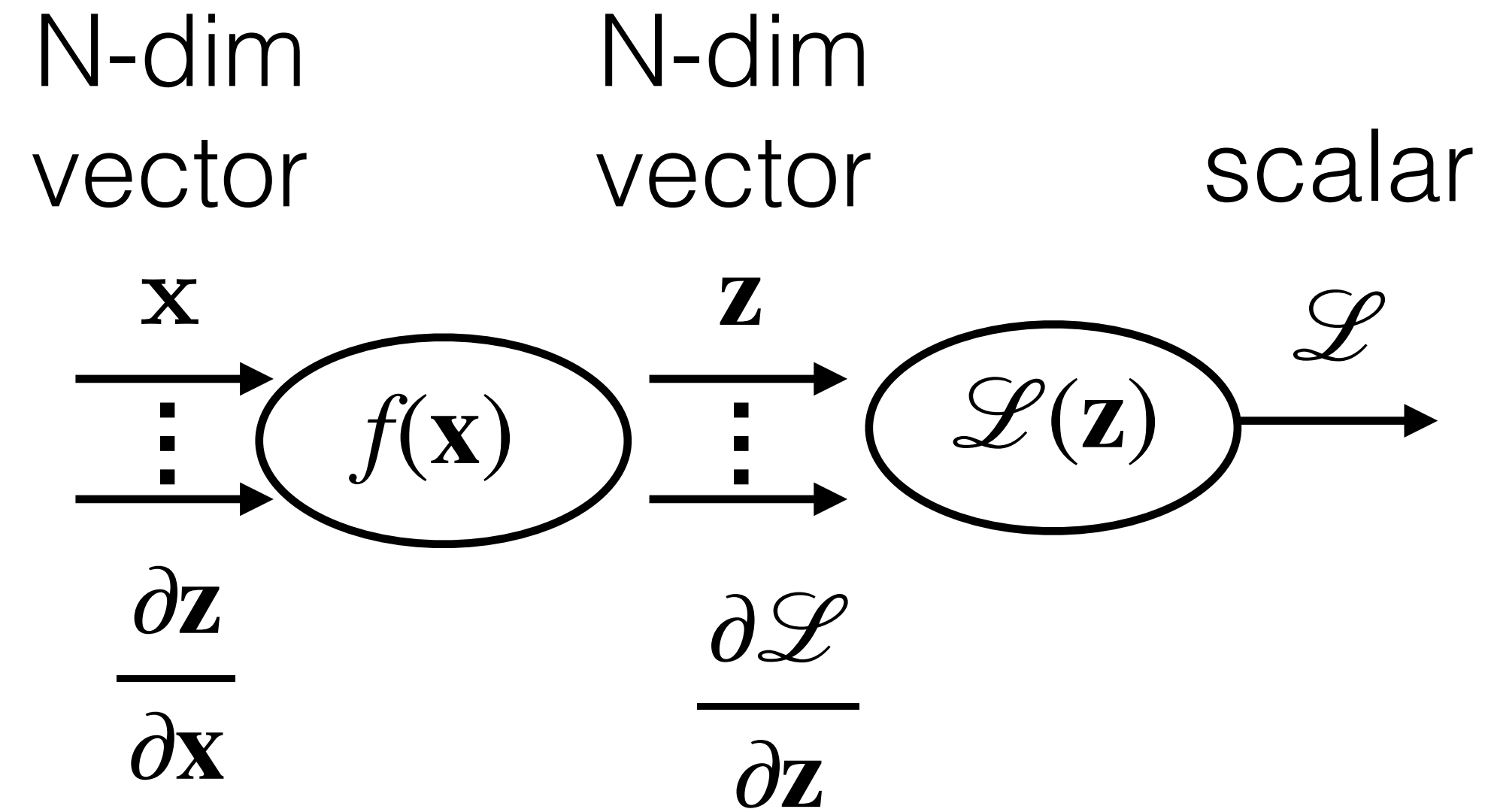
Chainrule and Jacobian



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$

Jacobian $1 \times M$ Jacobian $1 \times N$ Jacobian $N \times M$

Chainrule and Jacobian

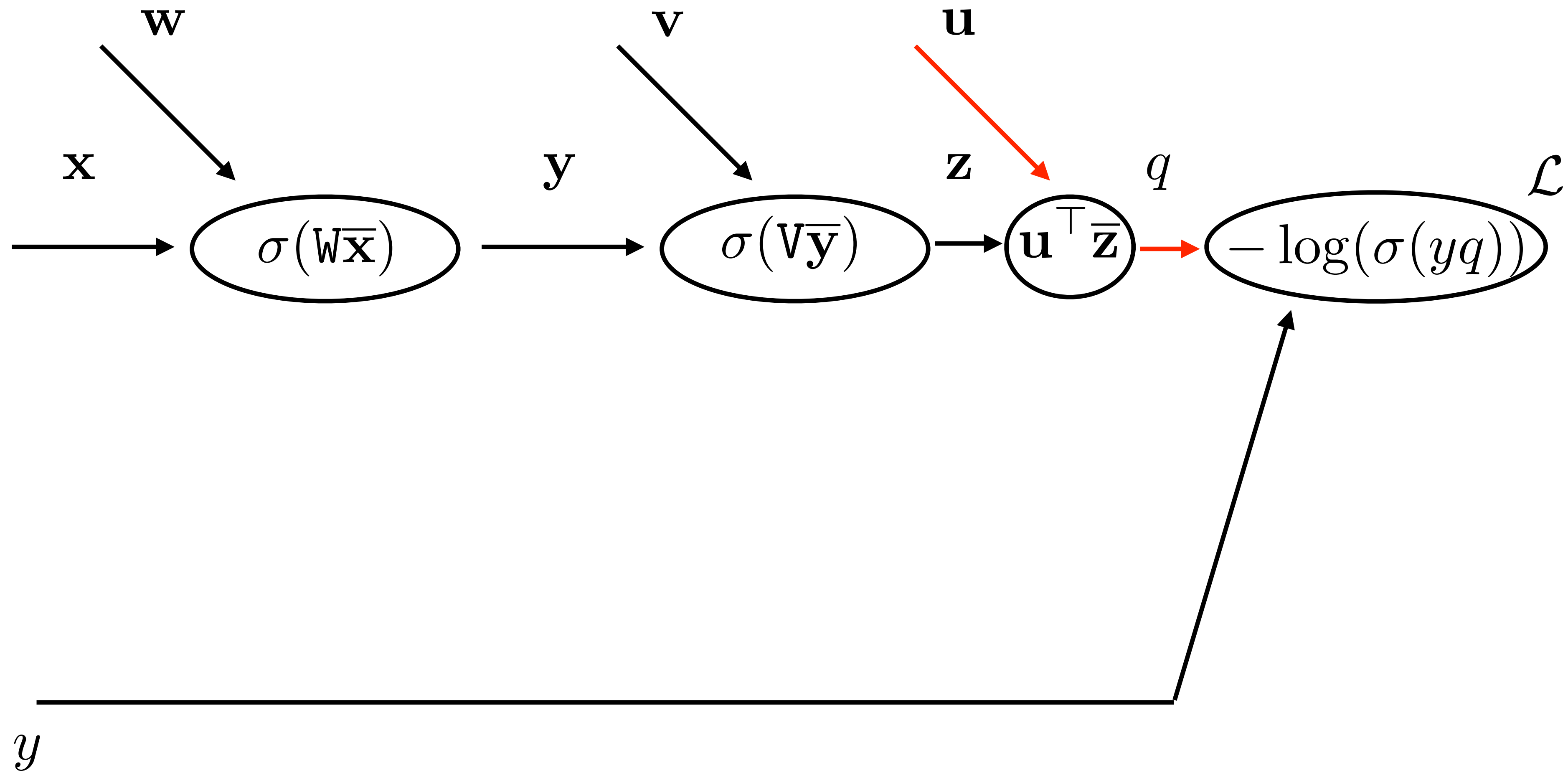


$N \times M$

$$\begin{array}{c}
 \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \\
 \text{Jacobian} \\
 1 \times M
 \end{array}
 =
 \begin{array}{c}
 1 \times N \\
 \text{Jacobian} \\
 1 \times N
 \end{array}
 \cdot
 \begin{array}{c}
 \text{Jacobian} \\
 N \times M
 \end{array}
 =
 \begin{array}{c}
 1 \times M
 \end{array}$$

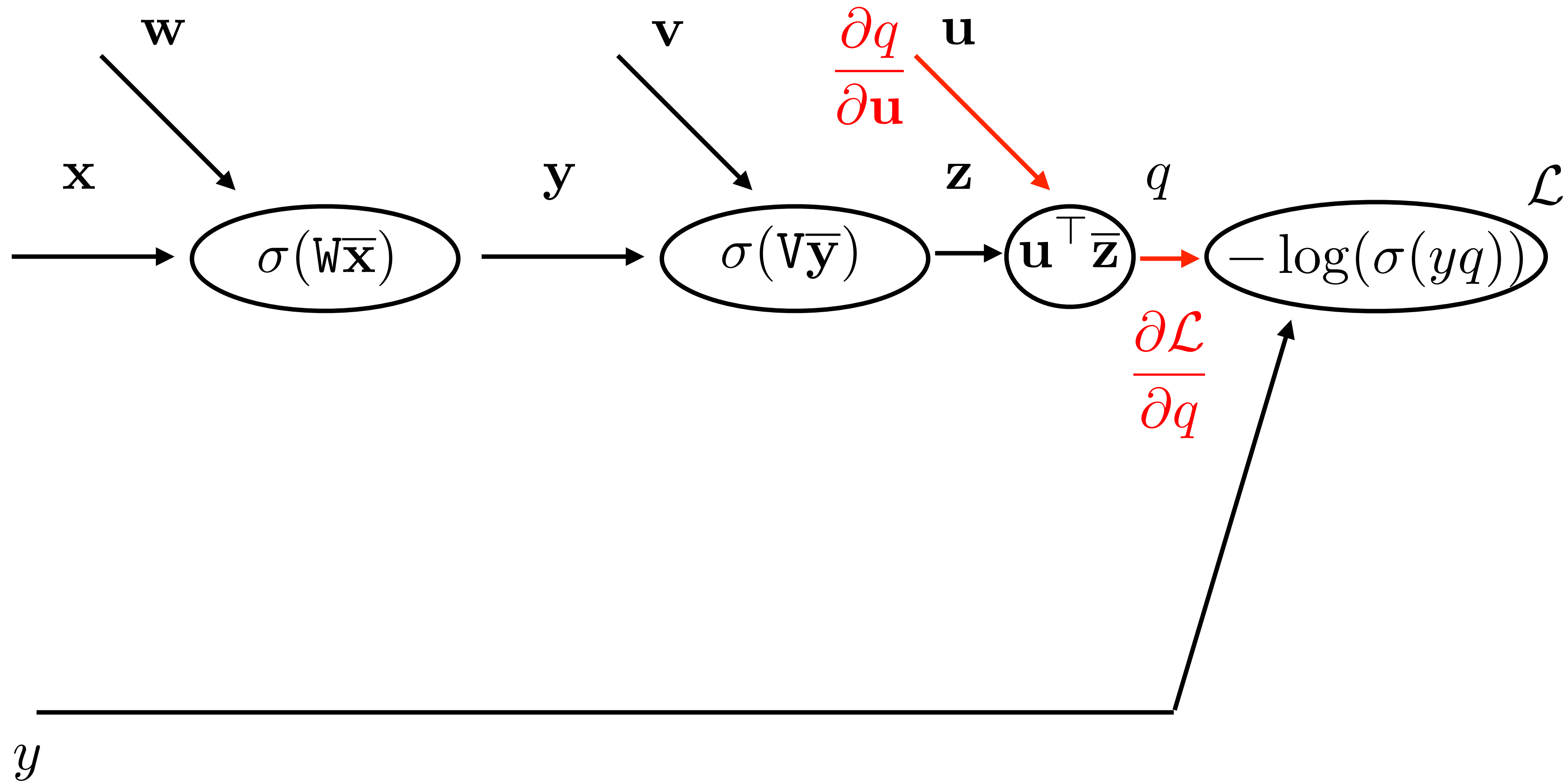
Chainrule in fully-connected neural network

Derivative wrt \mathbf{u} : $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = ?$

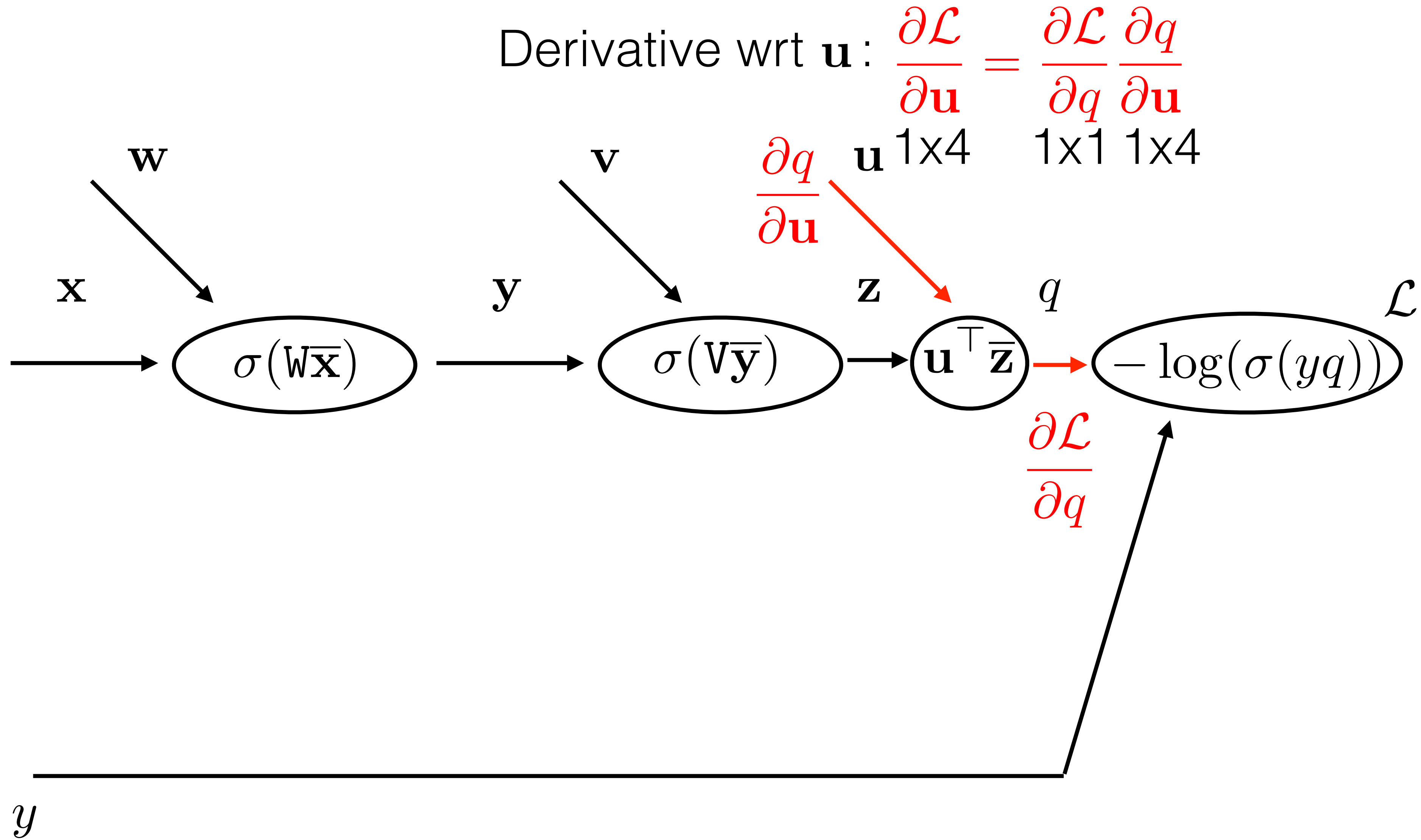


Chainrule in fully-connected neural network

Derivative wrt \mathbf{u} : $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{u}}$

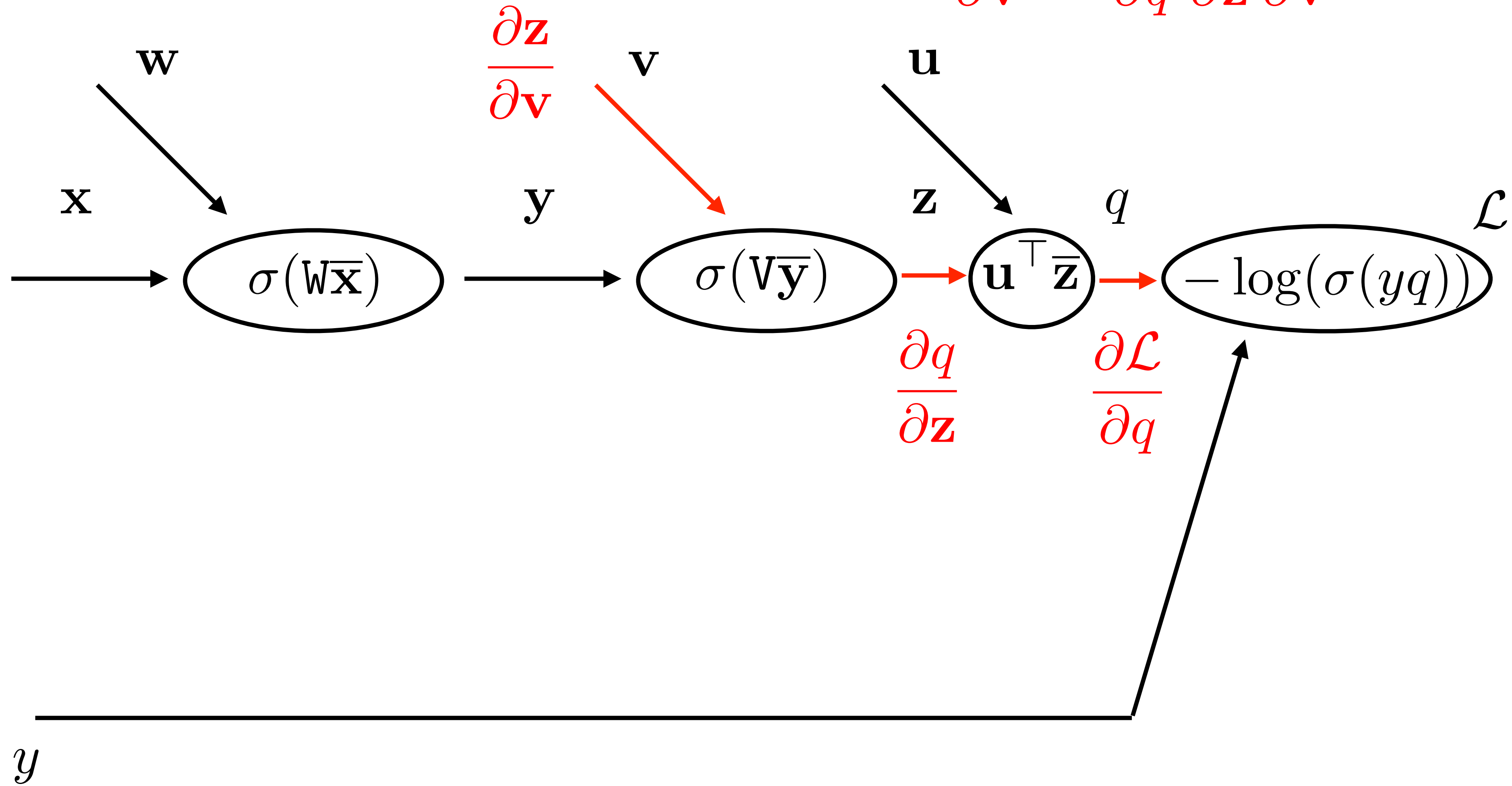


Chainrule in fully-connected neural network

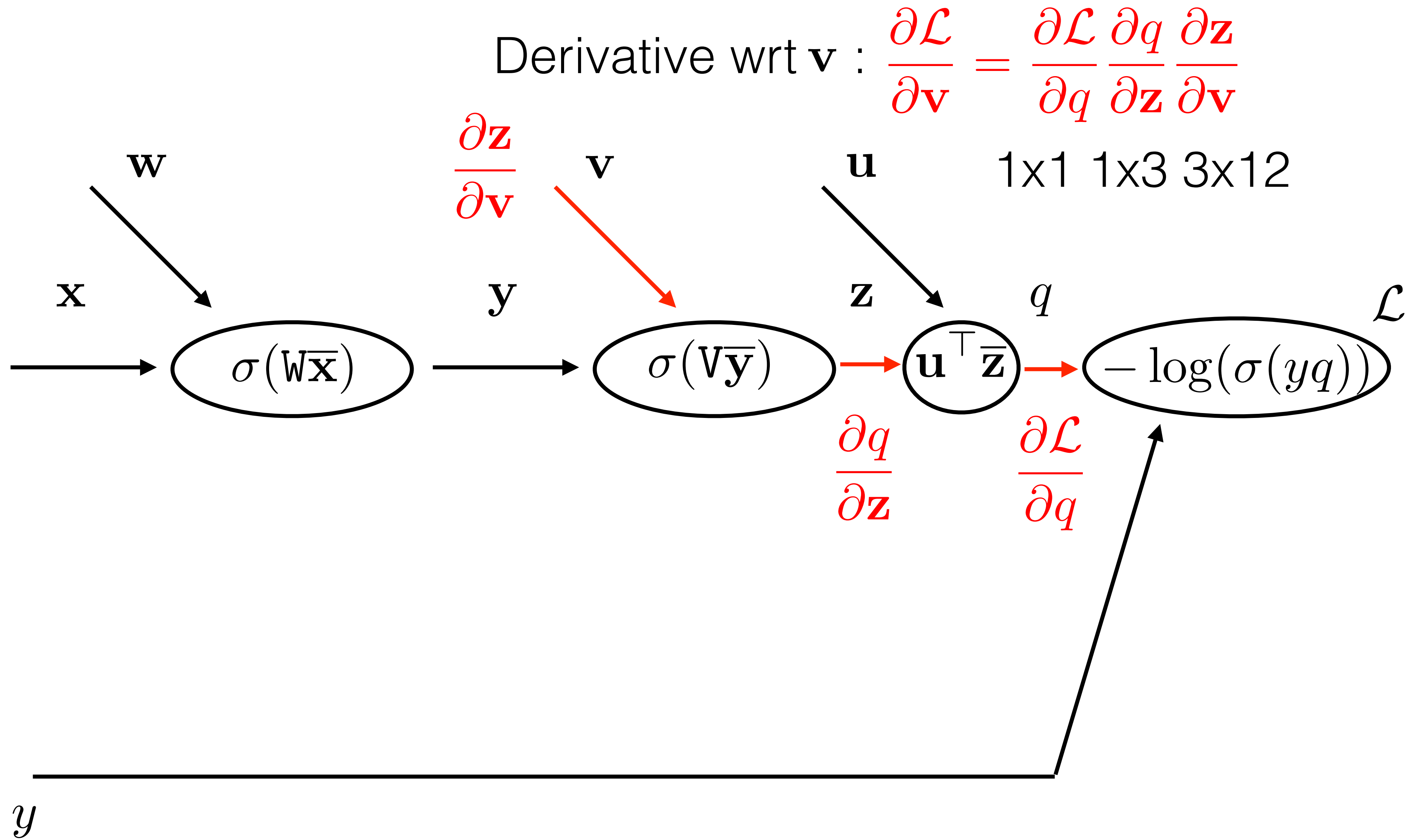


Chainrule in fully-connected neural network

Derivative wrt \mathbf{v} : $\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{v}}$

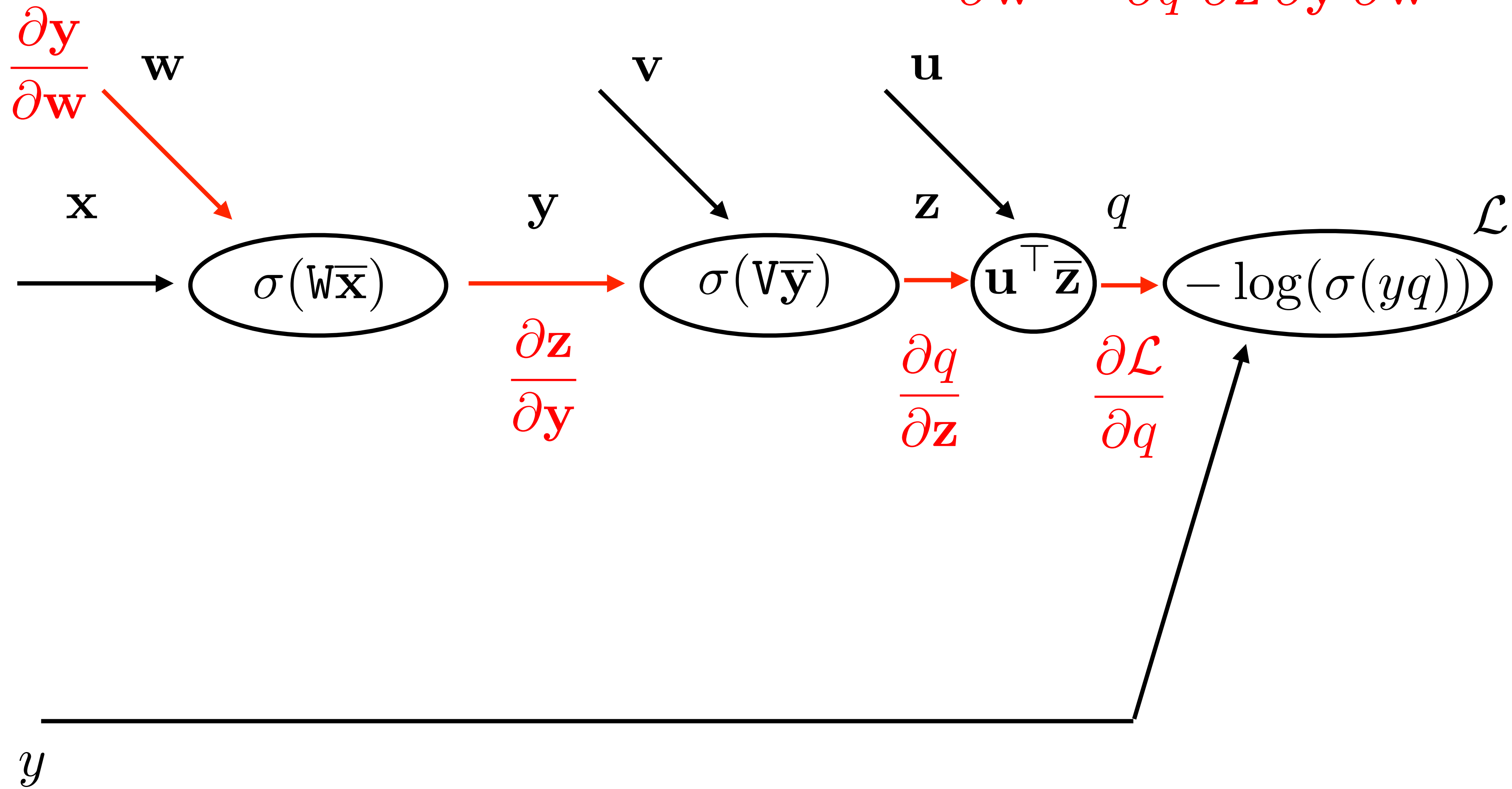


Chainrule in fully-connected neural network

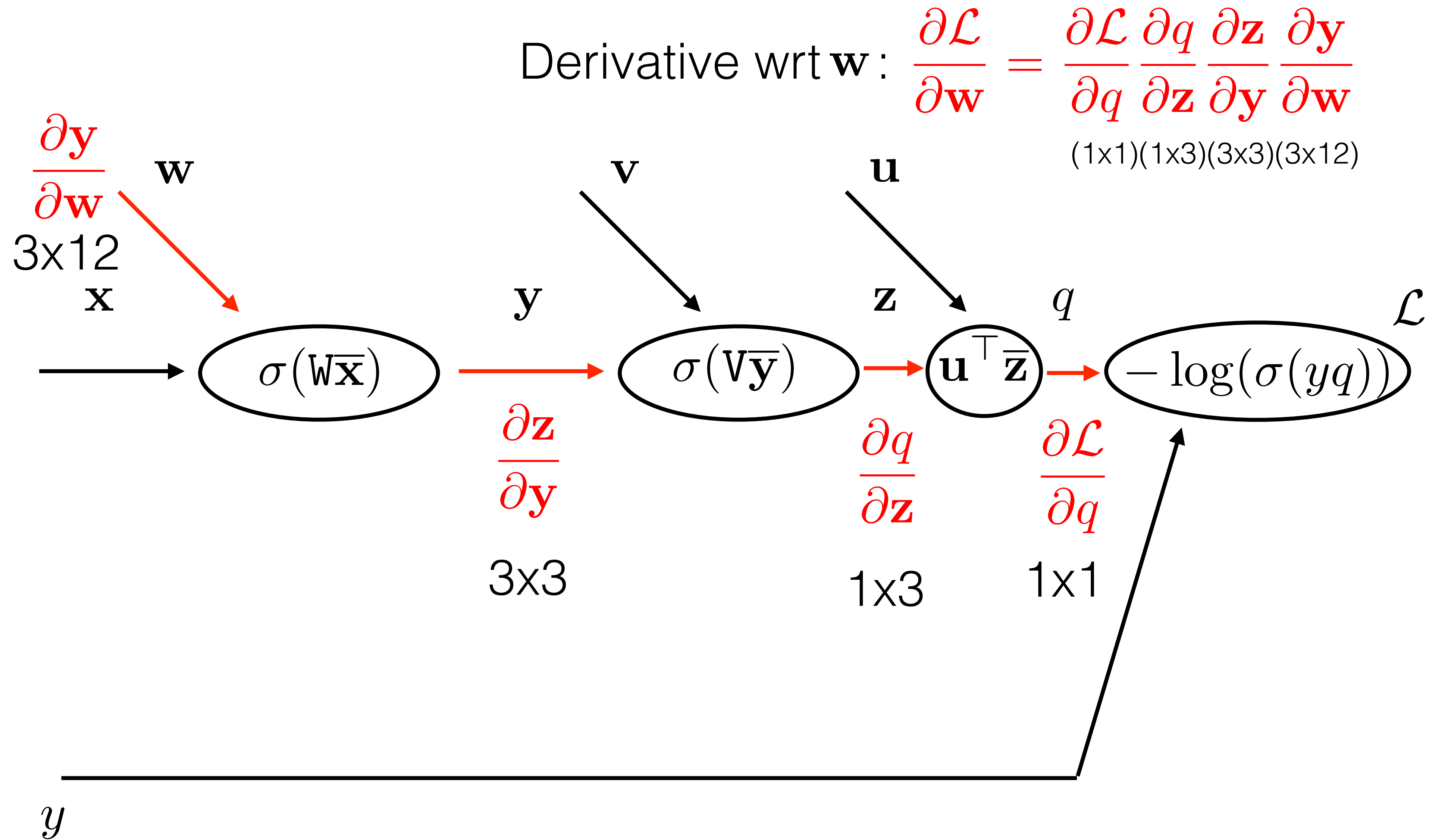


Chainrule in fully-connected neural network

Derivative wrt \mathbf{w} :
$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{w}}$$



Chainrule in fully-connected neural network



Chainrule in fully-connected neural network

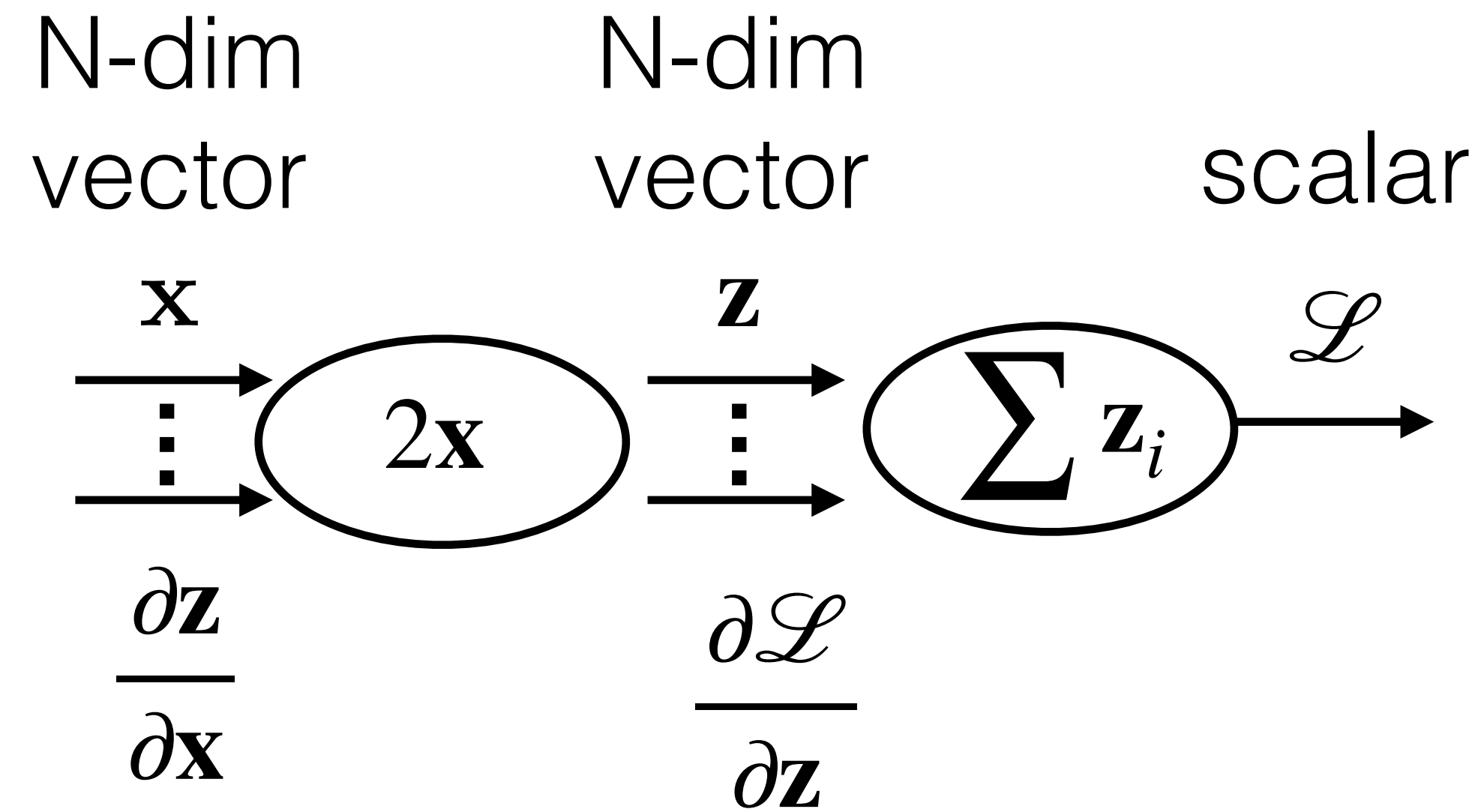
1. Estimate all required jacobians
2. Update weights:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{u}} \quad \mathbf{u} = \mathbf{u} - \alpha \left[\frac{\partial \mathcal{L}}{\partial \mathbf{u}} \right]^\top$$
$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{v}} \quad \mathbf{v} = \mathbf{v} - \alpha \left[\frac{\partial \mathcal{L}}{\partial \mathbf{v}} \right]^\top$$
$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{w}} \quad \mathbf{w} = \mathbf{w} - \alpha \left[\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right]^\top$$

3. Optionally update learning rate α
4. Repeat until convergence

Particular gradients can be obtained by multiplying jacobians, however it is inefficient => slightly different implementation

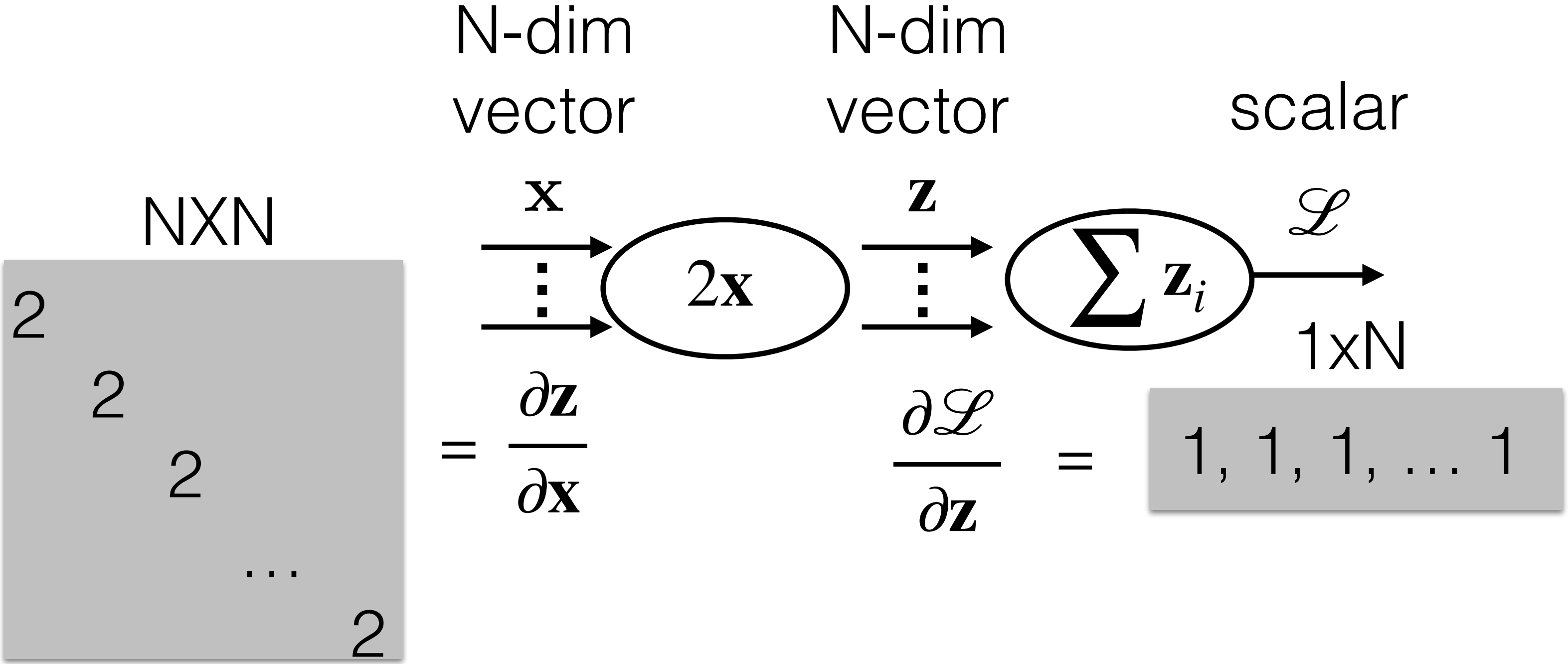
Backprop in Pytorch avoids Jacobians



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$

Jacobian Jacobian Jacobian
 $1 \times N$ $1 \times N$ $N \times N$

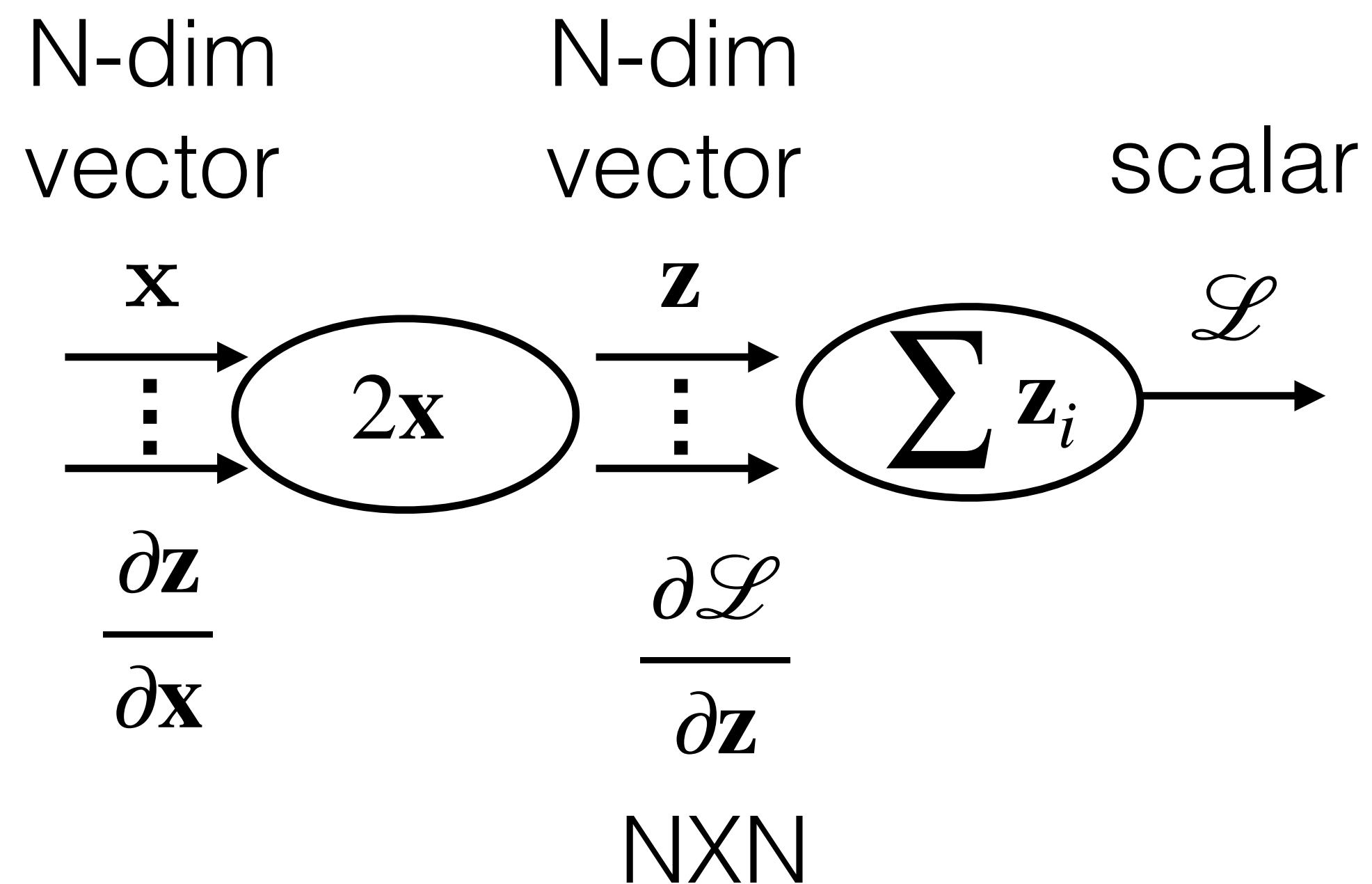
Backprop in Pytorch avoids Jacobians



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$

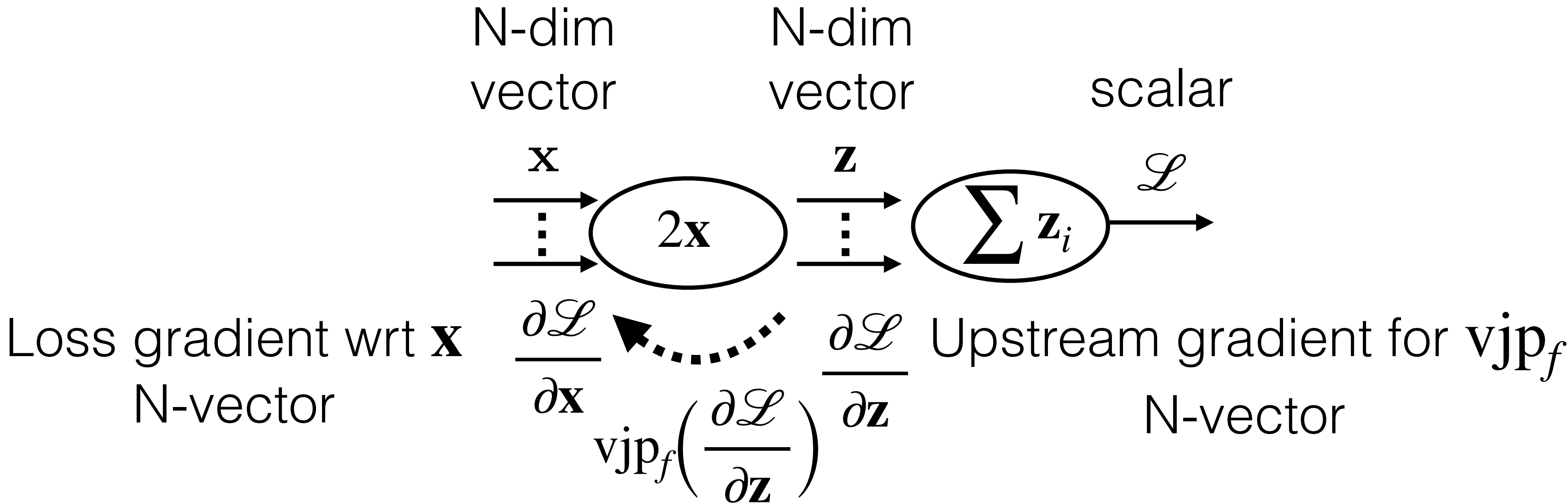
Jacobian $1 \times N$ Jacobian $1 \times N$ Jacobian $N \times N$

Backprop in Pytorch avoids Jacobians



$$\begin{array}{c}
 \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \\
 \text{Jacobian} \\
 1 \times N
 \end{array}
 =
 \begin{array}{c}
 1 \times N \\
 \boxed{1, 1, 1, \dots, 1} \\
 \text{Jacobian} \\
 1 \times N
 \end{array}
 \cdot
 \begin{array}{c}
 N \times N \\
 \boxed{\begin{array}{c} 2 \\ 2 \\ 2 \\ \dots \\ 2 \end{array}} \\
 \text{Jacobian} \\
 N \times N
 \end{array}
 =
 \begin{array}{c}
 1 \times N \\
 \boxed{2, 2, 2, \dots, 2}
 \end{array}$$

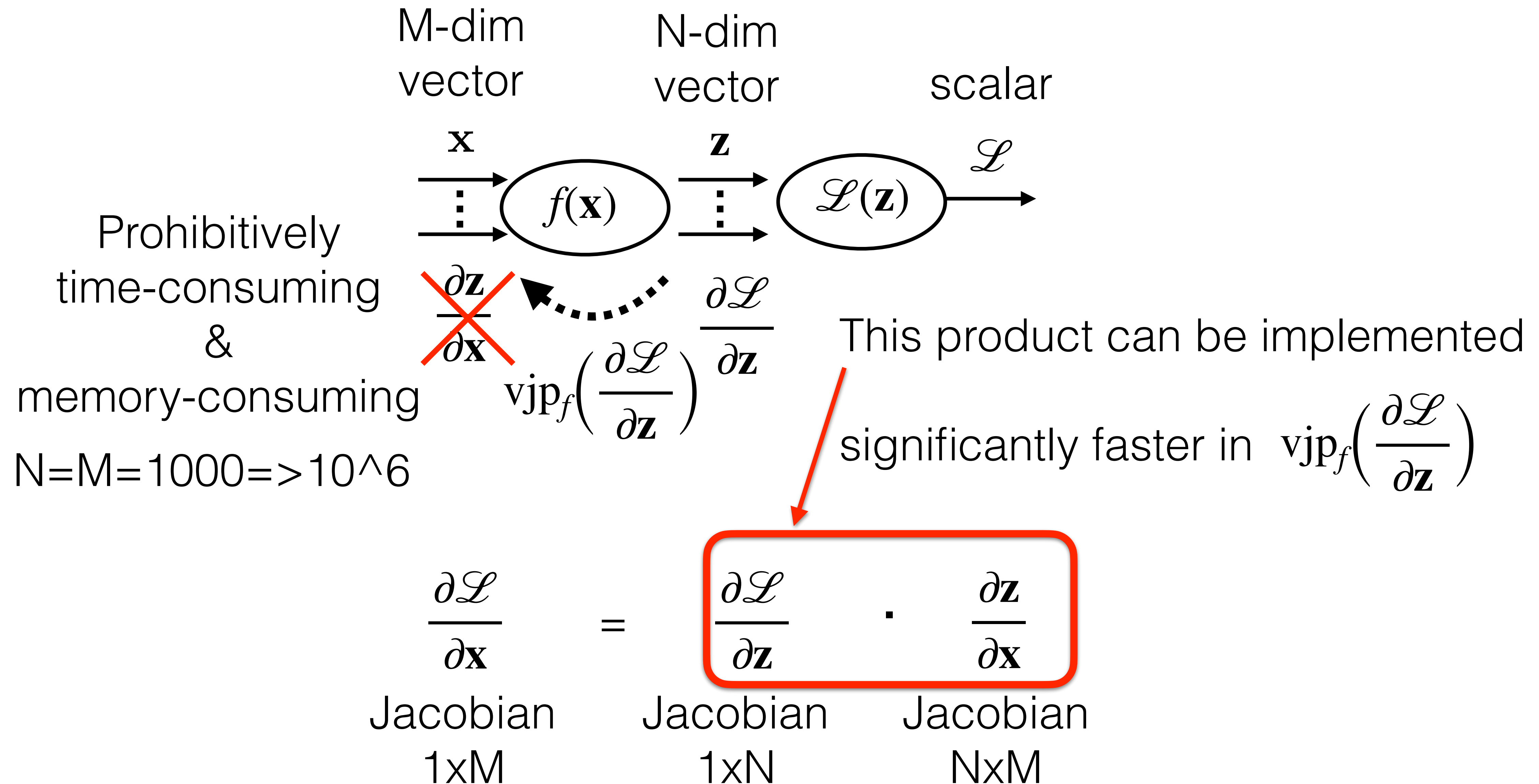
Backprop in Pytorch avoids Jacobians



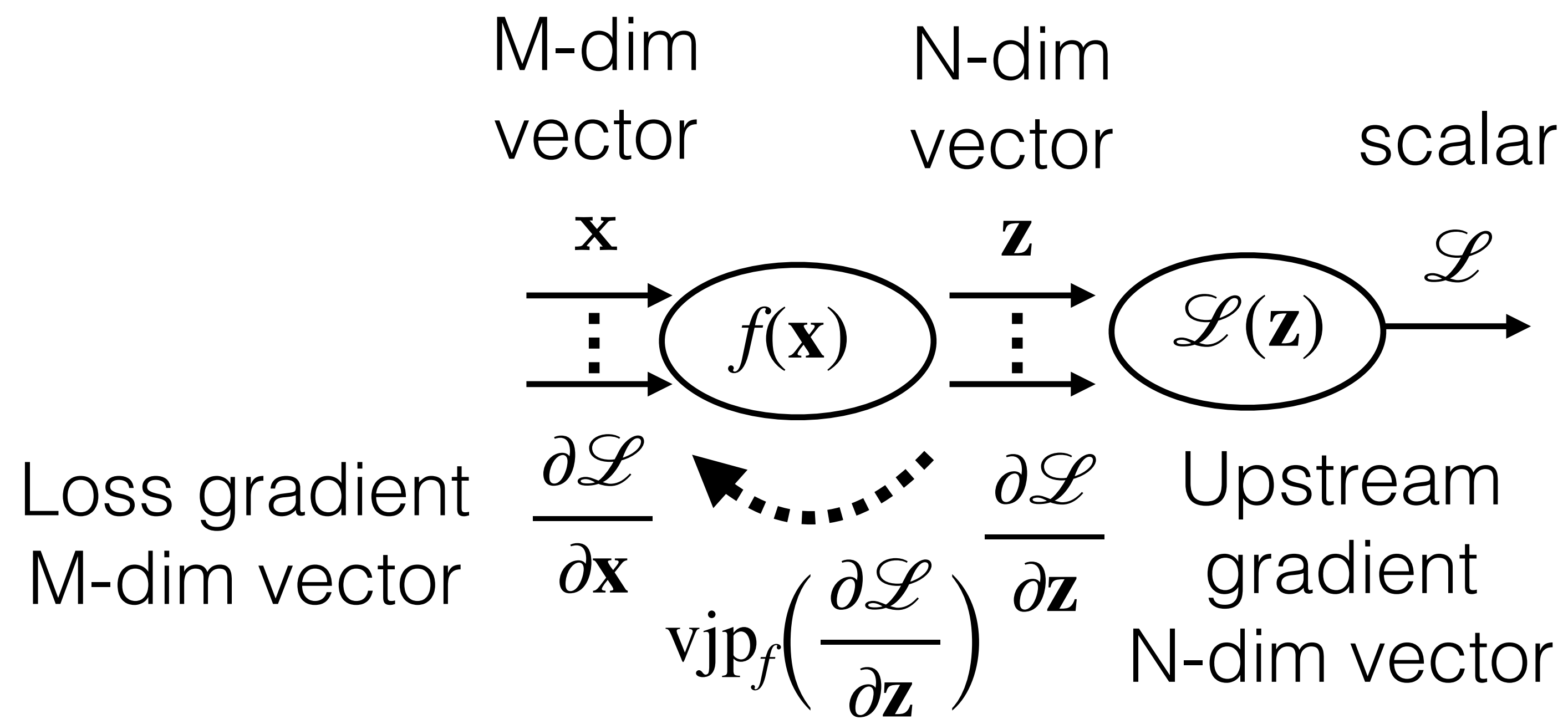
$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \text{vjp}_f\left(\begin{matrix} 1, 1, 1, \dots, 1 \end{matrix} \right) = \begin{matrix} 2, 2, 2, \dots, 2 \end{matrix}$$

Loss gradient N-vector vector-Jacobian product N-vector \rightarrow N-vector Loss gradient N-vector

Backprop in Pytorch avoids Jacobians



Backprop in Pytorch avoids Jacobians



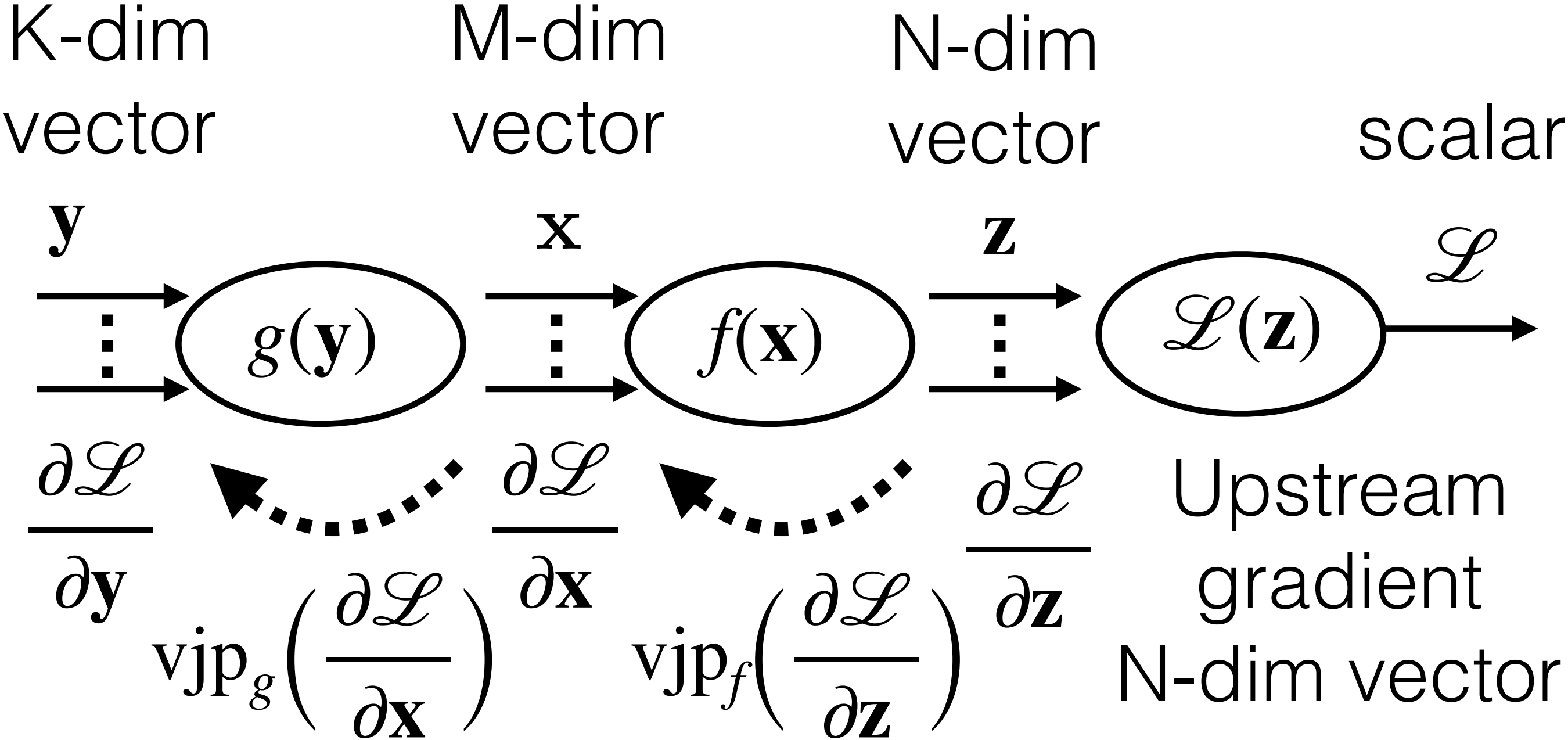
`grad_fn=<fBackward>`

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \text{vjp}_f\left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}}\right)$$

M-dim vector

vector-Jacobian product

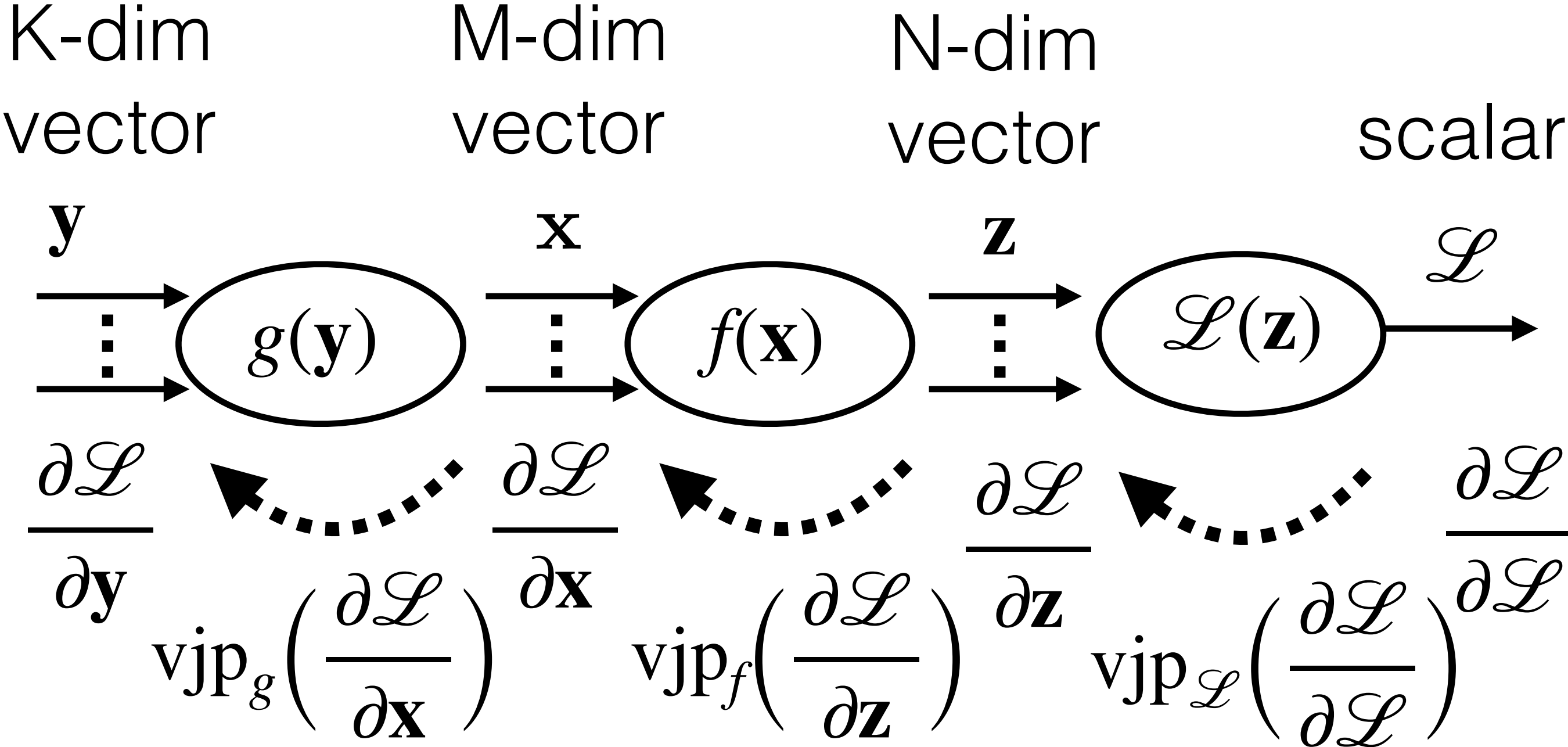
Backprop in Pytorch avoids Jacobians



$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \text{vjp}_g\left(\text{vjp}_f\left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}}\right)\right)$$

K-dim vector vector-Jacobian product

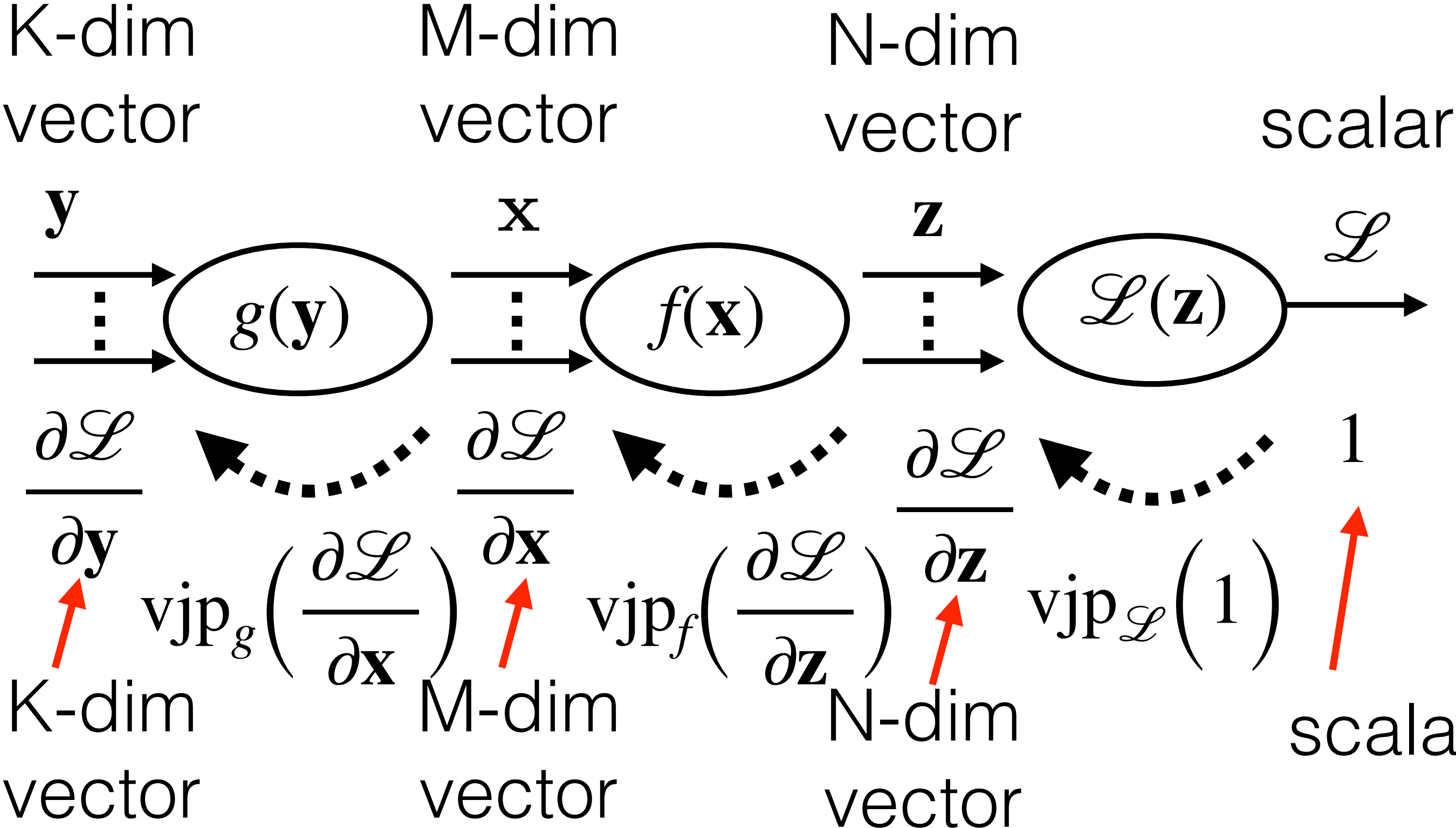
Backprop in Pytorch avoids Jacobians



$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \text{vjp}_g\left(\text{vjp}_f\left(\text{vjp}_{\mathcal{L}}\left(\frac{\partial \mathcal{L}}{\partial \mathcal{L}}\right)\right)\right)$$

K-dim vector vector-Jacobian product

Backprop in Pytorch avoids Jacobians



Loss gradient
K-dim vector

$\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$
K-dim vector

$$\text{vjp}_g\left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}}\right)$$

$\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$
M-dim vector

$$\text{vjp}_f\left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}}\right)$$

$\frac{\partial \mathcal{L}}{\partial \mathbf{z}}$
N-dim vector

$$\text{vjp}_{\mathcal{L}}(1)$$

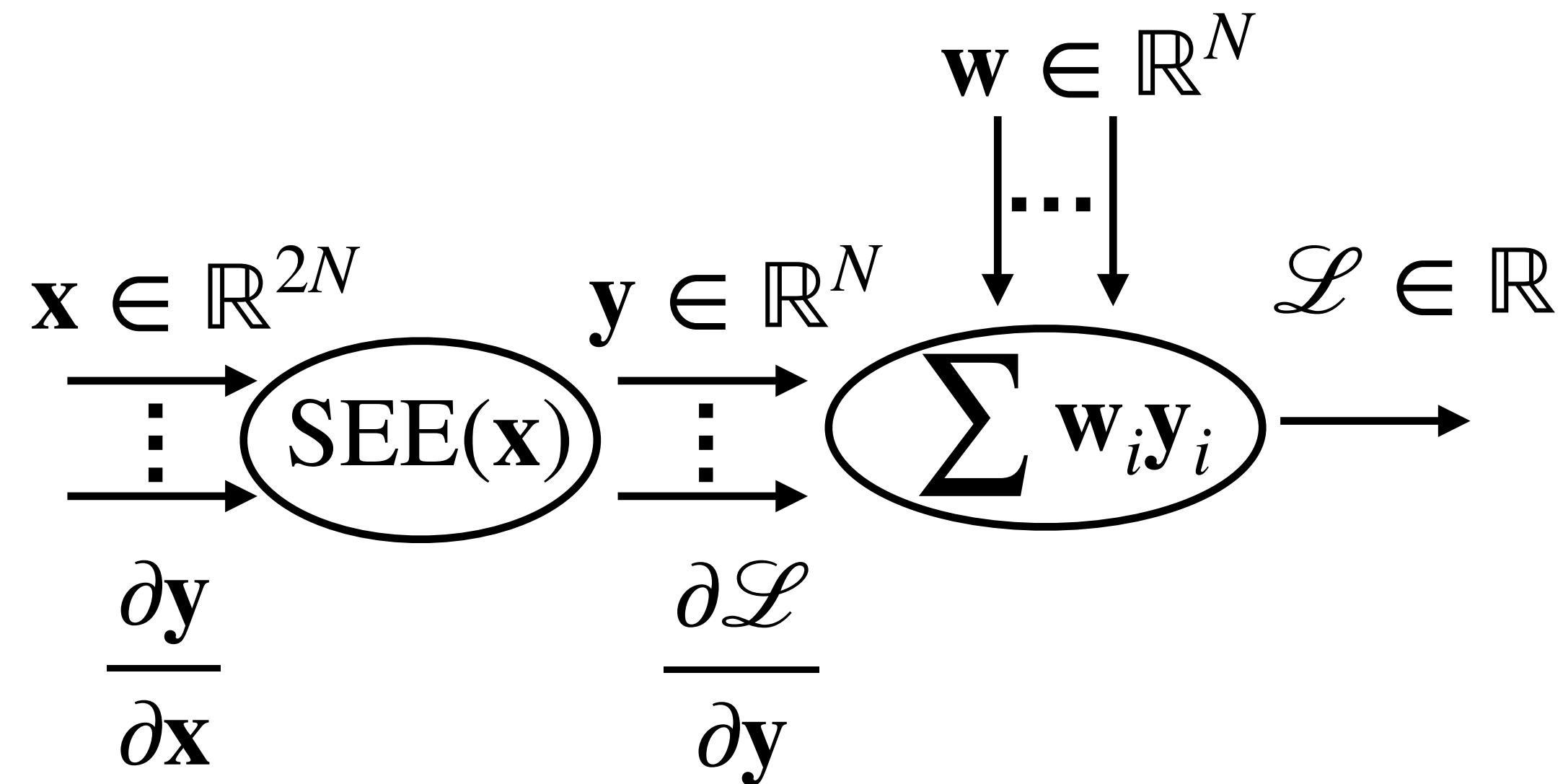
1
scalar

$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \text{vjp}_g\left(\text{vjp}_f\left(\text{vjp}_{\mathcal{L}}(1)\right)\right)$$

vector-Jacobian product

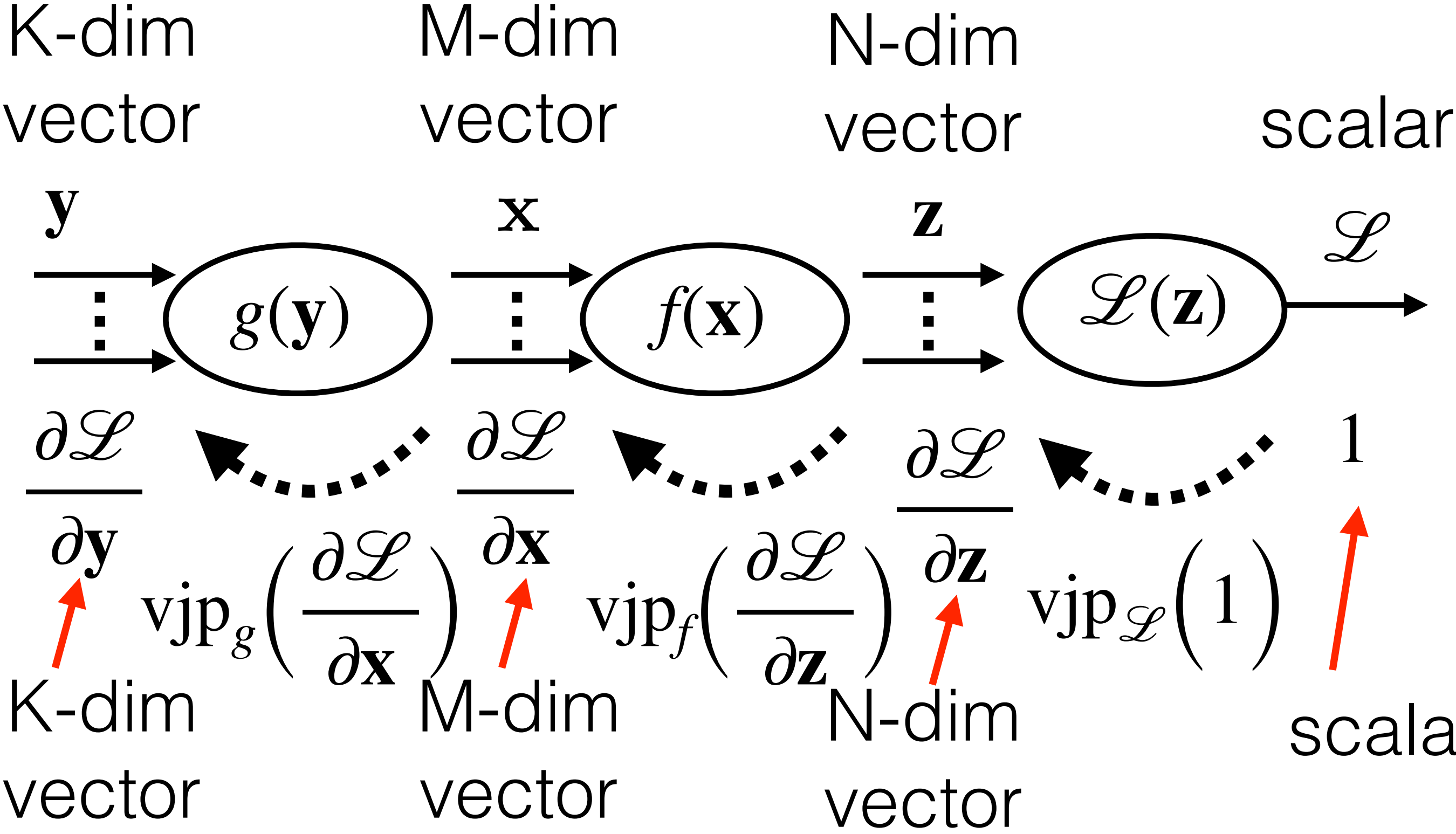
Example: Jacobian vs vector-jacobian-product function

$\mathbf{y} = \text{SEE}(\mathbf{x}) = \text{Select Even Element from input vector } \mathbf{x}$



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = ???$$

Backprop in Pytorch avoids Jacobians

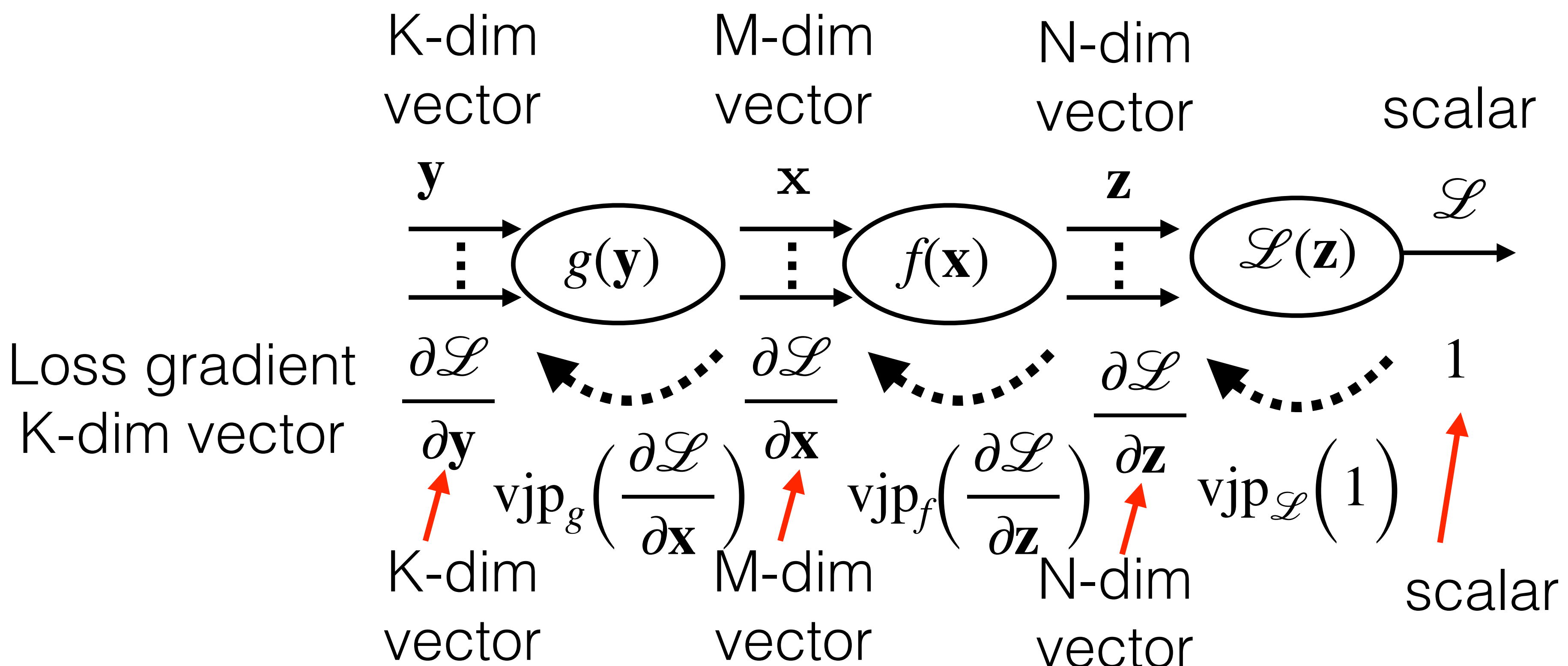


$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \text{vjp}_g\left(\text{vjp}_f\left(\text{vjp}_{\mathcal{L}}(1)\right)\right)$$

K-dim vector vector-Jacobian product

Backprop in Pytorch avoids Jacobians

Vector-Jacobian products allows to preserve data dimensionality (avoid vectorization)
 => Gradient in Pytorch is the tensor of the same dimensionality as the input.

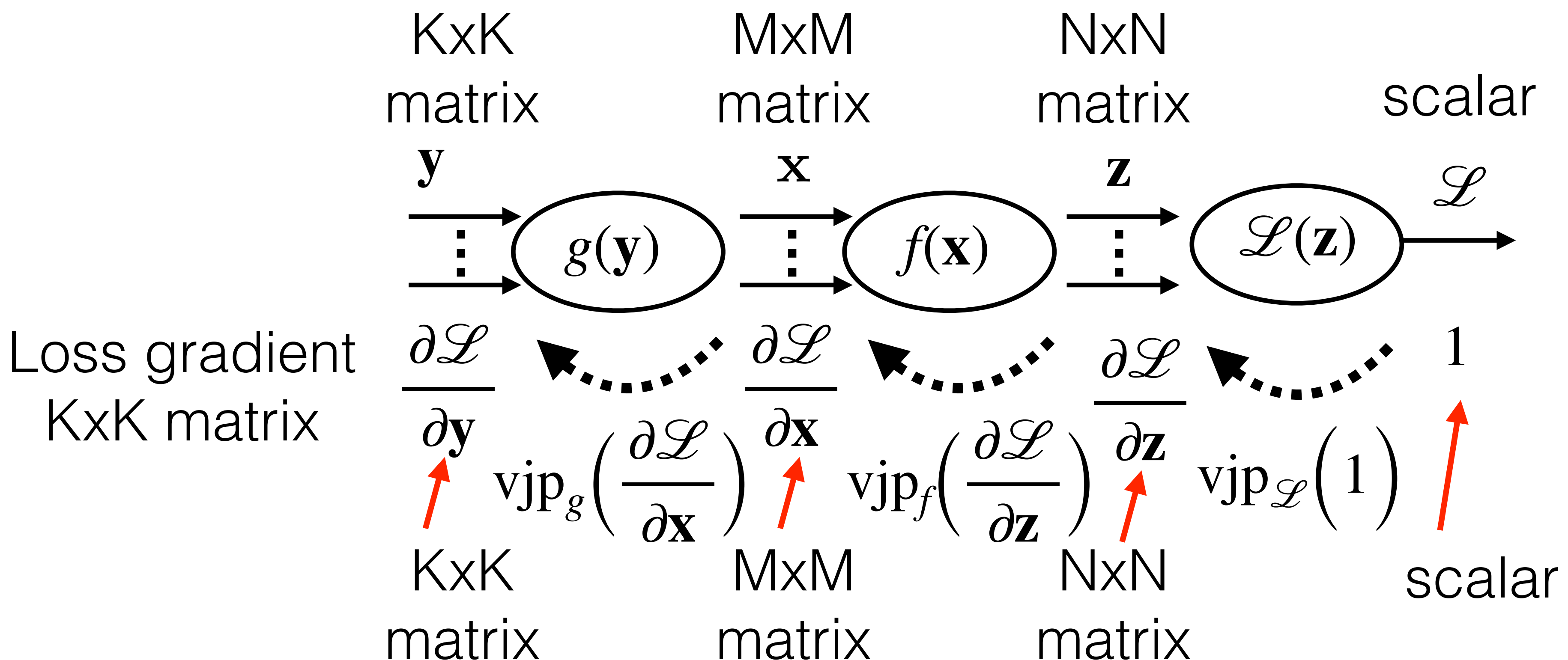


Backprop in Pytorch avoids Jacobians

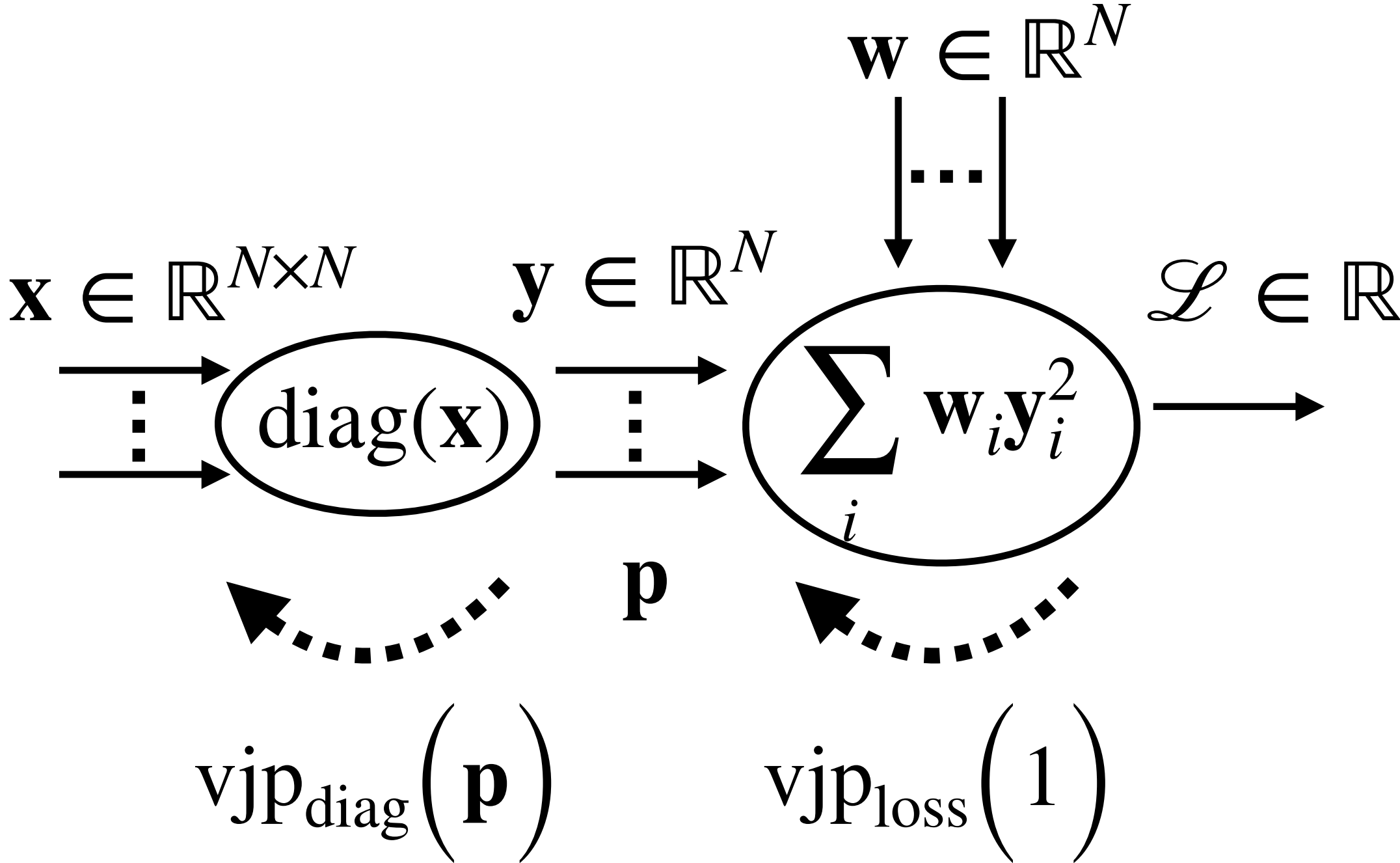
Vector-Jacobian products allows to preserve data dimensionality (avoid vectorization)

=> Gradient in Pytorch is the tensor of the same dimensionality as the input.

```
u = torch.tensor([[1,2], [3, 4]], dtype=torch.float32, requires_grad=True)
v = torch.sigmoid(u)
v.sum().backward()
print(u.grad)
tensor([[0.1966, 0.1050],
        [0.0452, 0.0177]])
```



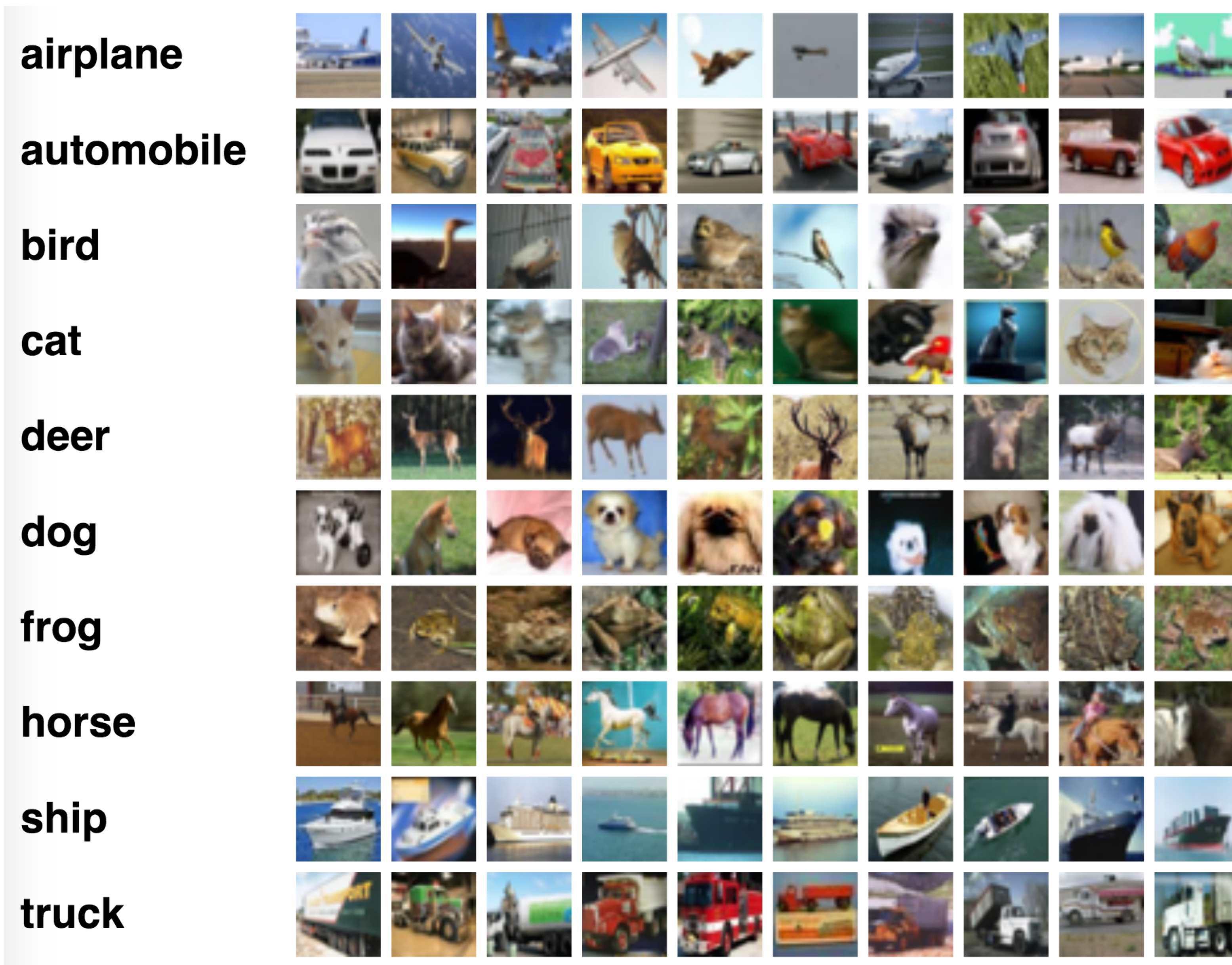
Example: Jacobian vs vector-jacobian-product function



Neural nets summary

- **Neural net** is a function created as concatenation of simpler functions (e.g. neurons or layers of neurons)
- **Fully connected neural net** is neural network created from neurons, where **all** outputs from previous layer are connected to **all** inputs of the following layer
- **Backpropagation** is gradient optimization of a neural net concatenated with a loss-layer on a training set.
- It is implemented as backward concatenation of vector-jacobian functions.
- **Deep learning frameworks** (Pytorch, Tensorflow) has many predefined layers.
- **Spoiler alert:** Fully connected NN does not work on structural data (images, sound) well.

3	4	2	1	9	5	6	2	1	8
8	9	1	2	5	0	0	6	6	4
6	7	0	1	6	3	6	3	7	0
3	7	7	9	4	6	6	1	8	2
2	9	3	4	3	9	8	7	2	5
1	5	9	8	3	6	5	7	2	3
9	3	1	9	1	5	8	0	8	4



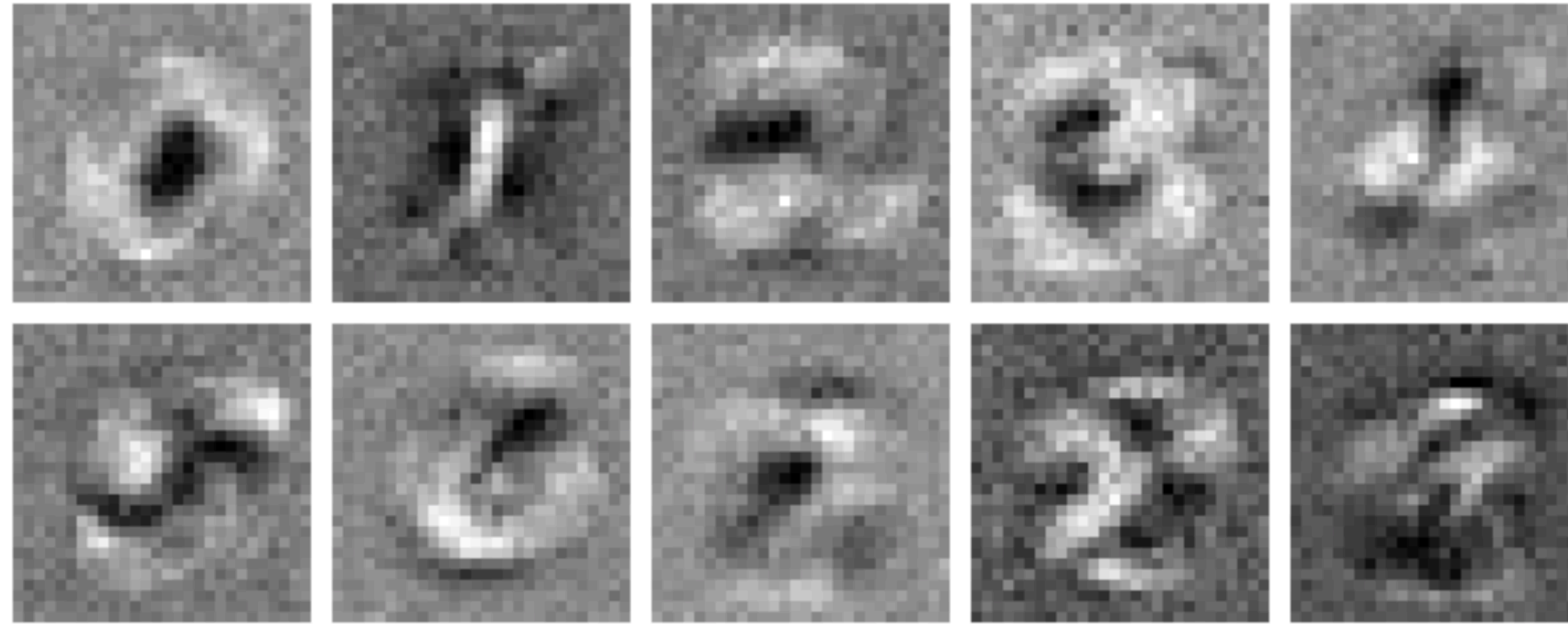
CIFAR-10: classify 32x32 RGB images into 10 categories
<https://www.cs.toronto.edu/~kriz/cifar.html>

Dataset

Learned weights of linear classifier

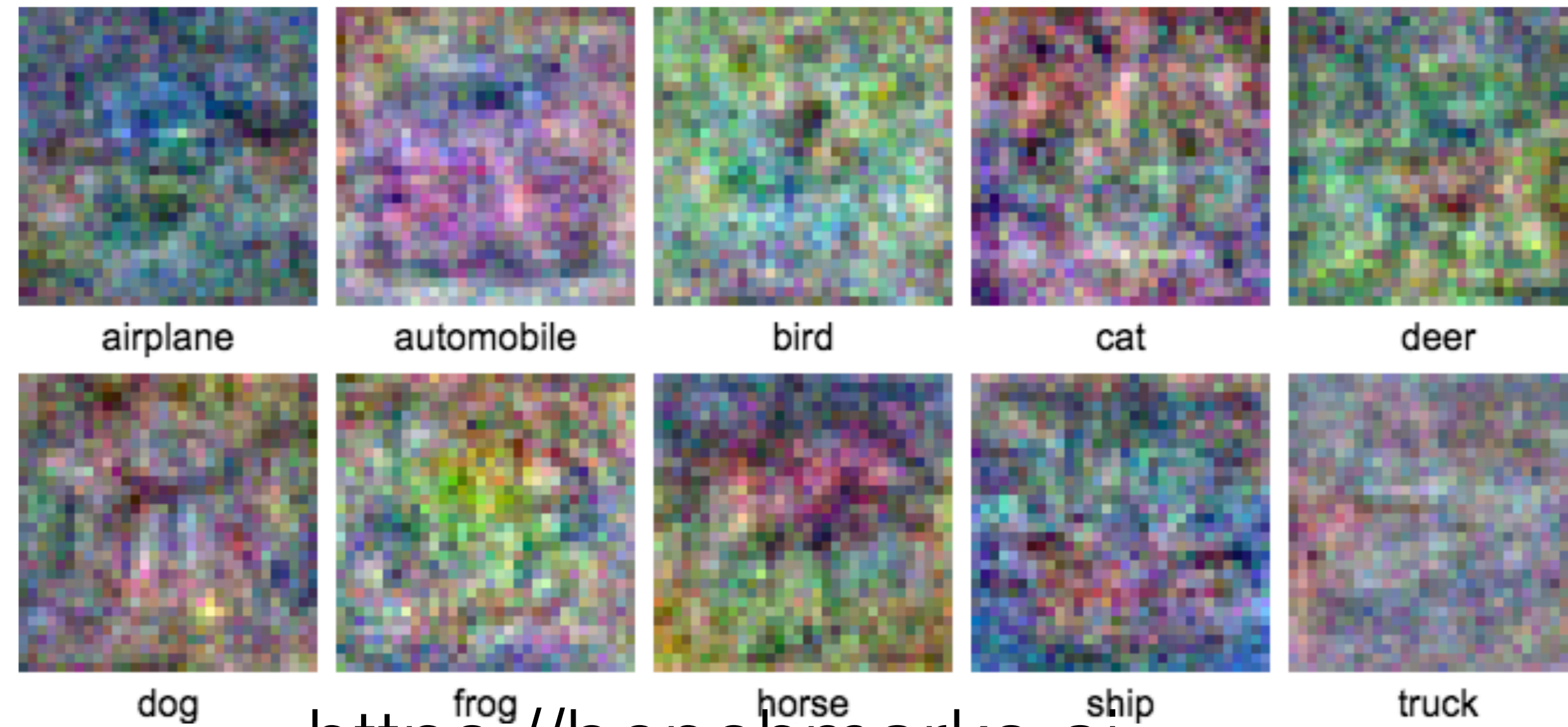
Error

MNIST



8%

CIFAR-10



63%

<https://benchmarks.ai>

Dataset

Error

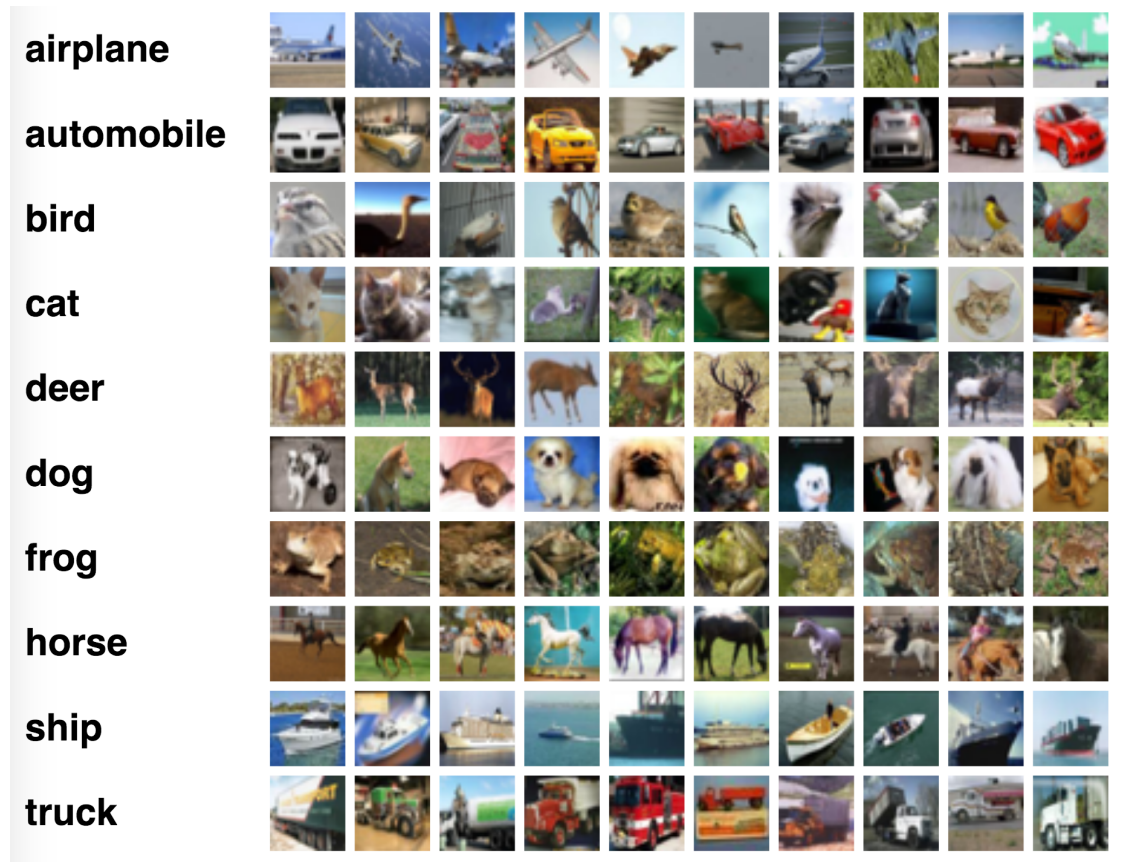
MNIST



Linear

8%

CIFAR-10



63%

<https://benchmarks.ai>

Dataset

Error

Linear

FCNN

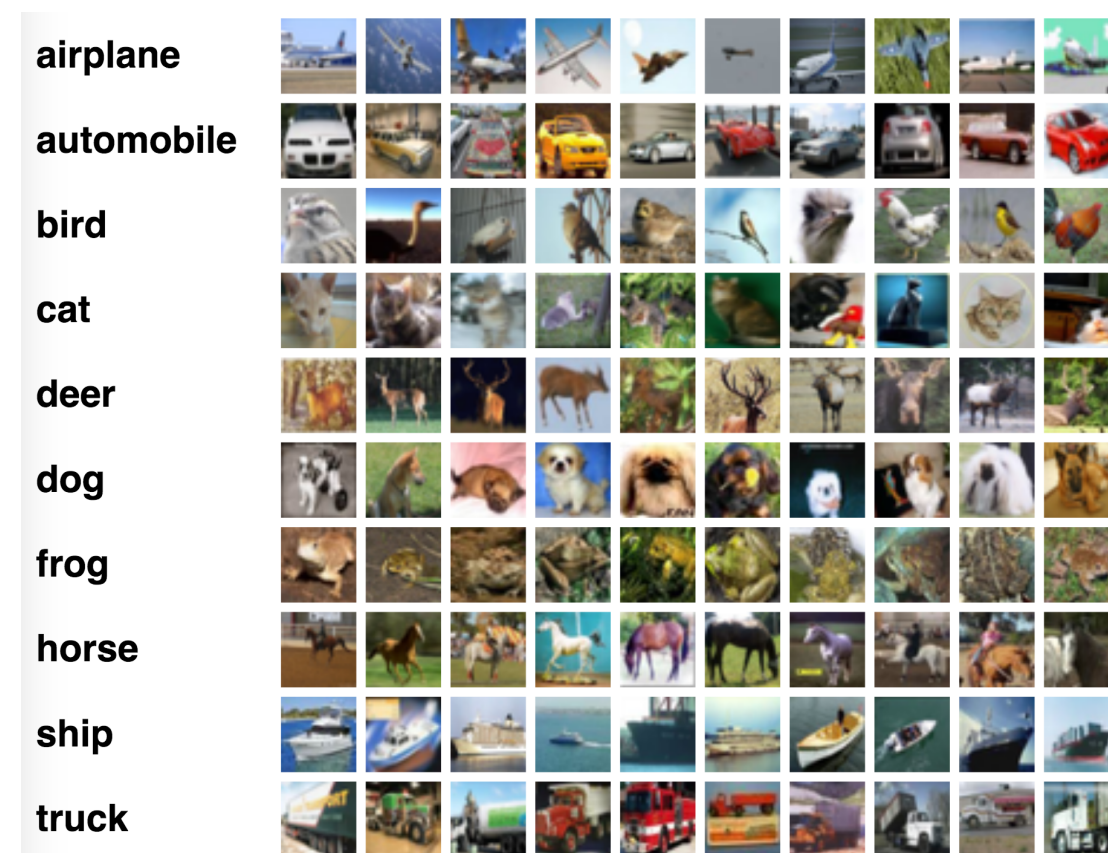
MNIST



8%

2%

CIFAR-10



63%

55%

<https://benchmarks.ai>

Dataset

Error

Linear

FCNN

ConvNet

MNIST



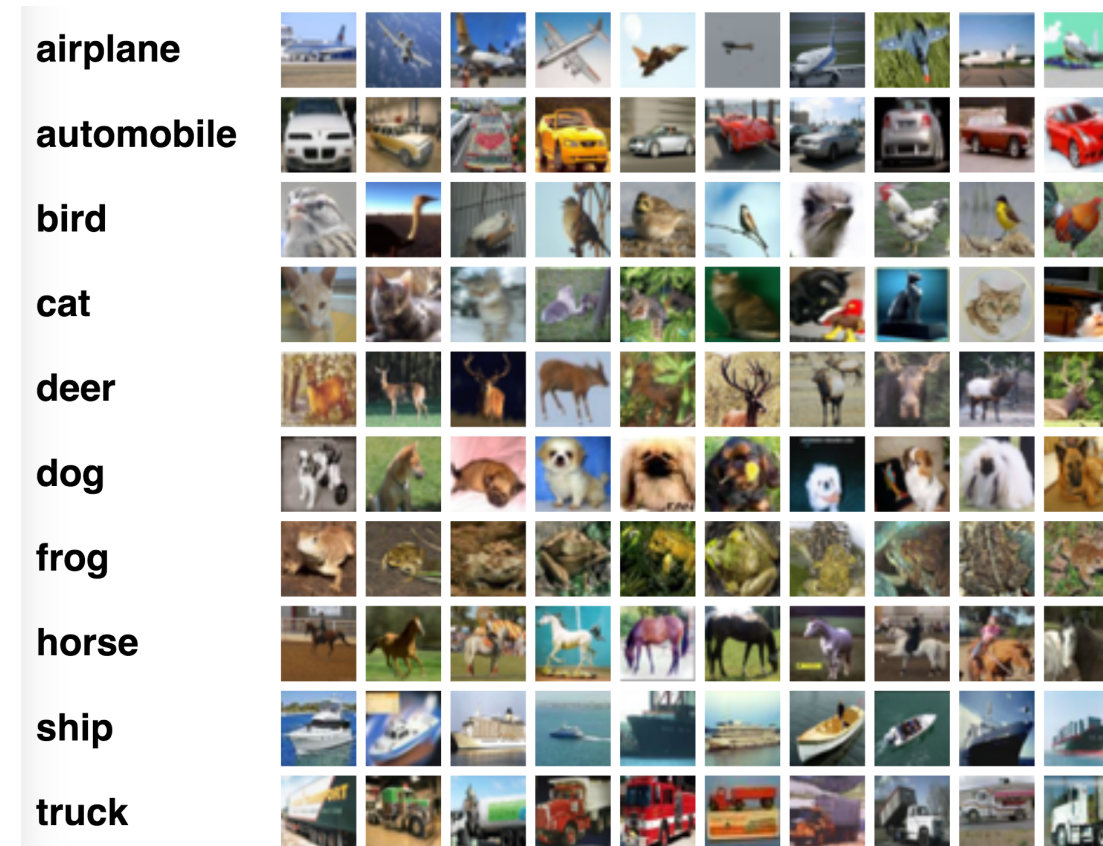
8%

2%

0.2%

[CVPR 2013]

CIFAR-10



63%

55%

1%

[EfficientNet,
2018]

<https://benchmarks.ai>

Dataset

MNIST

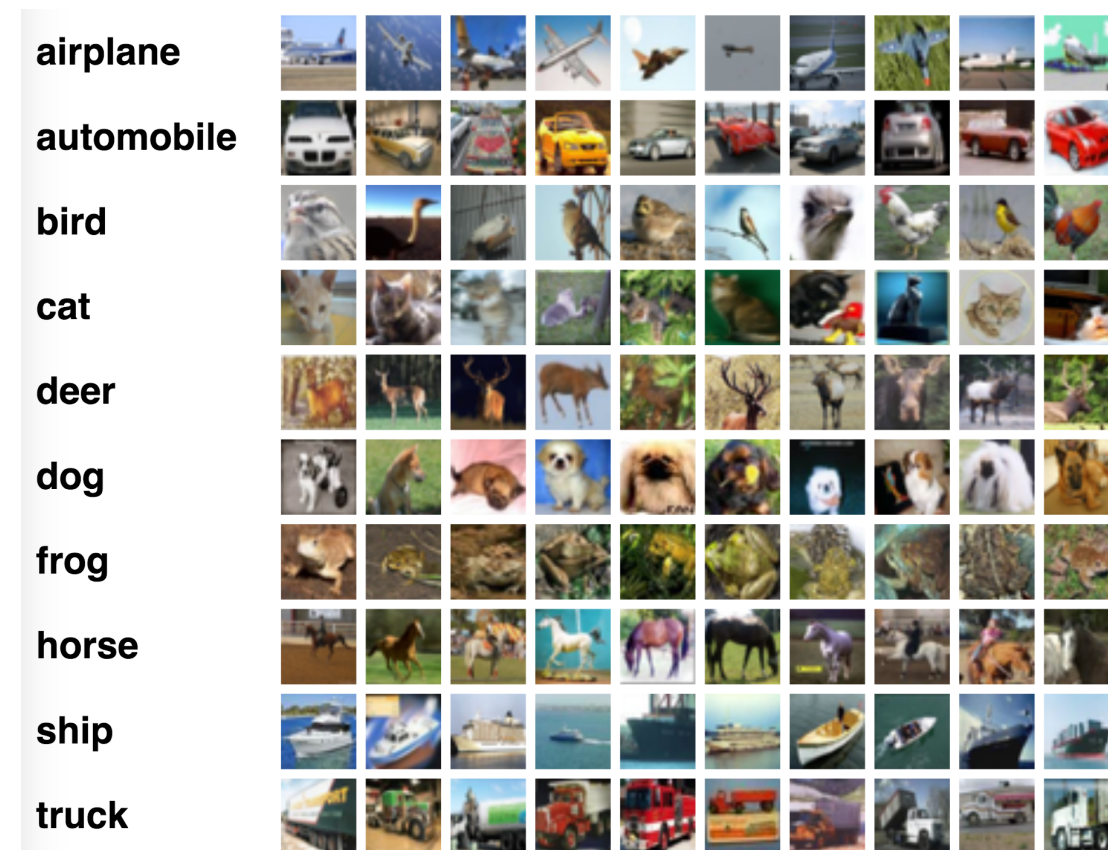


underfit

overfit

cortex inspired structure

CIFAR-10



<https://benchmarks.ai>

Error

Linear

FCNN

ConvNet

8%

2%

0.2%

[CVPR 2013]

63%

55%

1%

[EfficientNet,
2018]

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left(\sum_i -\log(p(y_i | \mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- Class of function represented by a FCNN is too general (it allows for many leprechauns).
- Naive regulariser helps a bit, but dimensionality/wildness is huge => curse-of-dimensionality, overfitting,...
- What is number of weights between two 1000-neuron layers?
- **Next lecture:** study animal cortex to find a stronger prior on the class of suitable functions.
- **Spoiler alert 2:**
reduce very general class of functions "neuron layer" to very specific sub-class of functions "convolution layer"

Competencies required for the test T1

- Ability to draw a computational graph.
- Compute edge gradients/jacobians.
- Perform one step of backpropagation in a vectorized form