

Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

dlvu.github.io



BACKPROPAGATION

Preliminaries:

neural networks

the gradient of the loss

learning by gradient descent

Today:

backpropagation

THE PLAN

part 1: review

part 2: scalar backpropagation

part 3: tensor backpropagation

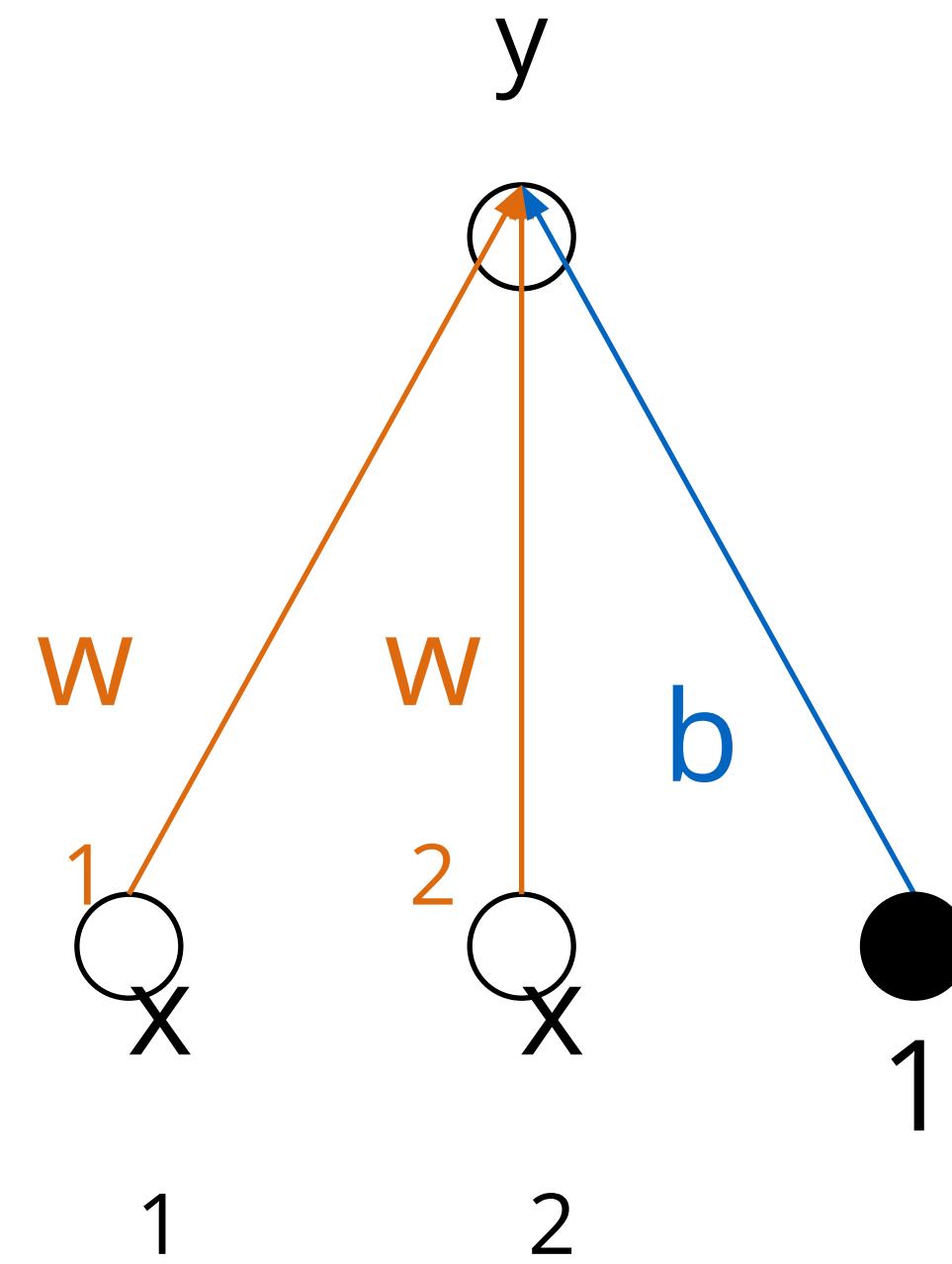
part 4: automatic differentiation

PART ONE: REVIEW

RECAP: NEURAL NETWORKS

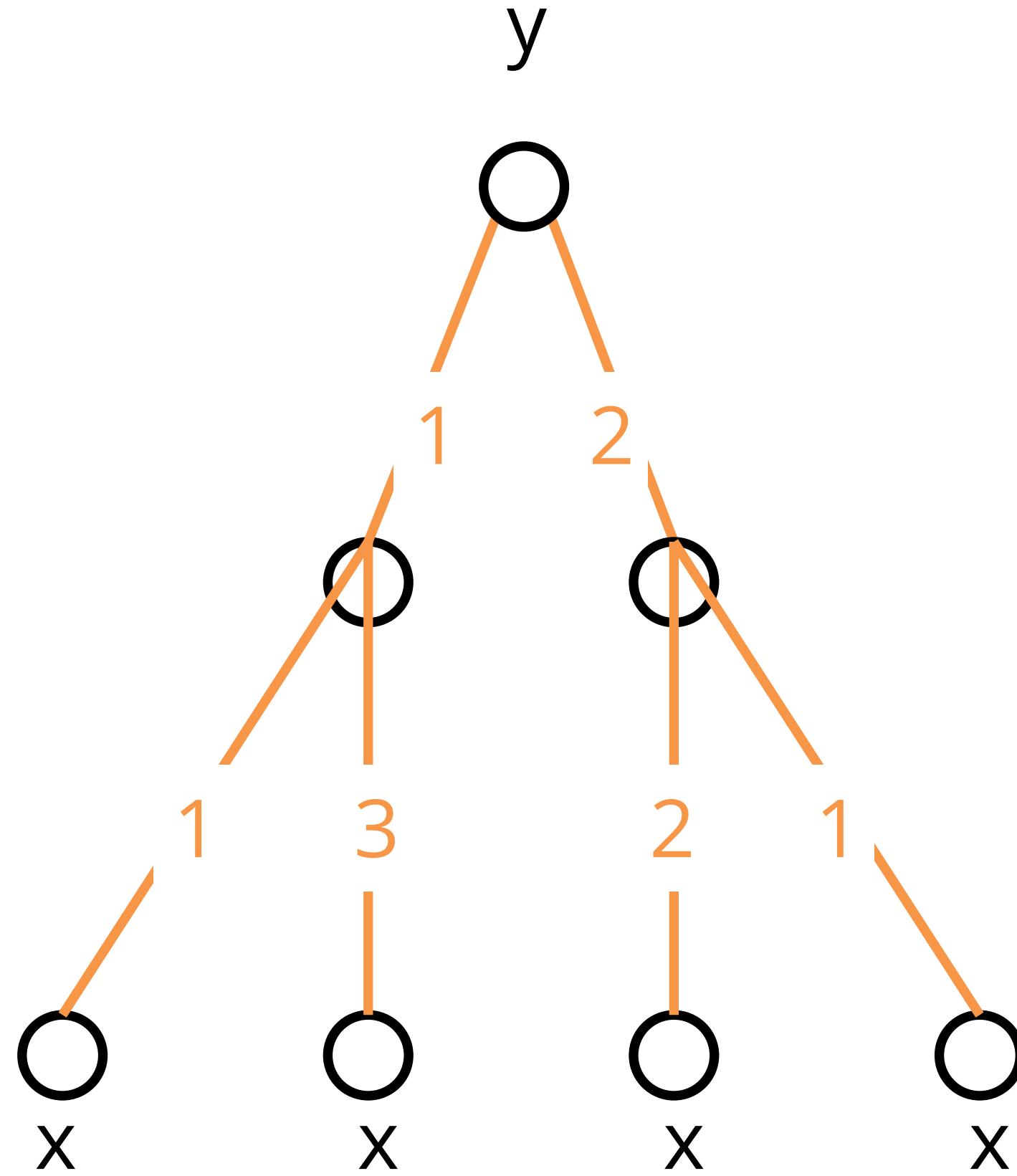


1957: THE PERCEPTRON



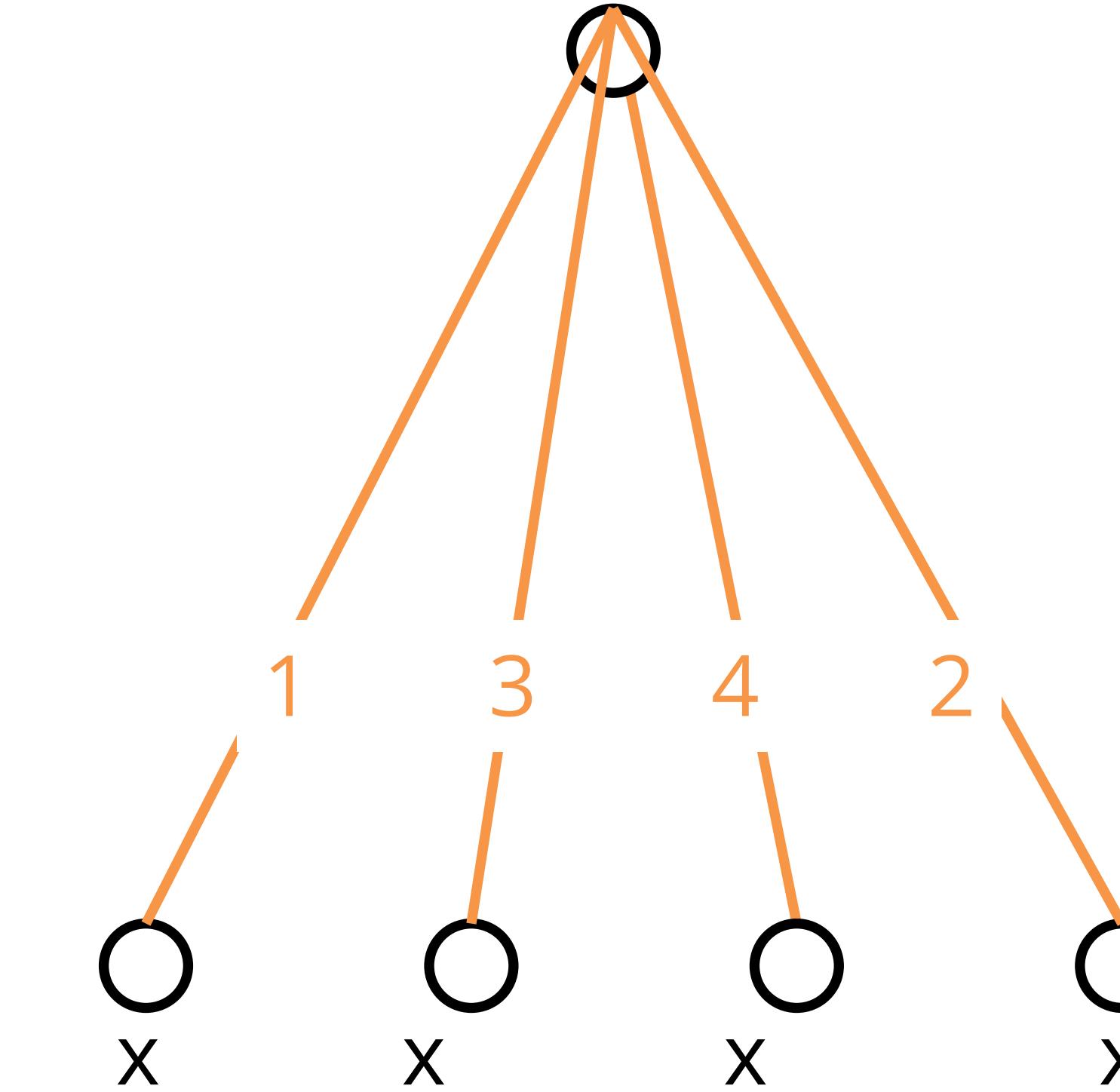
$$y = w_1 x_1 + w_2 x_2 + b$$

PROBLEM: COMPOSING NEURONS



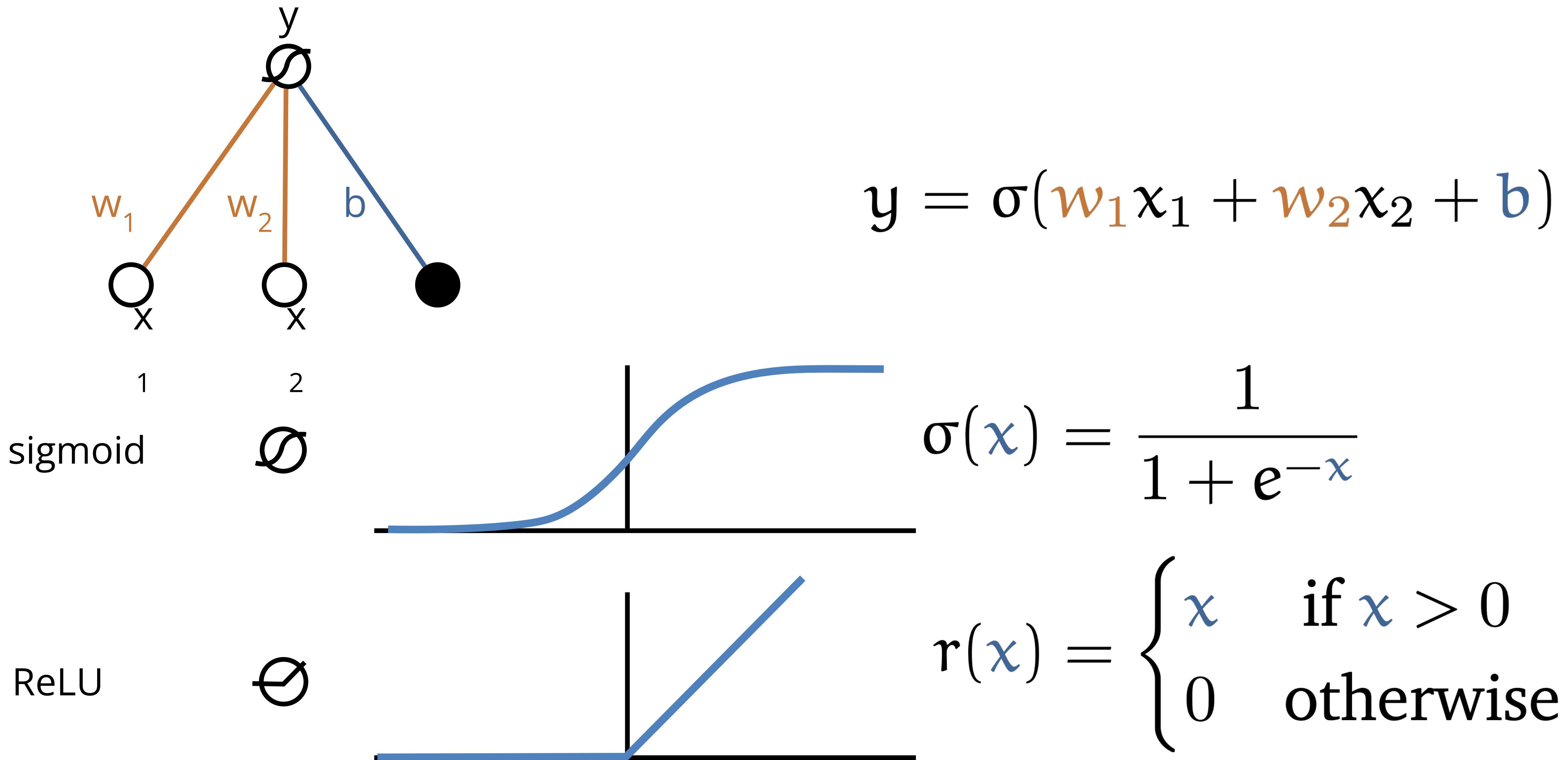
$$y = 1(1x_1 + 3x_2) + 2(2x_3 + 1x_4)$$

7



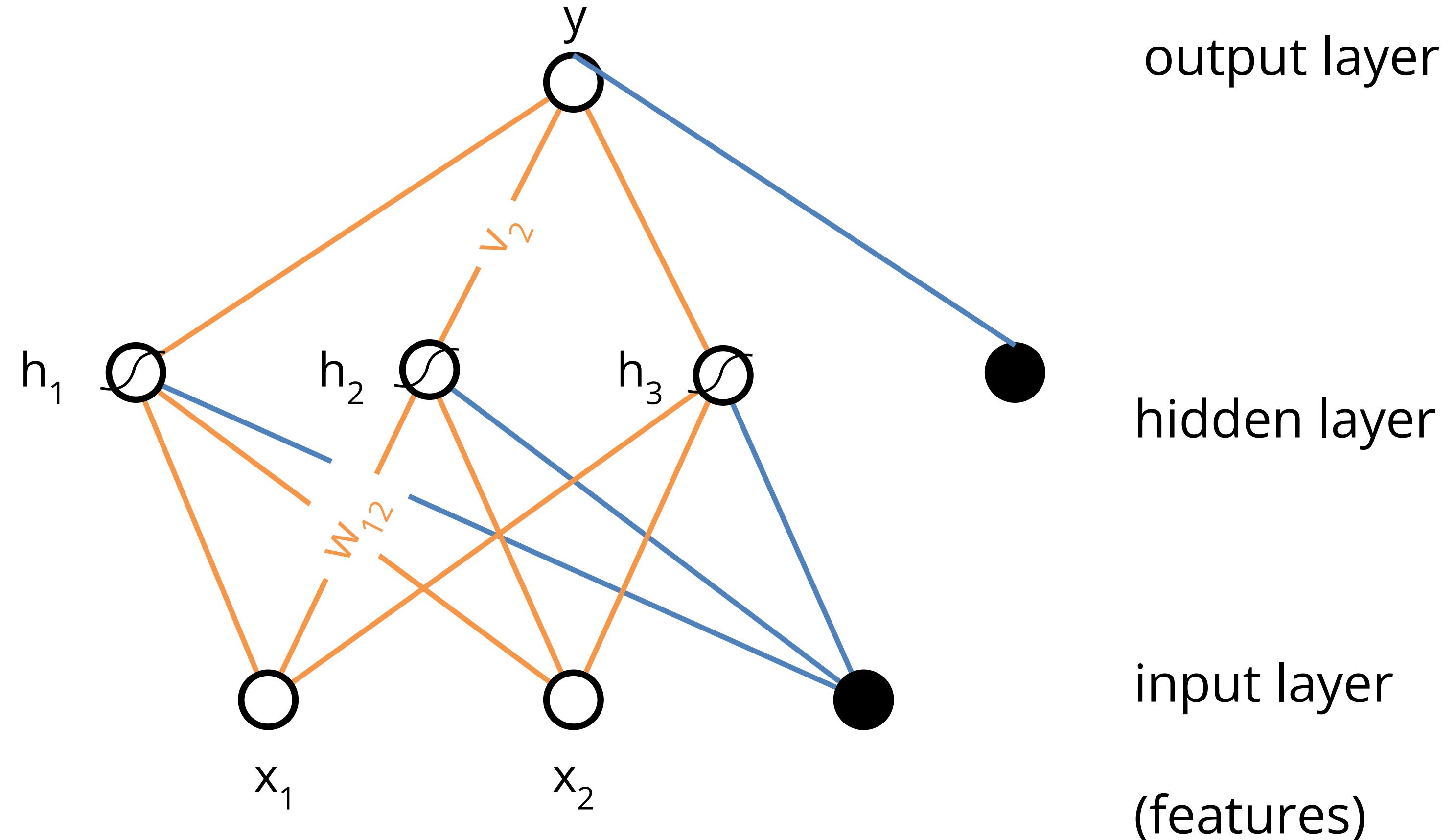
$$y = 1x_1 + 3x_2 + 4x_3 + 2x_4$$

NONLINEARITY



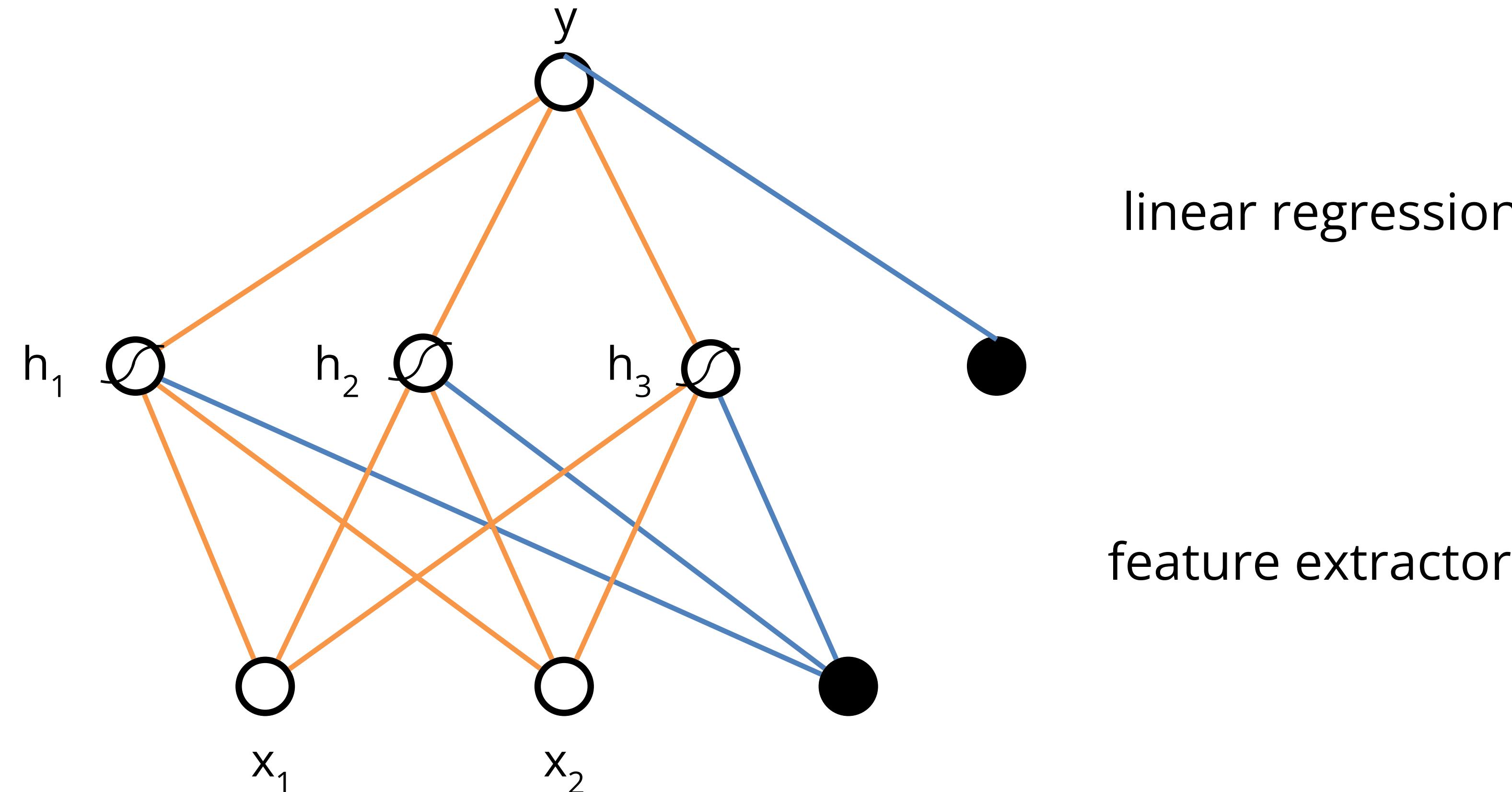
OTHER NONLINEARITIES

FEEDFORWARD NETWORK

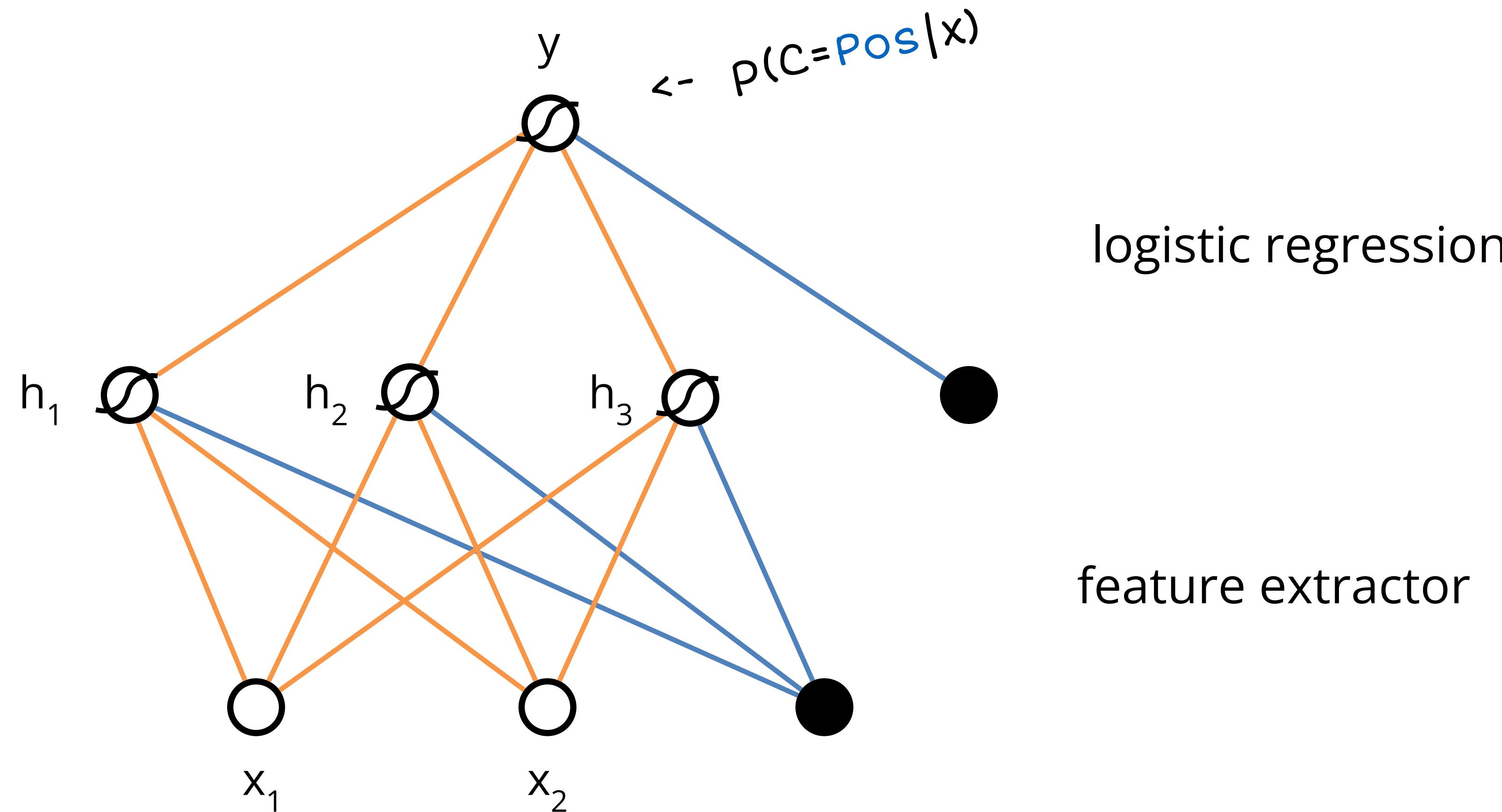


aka Multilayer Perceptron (MLP)

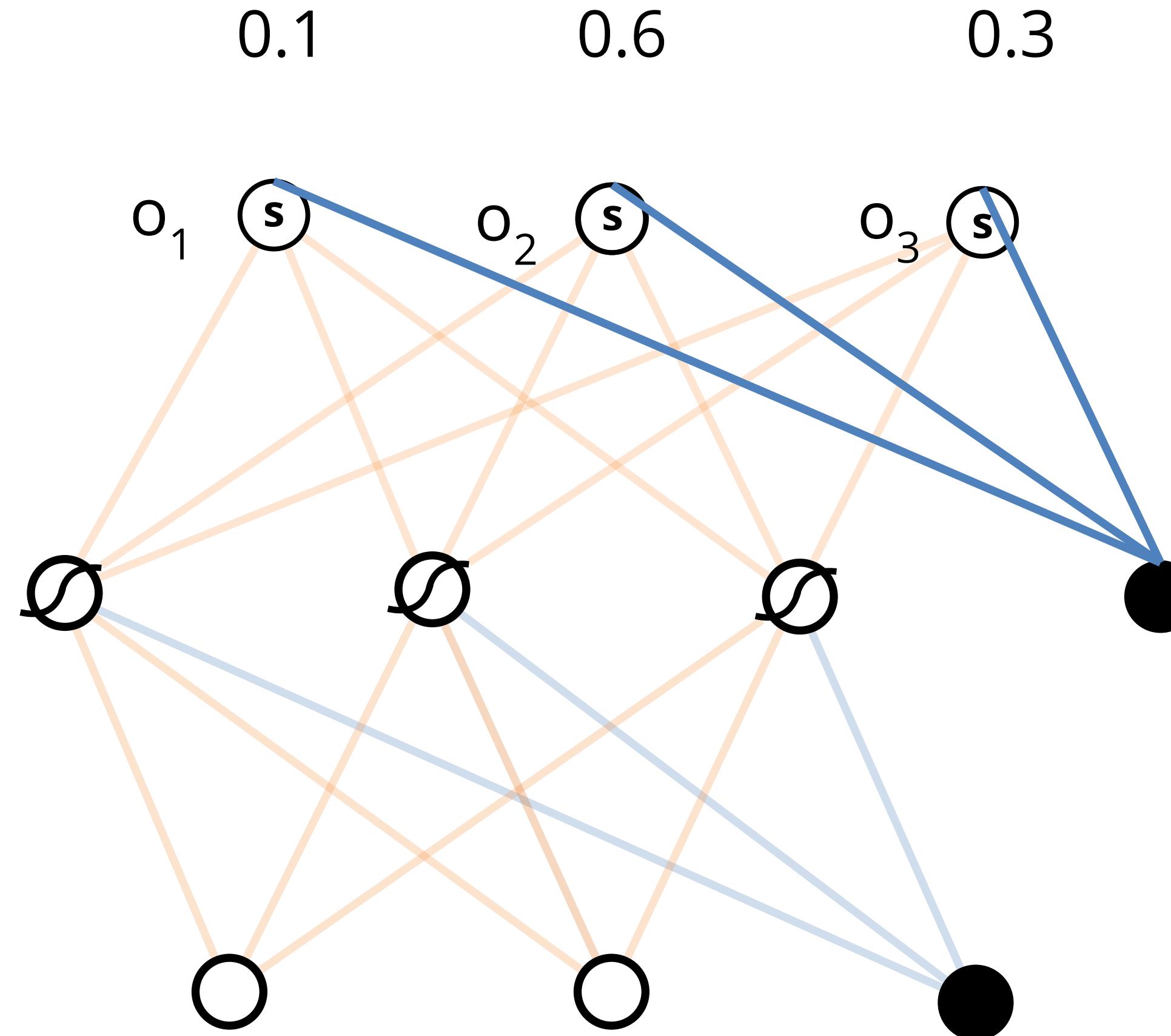
REGRESSION



BINARY CLASSIFICATION

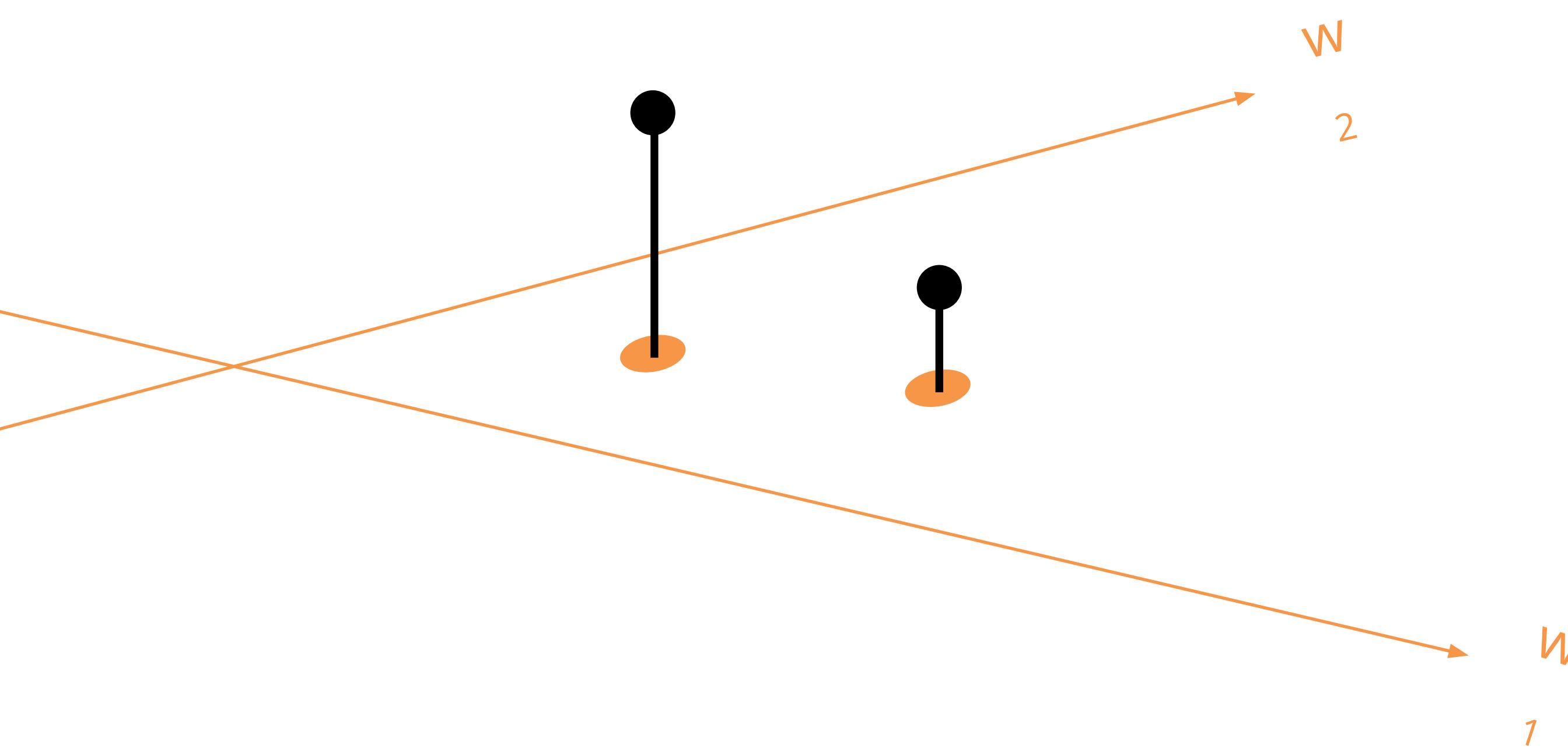


MULTI-CLASS CLASSIFICATION: THE SOFTMAX ACTIVATION



$$y_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

HOW DO WE FIND GOOD WEIGHTS?



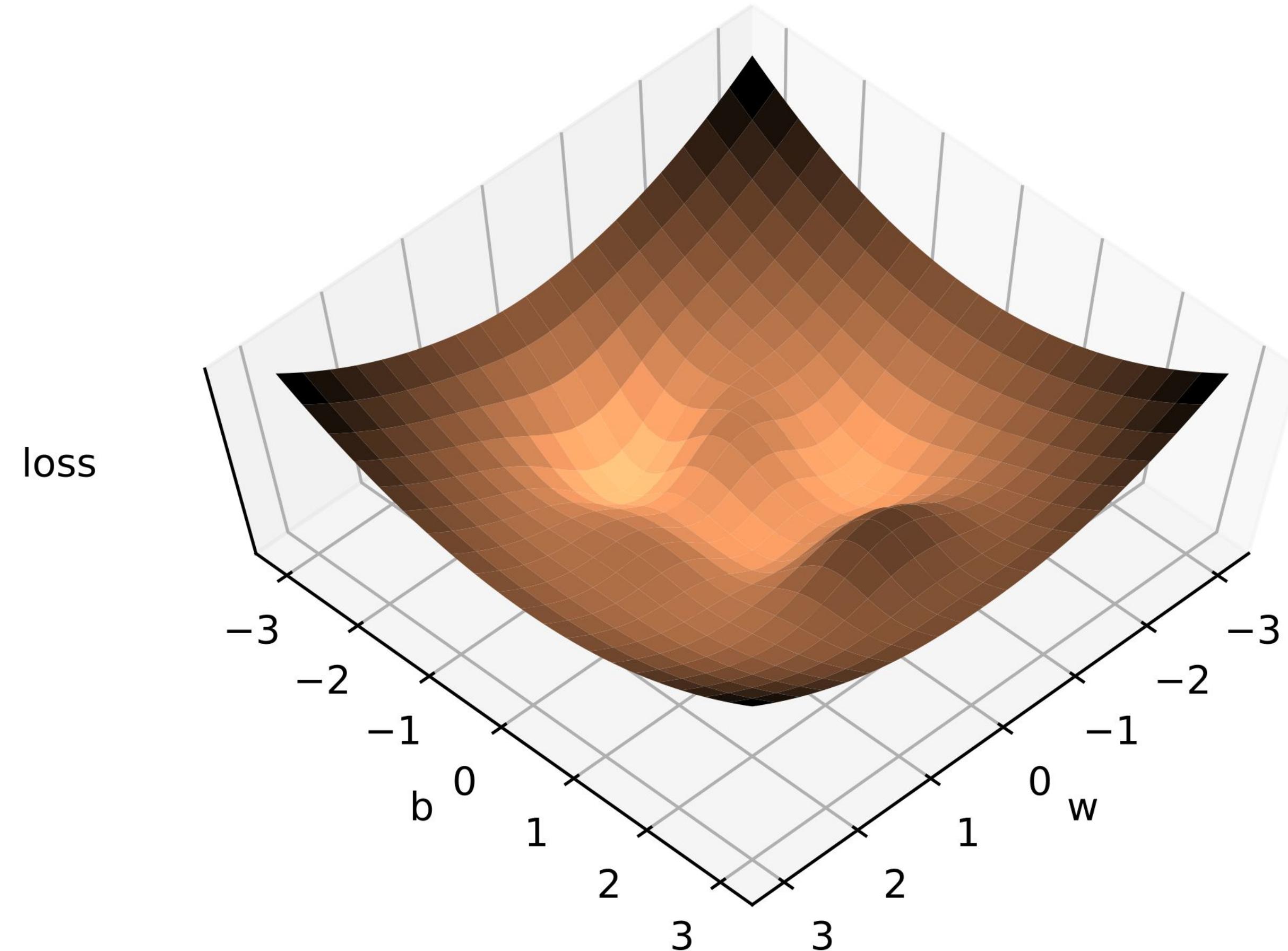
$$\text{model}_{\theta}(\mathbf{x}) = \mathbf{y}$$
$$\text{loss}_{\mathbf{x}, \mathbf{t}}(\theta) = \|\text{model}_{\theta}(\mathbf{x}) - \mathbf{t}\|$$

$$\arg \min_{\theta} \text{loss}_{\text{data}}(\theta)$$

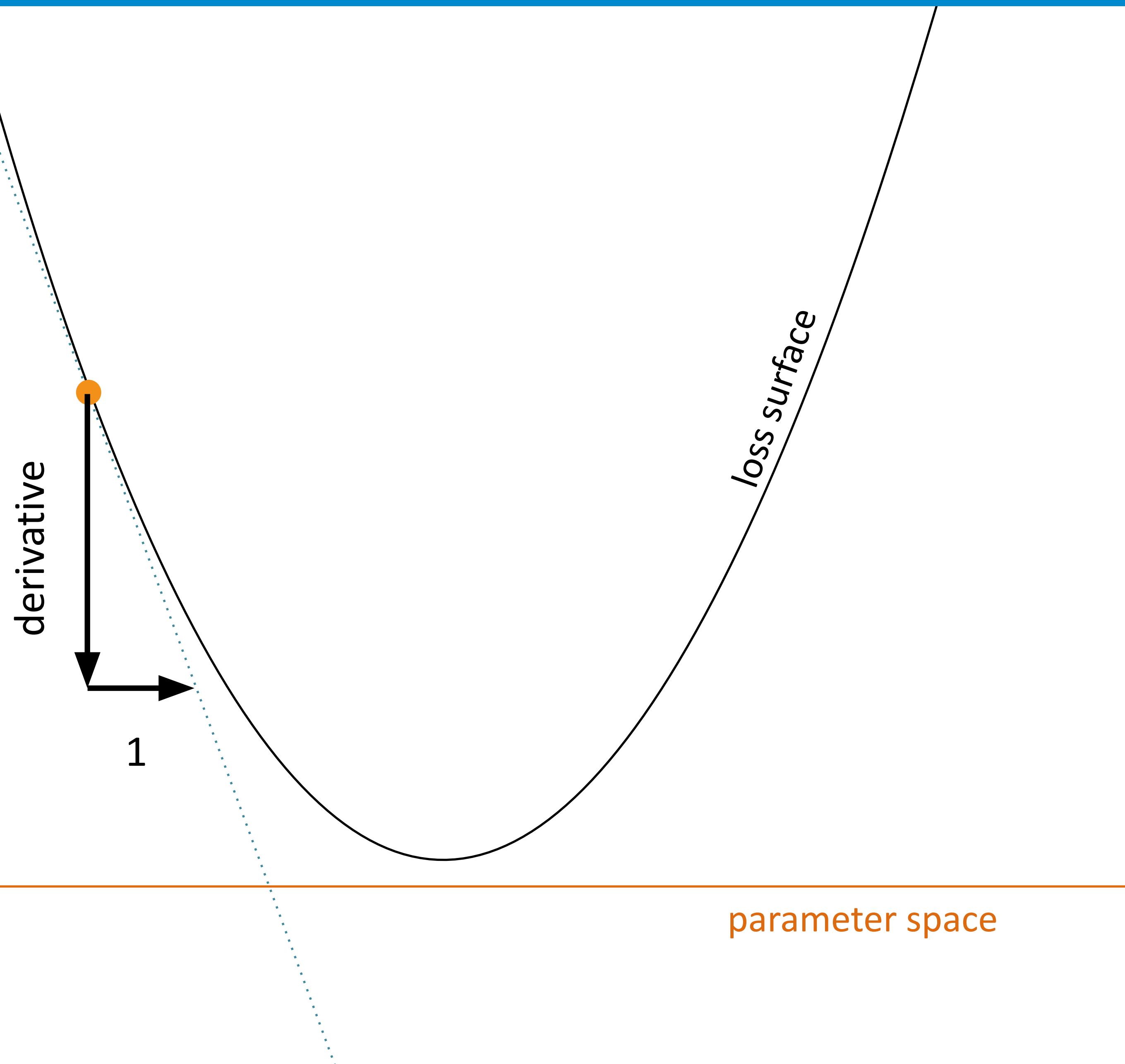
SOME COMMON LOSS FUNCTIONS

regression	$\ y - t\ $ with $y = \text{model}_\Theta(x)$
absolute errors	$\ y - t\ _1 = \sum_i \text{abs}(y_i - t_i)$
binary cross-entropy	$-\log p_\Theta(t)$ with $t \in \{0, 1\}$
cross-entropy	$-\log p_\Theta(t)$ with $t \in \{0, \dots, K\}$
hinge loss	$\max(0, 1 - ty)$ with $t \in \{-1, 1\}$

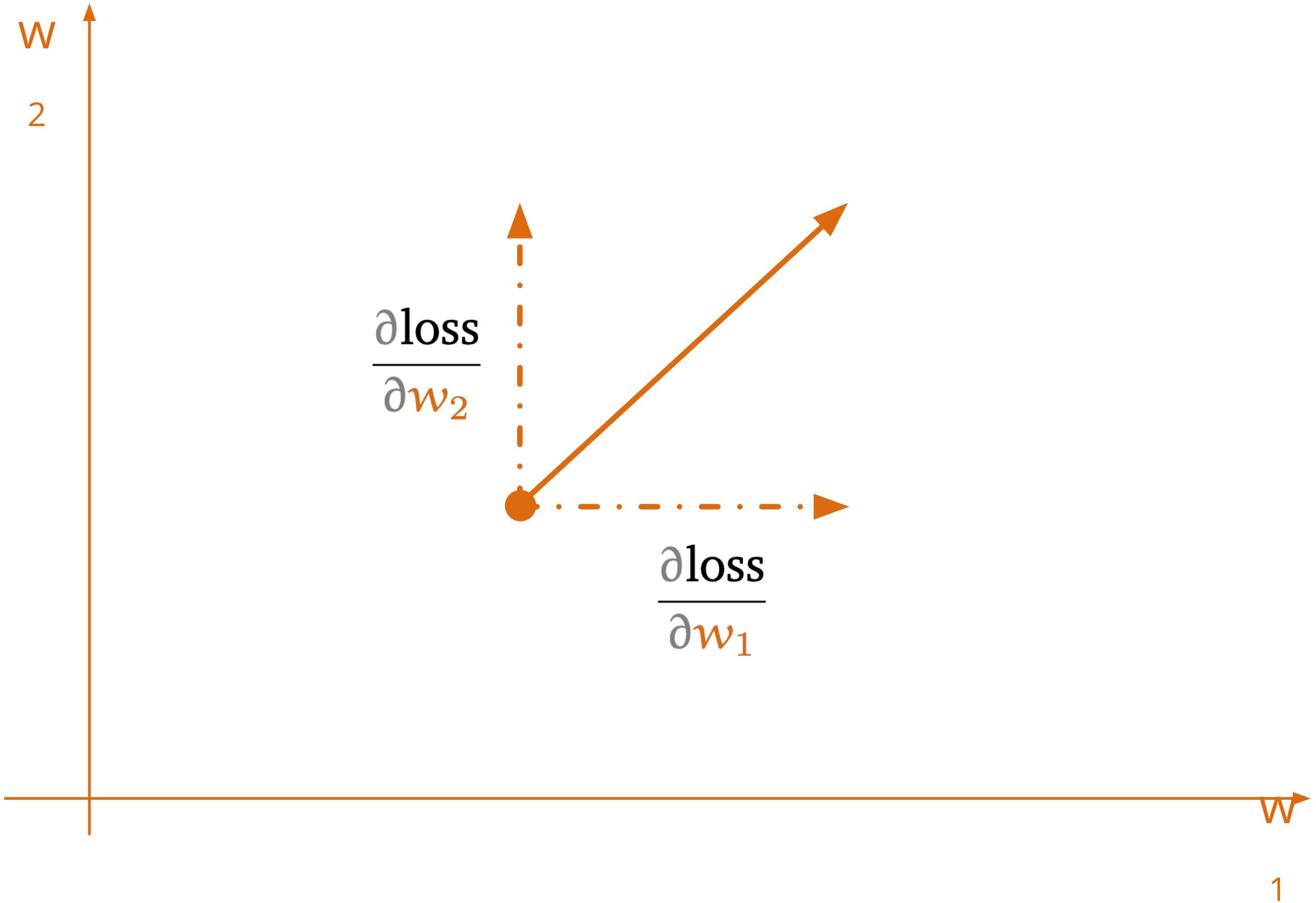
LOSS SURFACE



DERIVATIVE



GRADIENT



$$\nabla_{\theta} \text{loss}_{x,t}(\theta)$$

STOCHASTIC GRADIENT DESCENT

pick some initial weights θ (for the whole model)

loop:

for x, t in Data:
 $\theta \leftarrow \theta - \alpha \cdot \nabla_{\theta} \text{loss}_{x,t}(\theta)$



stochastic gradient descent: loss over one example per step (loop over data)

minibatch gradient descent: loss over a few examples per step (loop over data)

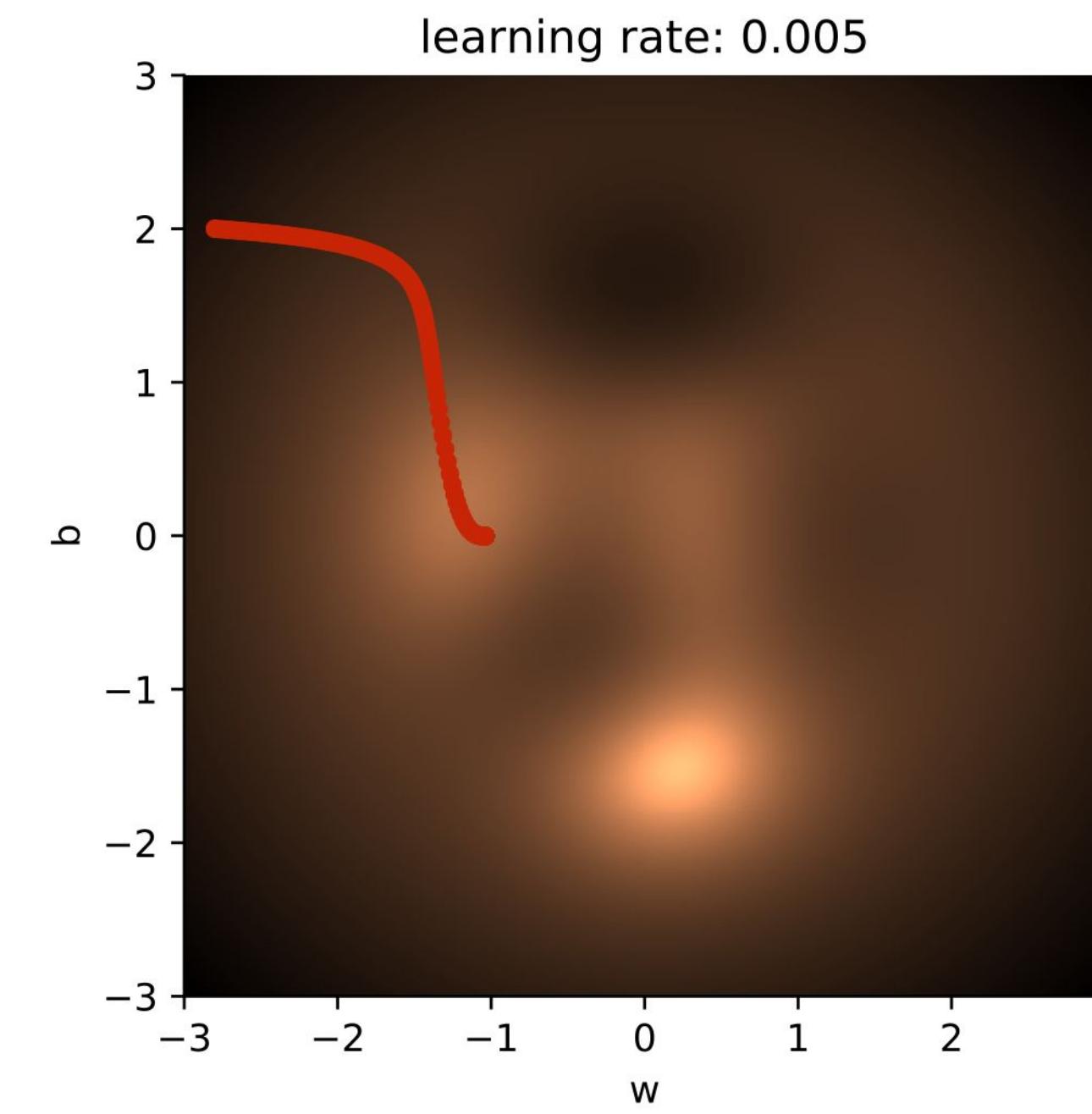
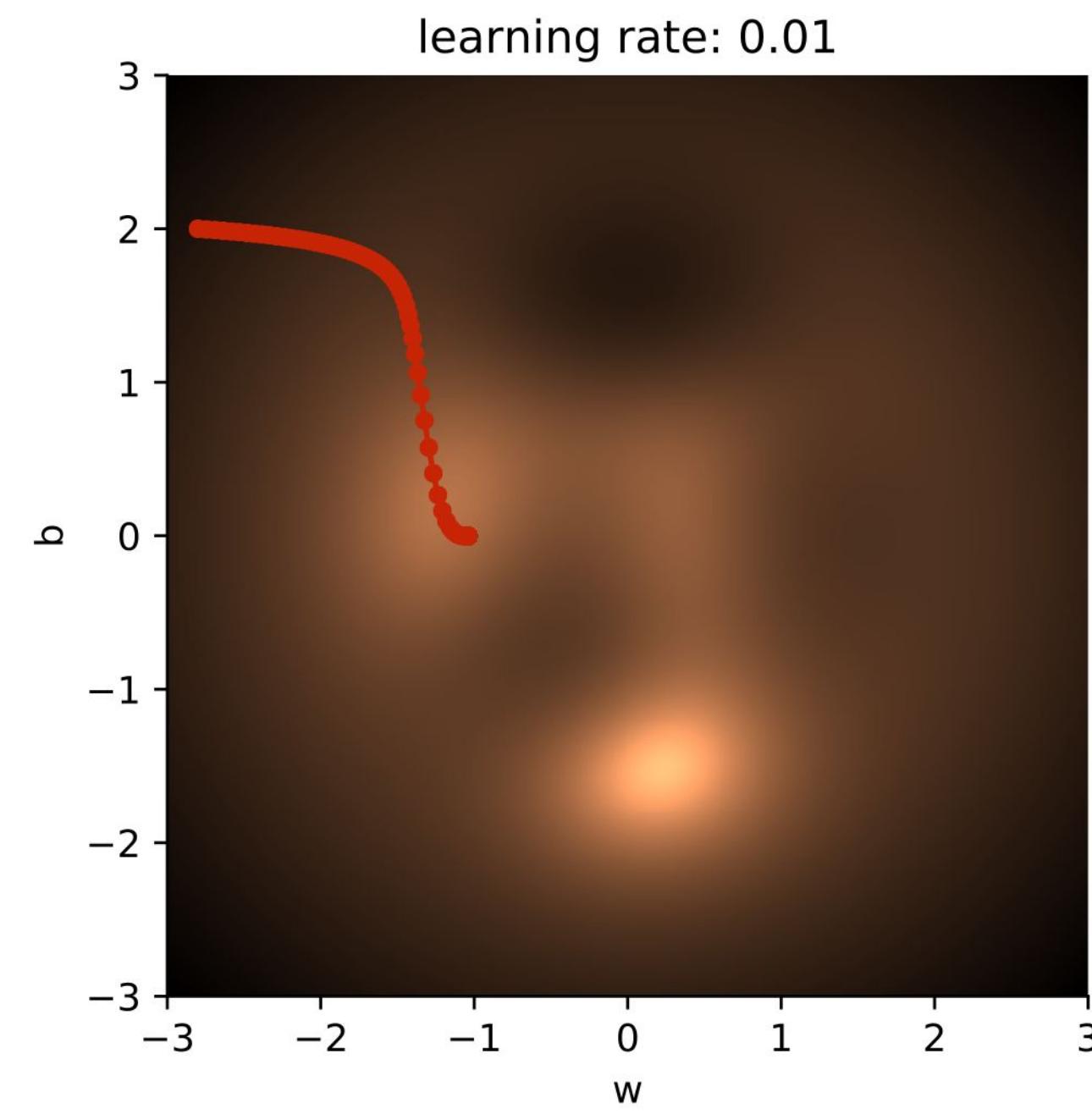
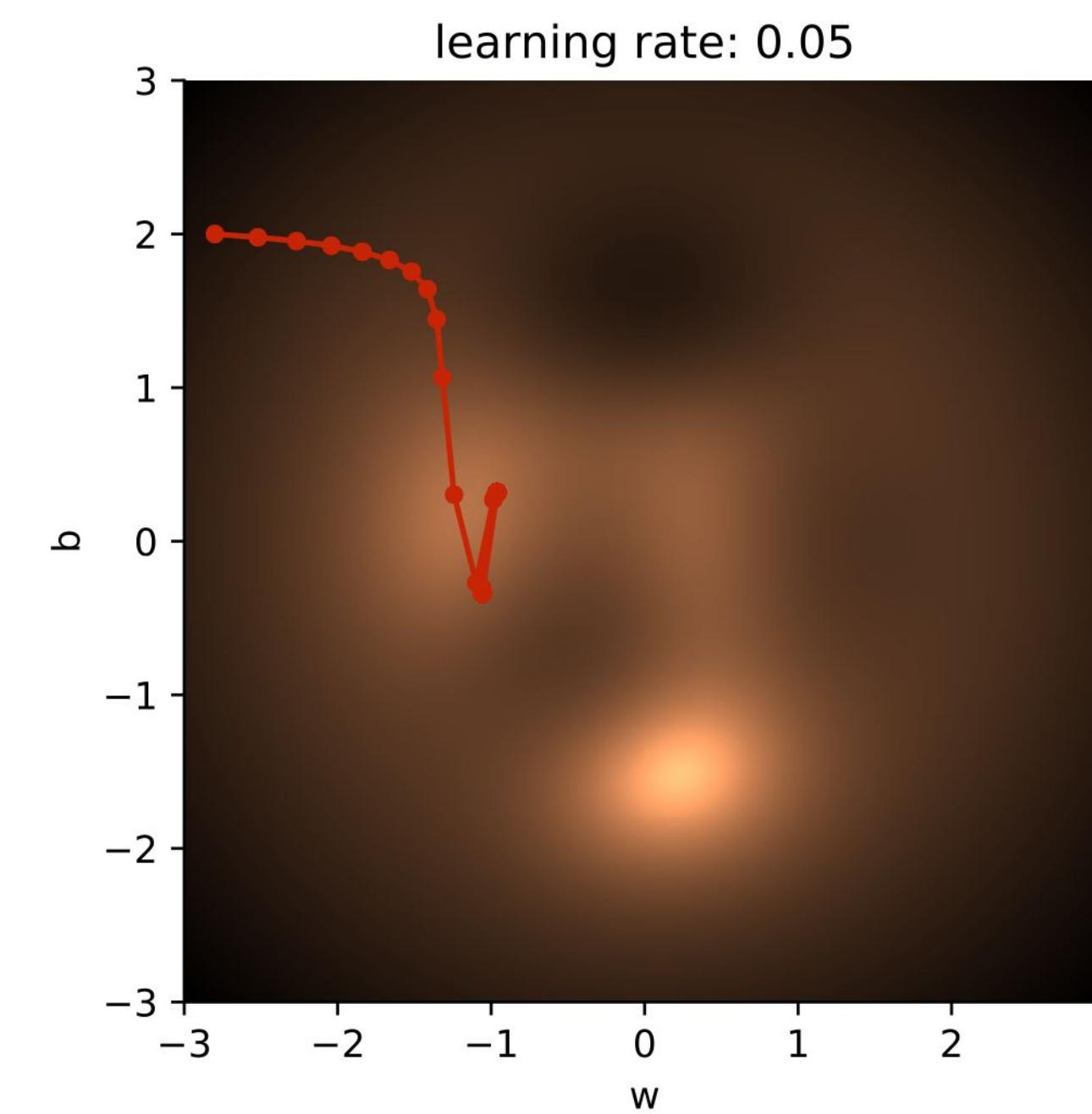
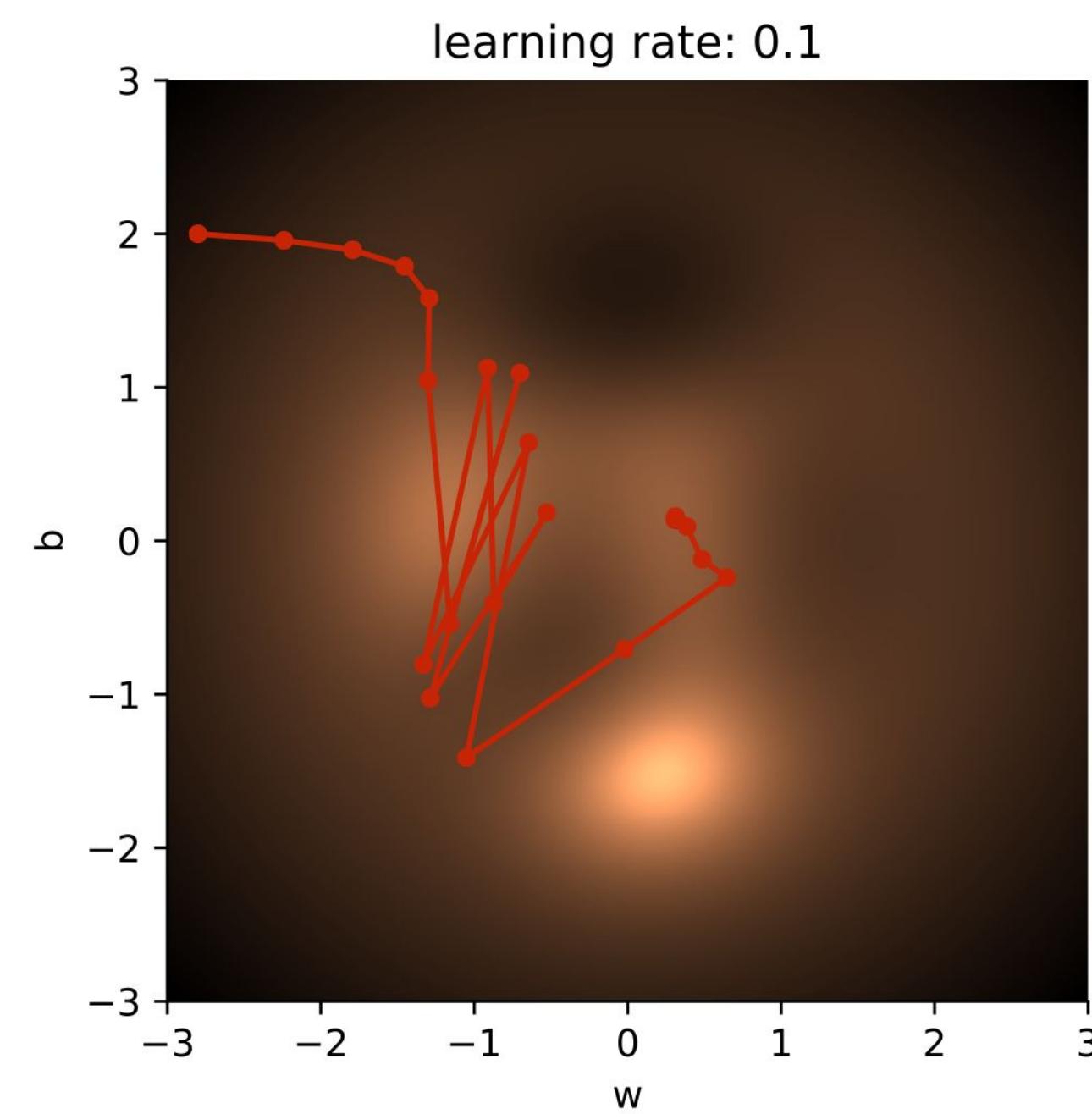
perceptron: linear combination of inputs

neural network: network of perceptrons, with scalar nonlinearities

training: minibatch gradient descent

But, how do we compute the gradient of a complex neural network?

Next video: **backpropagation.**



PART ONE

~~Neural Networks: review~~

Scalar backpropagation: definition

Scalar backpropagation on a simple neural network

Lecture 2: Backpropagation

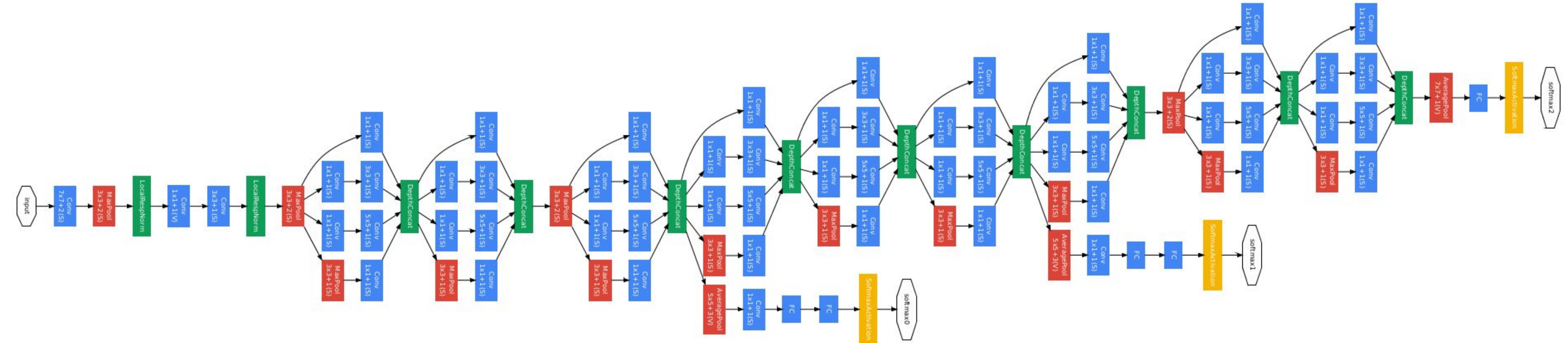
Peter Bloem
Deep Learning 2020

dlvu.github.io



PART TWO: SCALAR BACKPROPAGATION

How do we work out the gradient for a neural network?



GRADIENT COMPUTATION: THE SYMBOLIC APPROACH



derivative of $l(w) = (d + z * 1/(1+e^{-(c + v*(b + 1/(1+e^{-(x*w))))}))$



Σ Extended Keyboard

Upload

Examples

Random

Derivative:

Approximate form

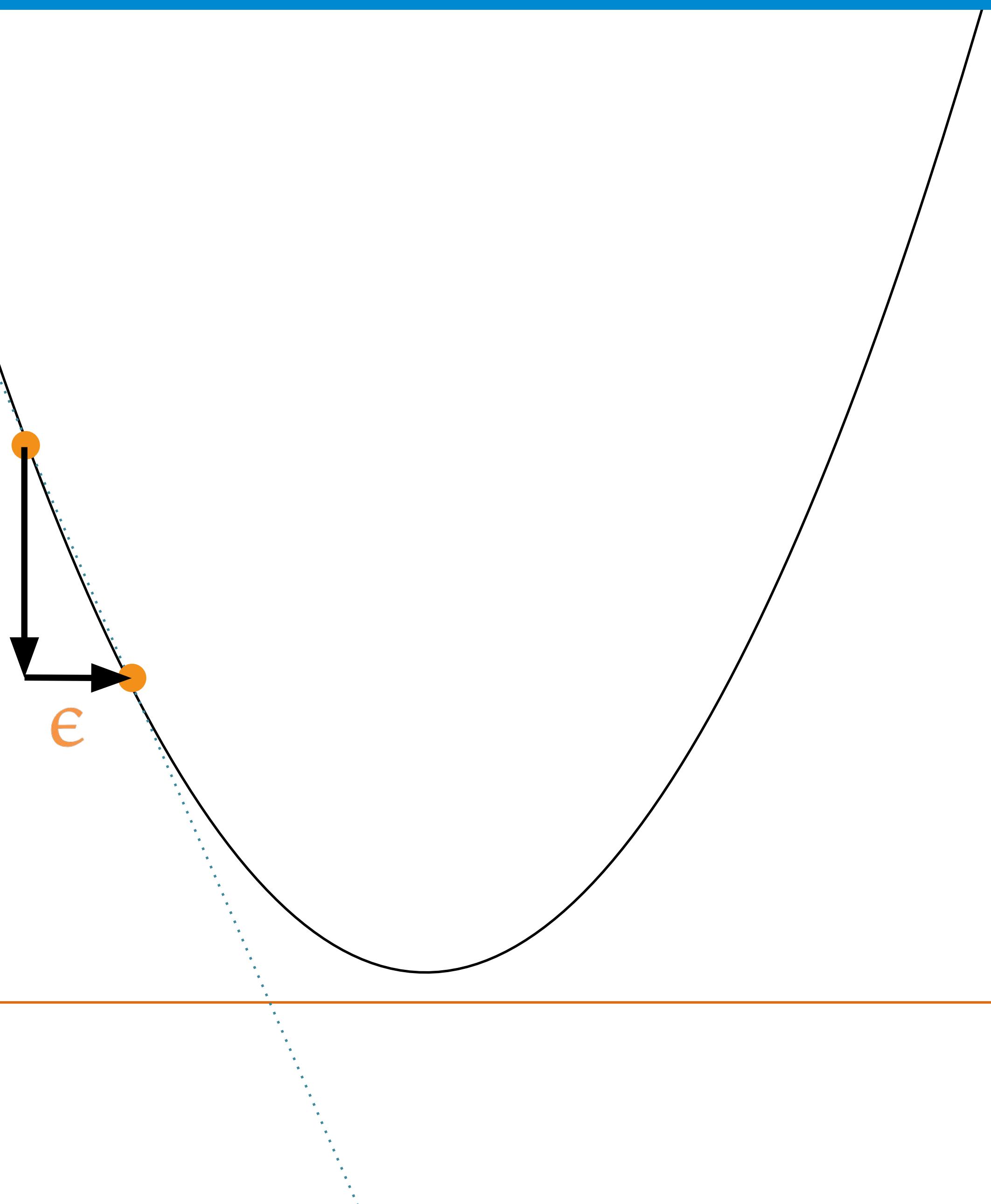
Step-by-step solution

$$\frac{\partial}{\partial w} \left(l(w) = d + \frac{z}{1 + e^{-(c+v(b+1/(1+e^{-(xw)})))}} \right) = \frac{\partial \text{Hold}\left[\frac{z}{e^{-(v(b+1/(e^{-(wx)+1))+c)+1)}} + d \right]}{\partial w}$$

Alternate form:

$$l'(w) = \frac{vxz e^{wx}}{(e^{wx} + 1)^2 (e^{v(b+1/(e^{-wx}+1))+c} + 1)} - \frac{vxz e^{wx}}{(e^{wx} + 1)^2 (e^{v(b+1/(e^{-wx}+1))+c} + 1)^2}$$

GRADIENT COMPUTATION: THE NUMERIC APPROACH



$$\text{deriv} \approx \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

the "method of finite differences"

BACKPROPAGATION: THE MIDDLE GROUND

Work out parts of the derivative **symbolically**
chain these together in a **numeric computation**.

secret ingredient: **the chain rule**.

THE CHAIN RULE

$$\frac{\partial \mathbf{f}(g(x))}{\partial x} = \frac{\partial \mathbf{f}(g(x))}{\partial g(x)} \frac{\partial g(x)}{\partial x}$$



$$\frac{\partial \mathbf{f}}{\partial x} = \frac{\partial \mathbf{f}}{\partial g} \frac{\partial g}{\partial x}$$

INTUITION FOR THE CHAIN RULE

slope of f over g

$$f(g) = s_f g + b_f$$
$$g(x) = s_g x + b_g$$

with $s_f = \frac{\partial f}{\partial g}$

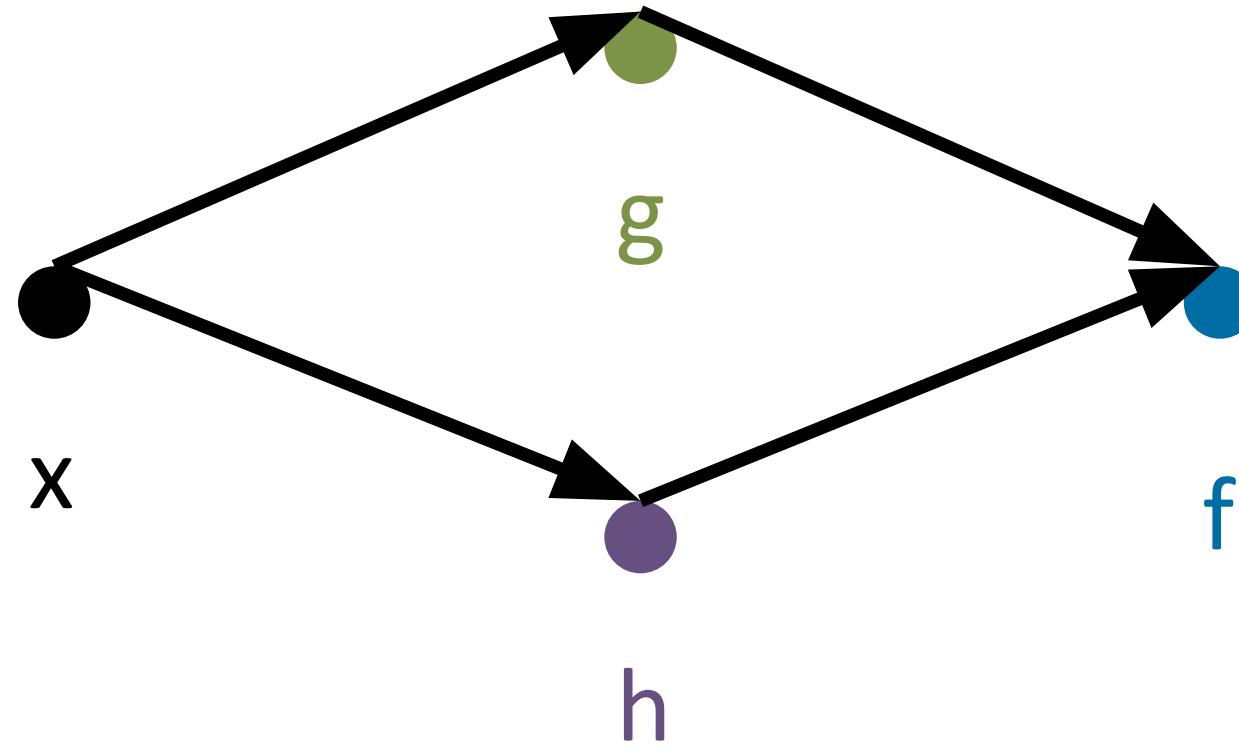
$$s_g = \frac{\partial g}{\partial x}$$

$$\begin{aligned} f(g(x)) &= s_f(s_g x + b_g) + b_f \\ &= s_f s_g x + s_f b_g + s_g \end{aligned}$$



$$\frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

MULTIVARIATE CHAIN RULE



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} + \frac{\partial f}{\partial h} \frac{\partial h}{\partial x}$$

INTUITION FOR THE MULTIVARIATE CHAIN RULE

$$f(g, h) = s_1 g + s_2 h + b_f$$

$$g(x) = s_g x + b_g$$

$$h(x) = s_h x + b_h$$

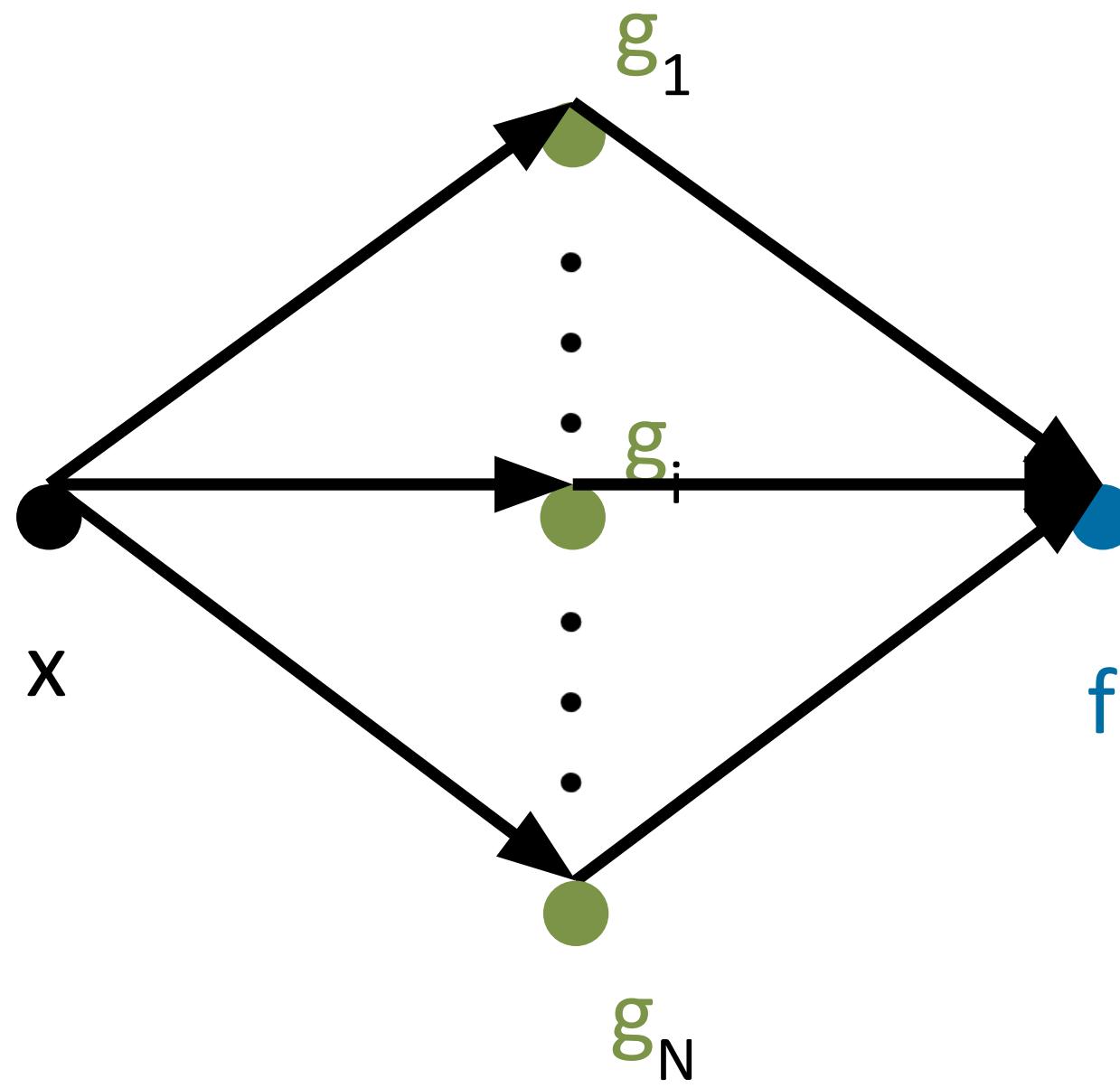
with $s_1 = \frac{\partial f}{\partial g}$ $s_2 = \frac{\partial f}{\partial h}$ $s_g = \frac{\partial g}{\partial x}$ $s_h = \frac{\partial h}{\partial x}$

$$\begin{aligned} f(g(x), h(x)) &= s_1(s_g x + b_g) + s_2(s_h x + b_h) + b_f \\ &= s_1 s_g x + s_1 b_g + s_2 s_h x + s_2 b_h + b_f \\ &= (s_1 s_g + s_2 s_h)x + s_1 b_g + s_2 b_h + b_f \end{aligned}$$



$$\frac{\partial f}{\partial g} \frac{\partial g}{\partial x} + \frac{\partial f}{\partial h} \frac{\partial h}{\partial x}$$

MULTIVARIATE CHAIN RULE



$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \sum_i \frac{\partial \mathbf{f}}{\partial \mathbf{g}_i} \frac{\partial \mathbf{g}_i}{\partial \mathbf{x}}$$

EXAMPLE

$$f(x) = \frac{2}{\sin(e^{-x})}$$

ITERATING THE CHAIN RULE

$$f(x) = \frac{2}{\sin(e^{-x})}$$

$$f(x) = d(c(b(a(x))))$$

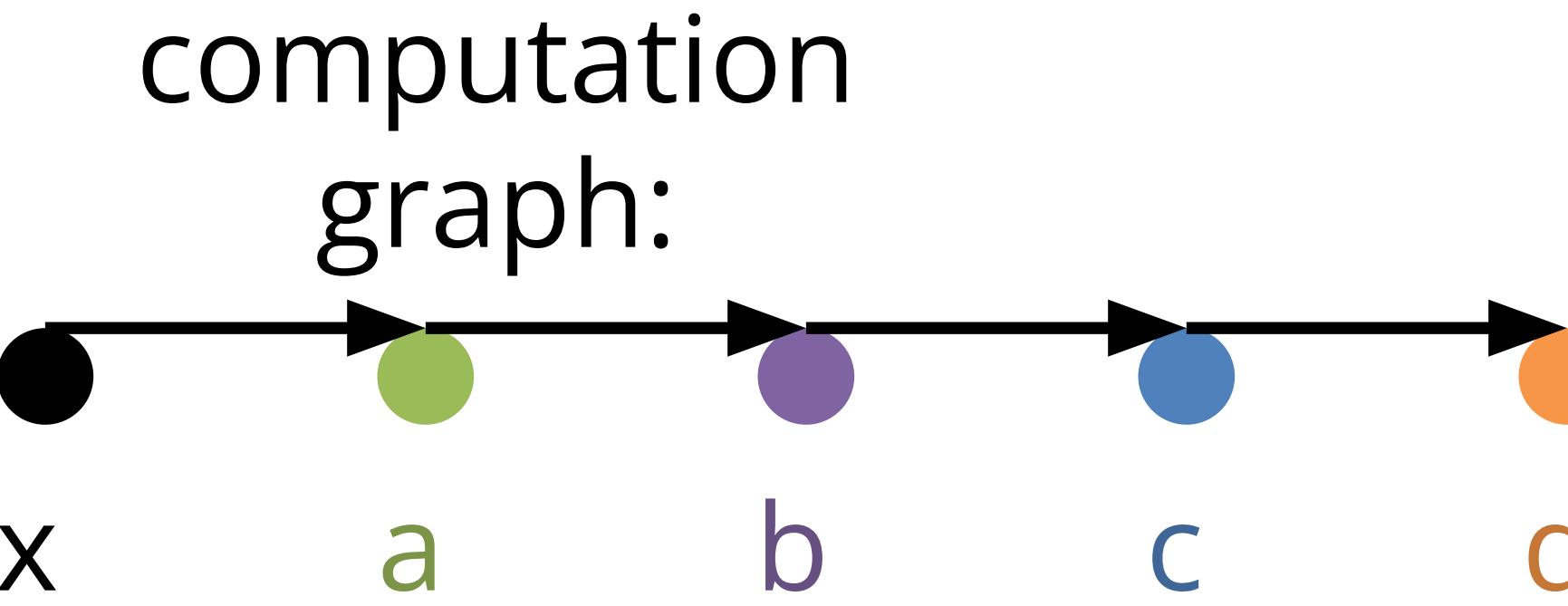
operations
:

$$d(c) = \frac{2}{c}$$

$$c(b) = \sin b$$

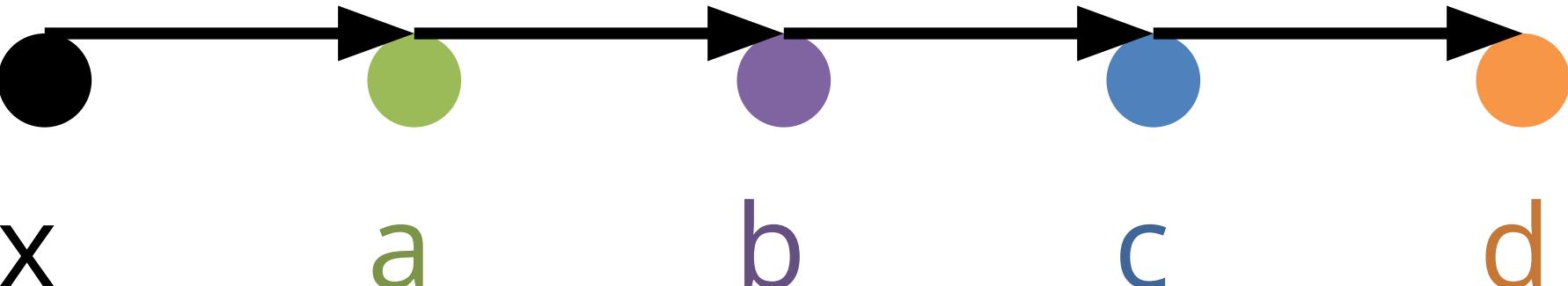
$$b(a) = e^a$$

$$a(x) = -x$$



ITERATING THE CHAIN RULE

$$\begin{aligned}\frac{\partial f}{\partial x} &= \frac{\partial d}{\partial x} \\&= \frac{\partial d}{\partial c} \frac{\partial c}{\partial x} \\&= \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial x} \\&= \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}\end{aligned}$$



GLOBAL AND LOCAL DERIVATIVES

$$\frac{\partial f}{\partial x} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

global derivative local derivatives

BACKPROPAGATION

The BACKPROPAGATION algorithm:

- break your computation up into a sequence of operations
what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically
by computing the local derivatives and multiplying them

WORK OUT THE LOCAL DERIVATIVES SYMBOLICALLY

$$f(x) = \frac{2}{\sin(e^{-x})}$$

operations
:

$$d(c) = \frac{2}{c}$$

$$c(b) = \sin b$$

$$b(a) = e^a$$

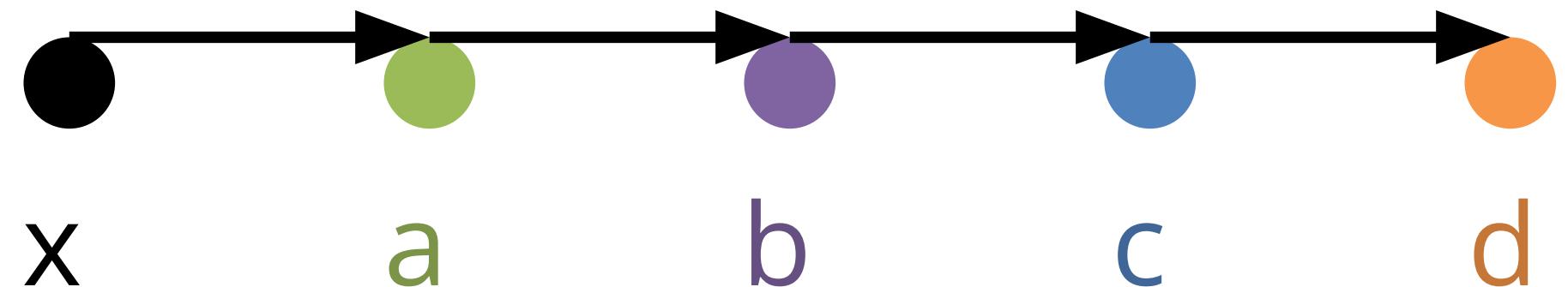
$$a(x) = -x$$

$$\frac{\partial f}{\partial x} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$= -\frac{2}{c^2} \cdot \cos b \cdot e^a \cdot -1$$

COMPUTE A FORWARD PASS ($X = -4.499$)

$$f(-4.499) = 2$$



$$d = \frac{2}{c} = 2$$

$$c = \sin b = 1$$

$$b = e^a = 90$$

$$a = -x = 4.499$$

COMPUTE A FORWARD PASS ($X = -4.499$)

$$f(-4.499) = 2$$

$$\frac{\partial f}{\partial x} = -\frac{2}{c^2} \cdot \cos b \cdot e^a \cdot -1$$

$$d = \frac{2}{c} = 2$$

$$= -\frac{2}{1^2} \cdot \cos 90 \cdot e^{4.499} \cdot -1$$

$$c = \sin b = 1$$

$$= -\frac{1}{2} \cdot 0 \cdot 90 \cdot -1 = 0$$

$$b = e^a = 90$$

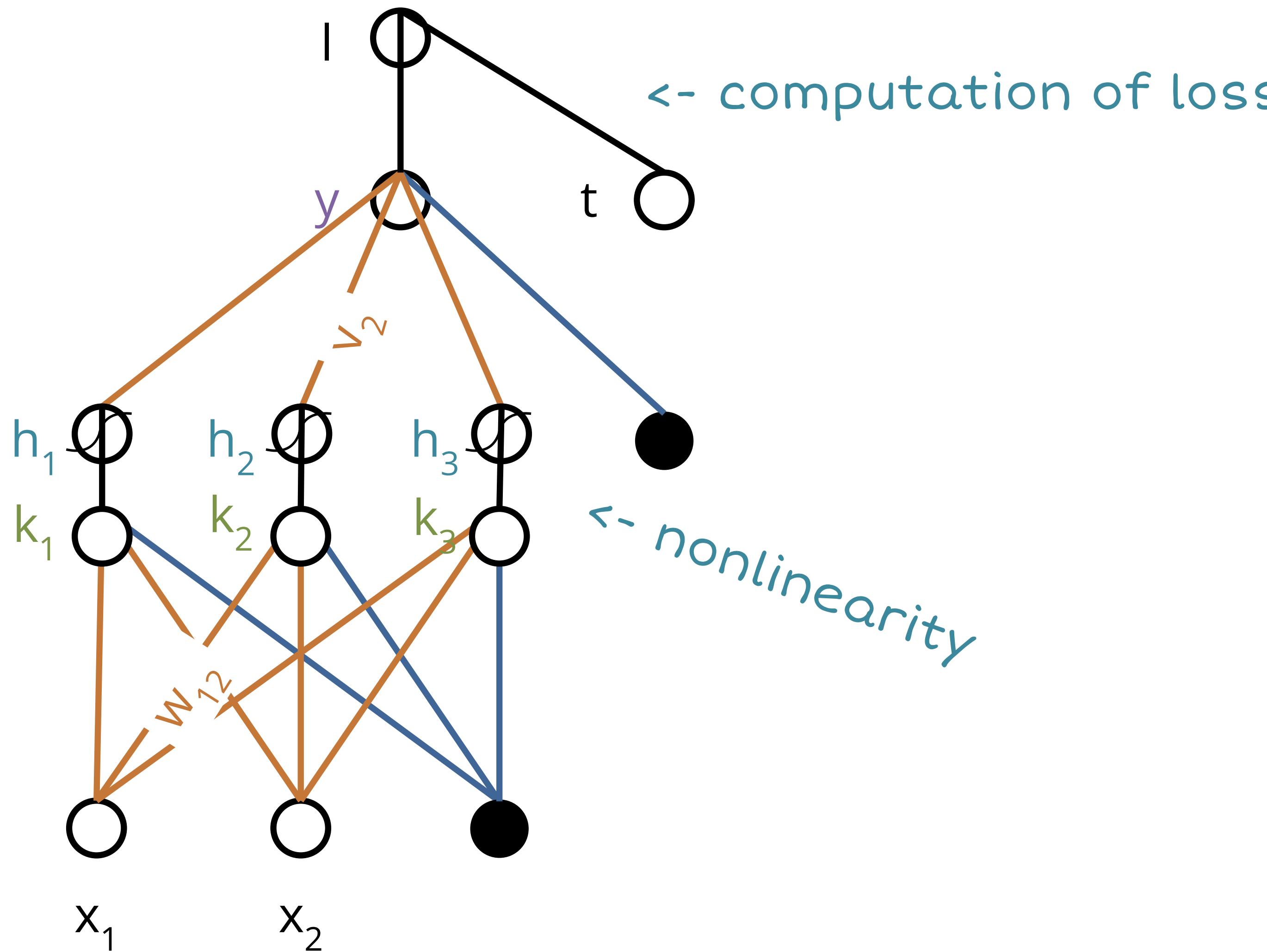
$$a = -x = 4.499$$

BACKPROPAGATION

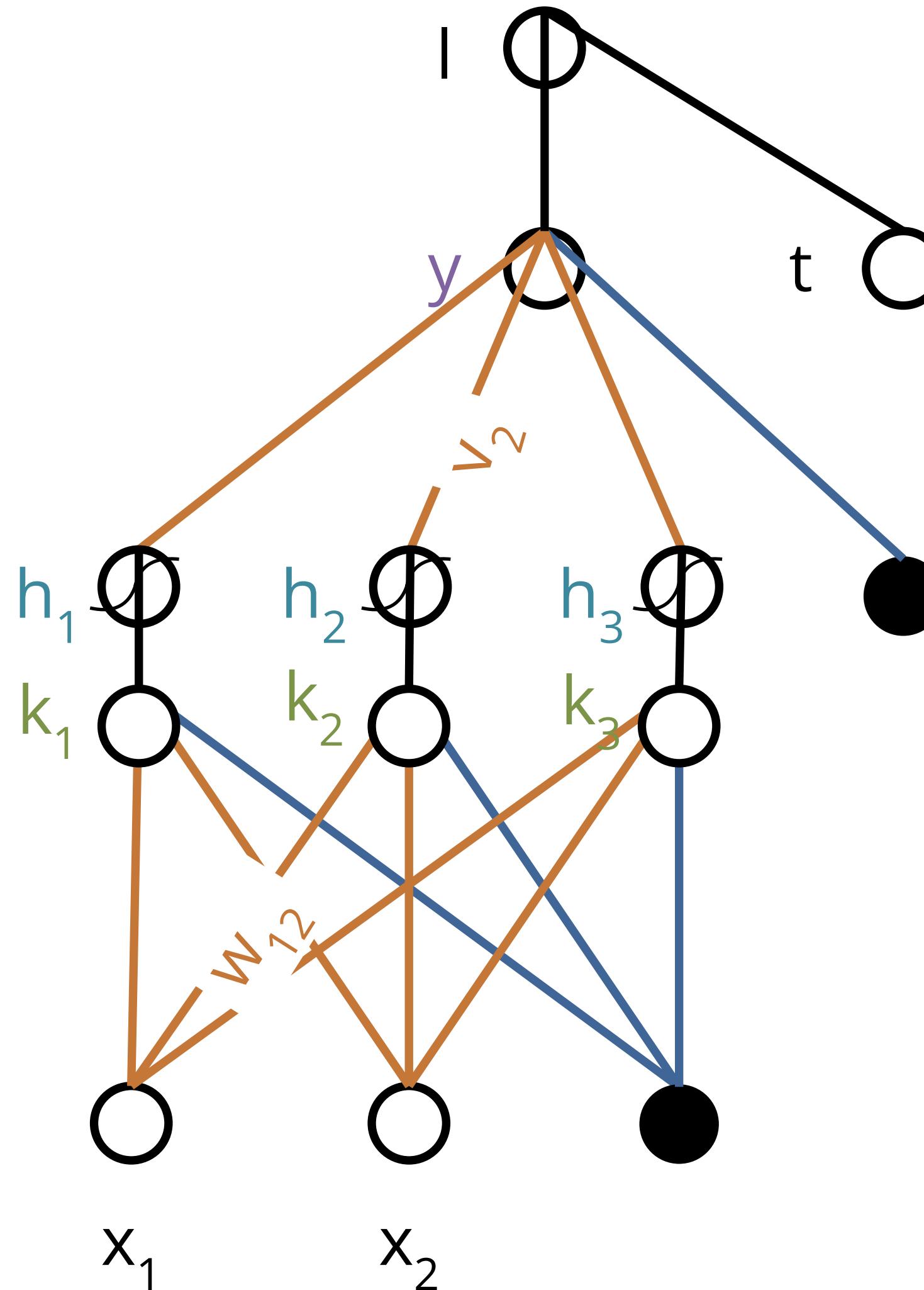
The BACKPROPAGATION algorithm:

- break your computation up into a sequence of operations
what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically
by computing the local derivatives and multiplying them
- Much more accurate than finite differences
only source of inaccuracy is the numeric computation of the operations.
- Much faster than symbolic differentiation
The backward pass has (broadly) the same complexity as the forward.

BACKPROPAGATION IN A FEEDFORWARD NETWORK



BACKPROPAGATION IN A FEEDFORWARD NETWORK



$$\frac{\partial l}{\partial v_2}, \frac{\partial l}{\partial w_{12}}$$

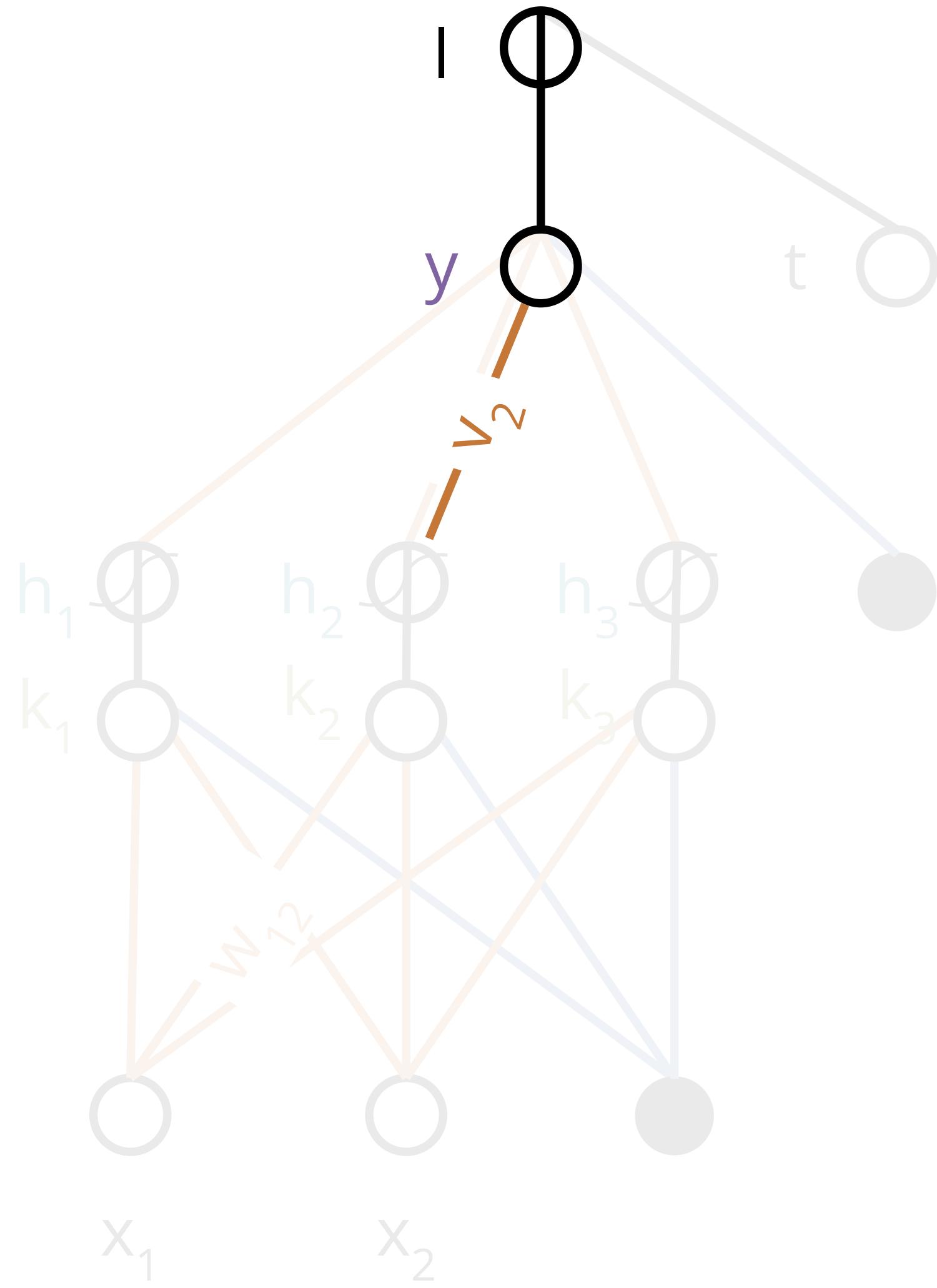
$$l = (y - t)^2$$

$$y = v_1 h_1 + v_2 h_2 + v_3 h_3 + b_h$$

$$h_2 = \frac{1}{1 + \exp(-k_2)}$$

$$k_2 = w_{21}x_1 + w_{22}x_2 + b_x$$

BACKPROPAGATION IN A FEEDFORWARD NETWORK

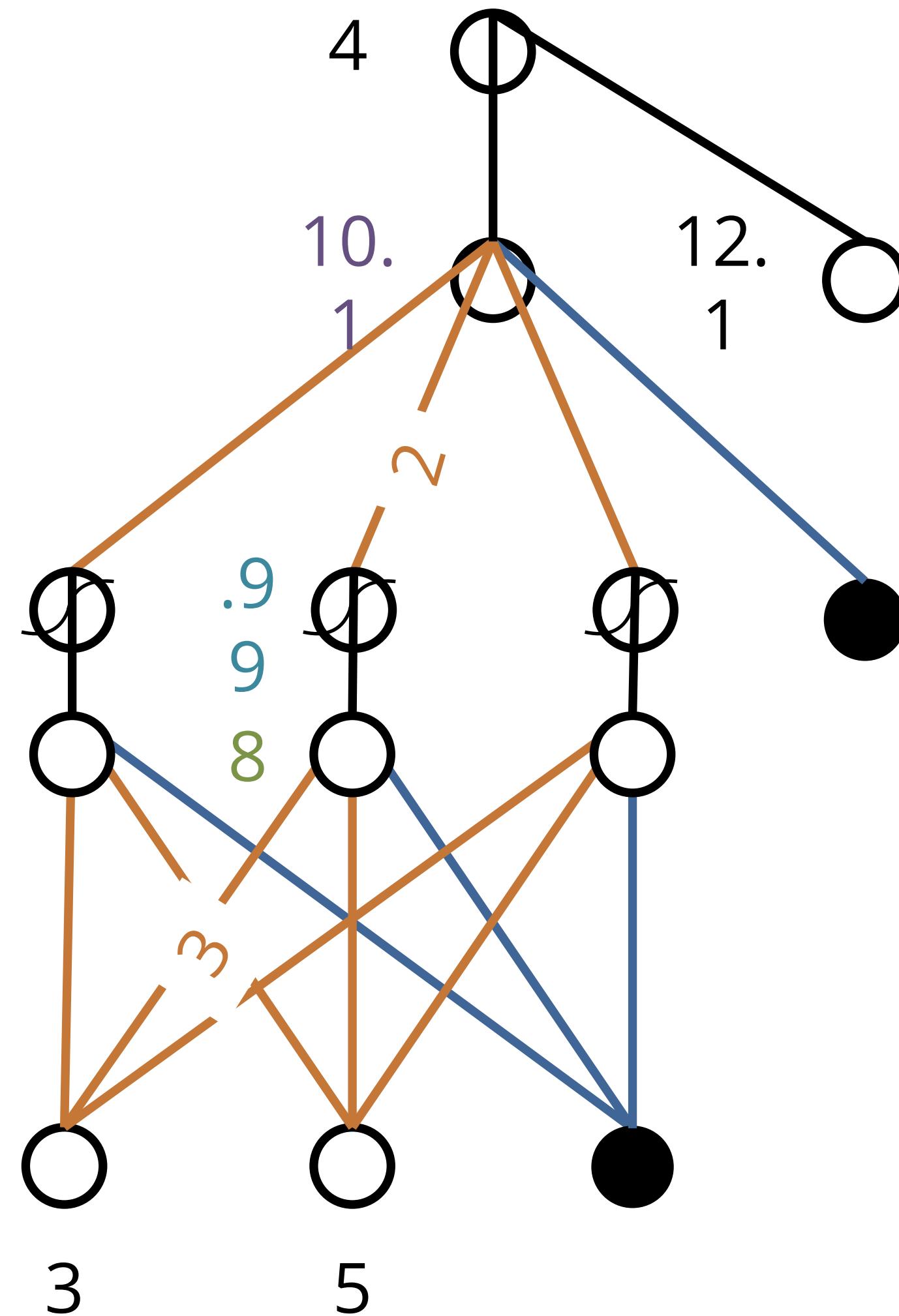


$$l = (\textcolor{violet}{y} - t)^2$$

$$y = \textcolor{brown}{v}_1 h_1 + \textcolor{brown}{v}_2 h_2 + \textcolor{brown}{v}_3 h_3 + b_h$$

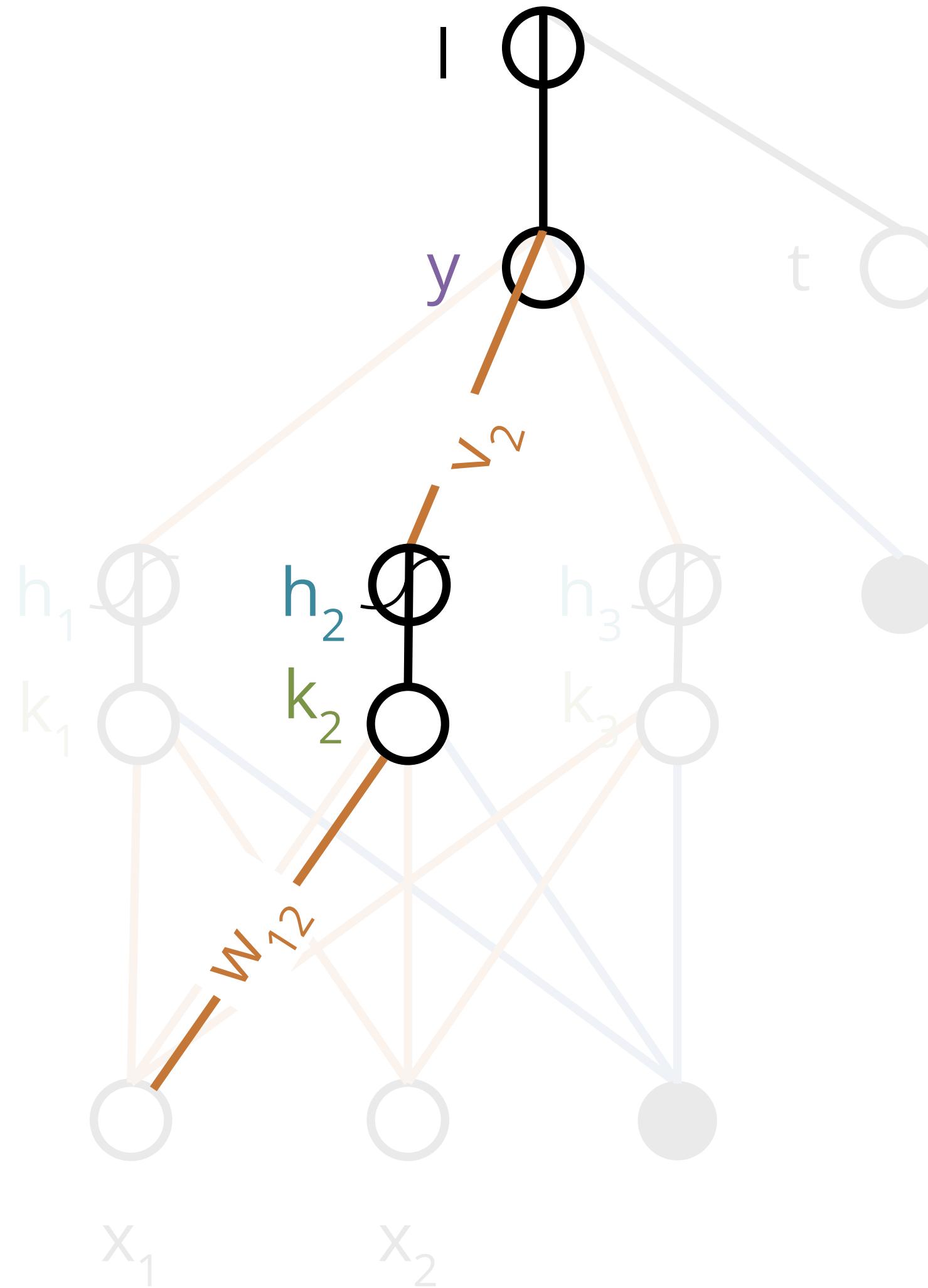
$$\begin{aligned}\frac{\partial l}{\partial \textcolor{brown}{v}_2} &= \frac{\partial l}{\partial y} \frac{\partial y}{\partial \textcolor{brown}{v}_2} \\ &= 2(\textcolor{violet}{y} - t) \cdot \textcolor{teal}{h}_2\end{aligned}$$

BACKPROPAGATION IN A FEEDFORWARD NETWORK



$$\begin{aligned}\frac{\partial l}{\partial v_2} &= \frac{\partial l}{\partial y} \frac{\partial y}{\partial v_2} \\ &= 2(y - t) \cdot h_2 \\ &= 2(10.1 - 12.1) \cdot .99 = -3.96 \\ v_2 &\leftarrow v_2 - \alpha \cdot -3.96\end{aligned}$$

BACKPROPAGATION IN A FEEDFORWARD NETWORK

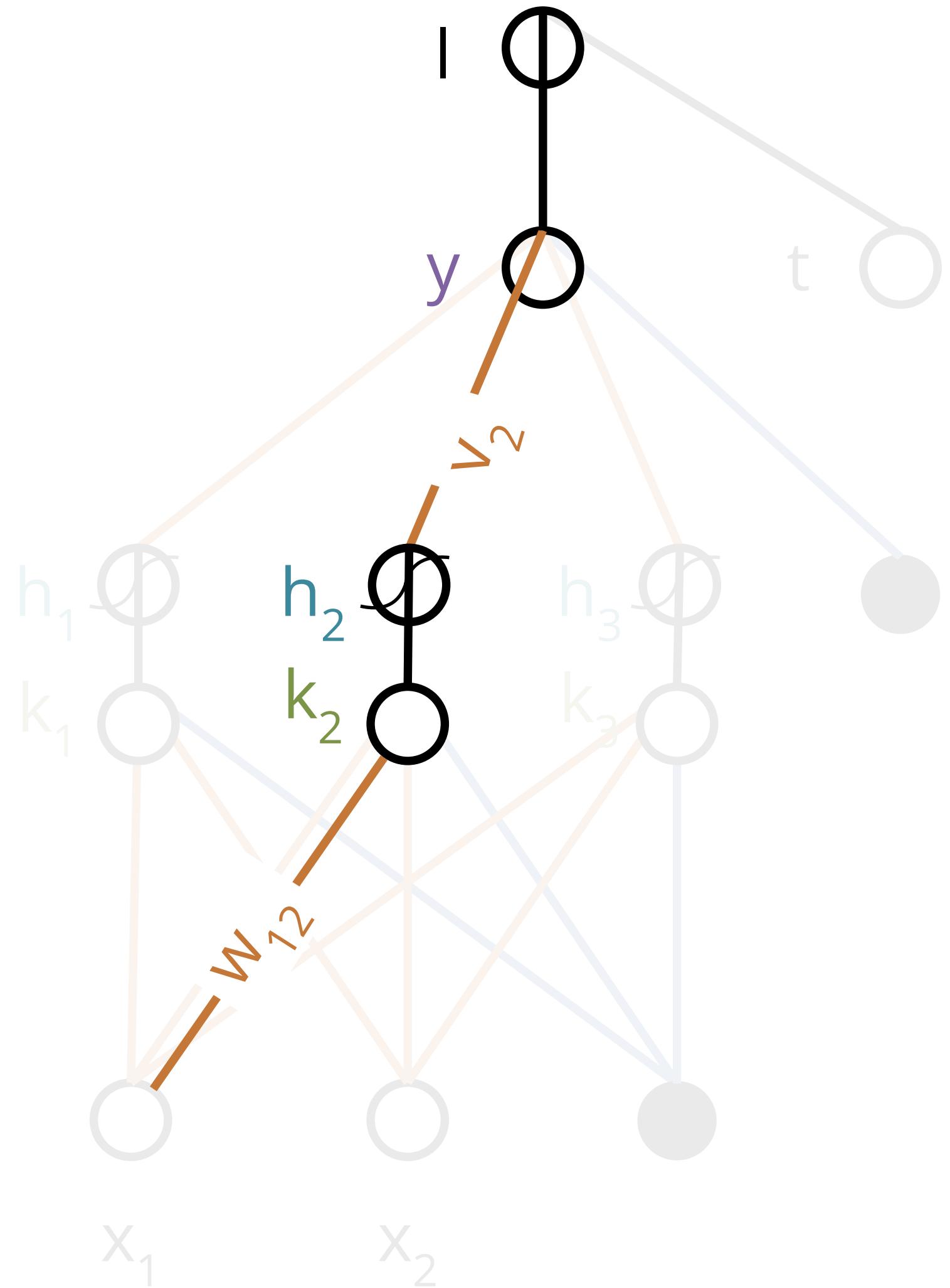


$$\begin{aligned} l &= (y - t)^2 \\ y &= v_1 h_1 + v_2 h_2 + v_3 h_3 + b_h \\ h_2 &= \frac{1}{1 + \exp(-k_2)} \\ k_2 &= w_{21} x_1 + w_{22} x_2 + b_x \end{aligned}$$

$$\begin{aligned} \frac{\partial l}{\partial w_{12}} &= \frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} \frac{\partial k_2}{\partial w_{12}} \\ &= 2(y - t) \cdot v_2 \cdot h_2(1 - h_2) \cdot x_1 \end{aligned}$$

derivative of sigmoid: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

BACKPROPAGATION



$$\frac{\partial l}{\partial y}$$

$$\frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} = \frac{\partial l}{\partial h_2}$$

$$\frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} = \frac{\partial l}{\partial k_2}$$

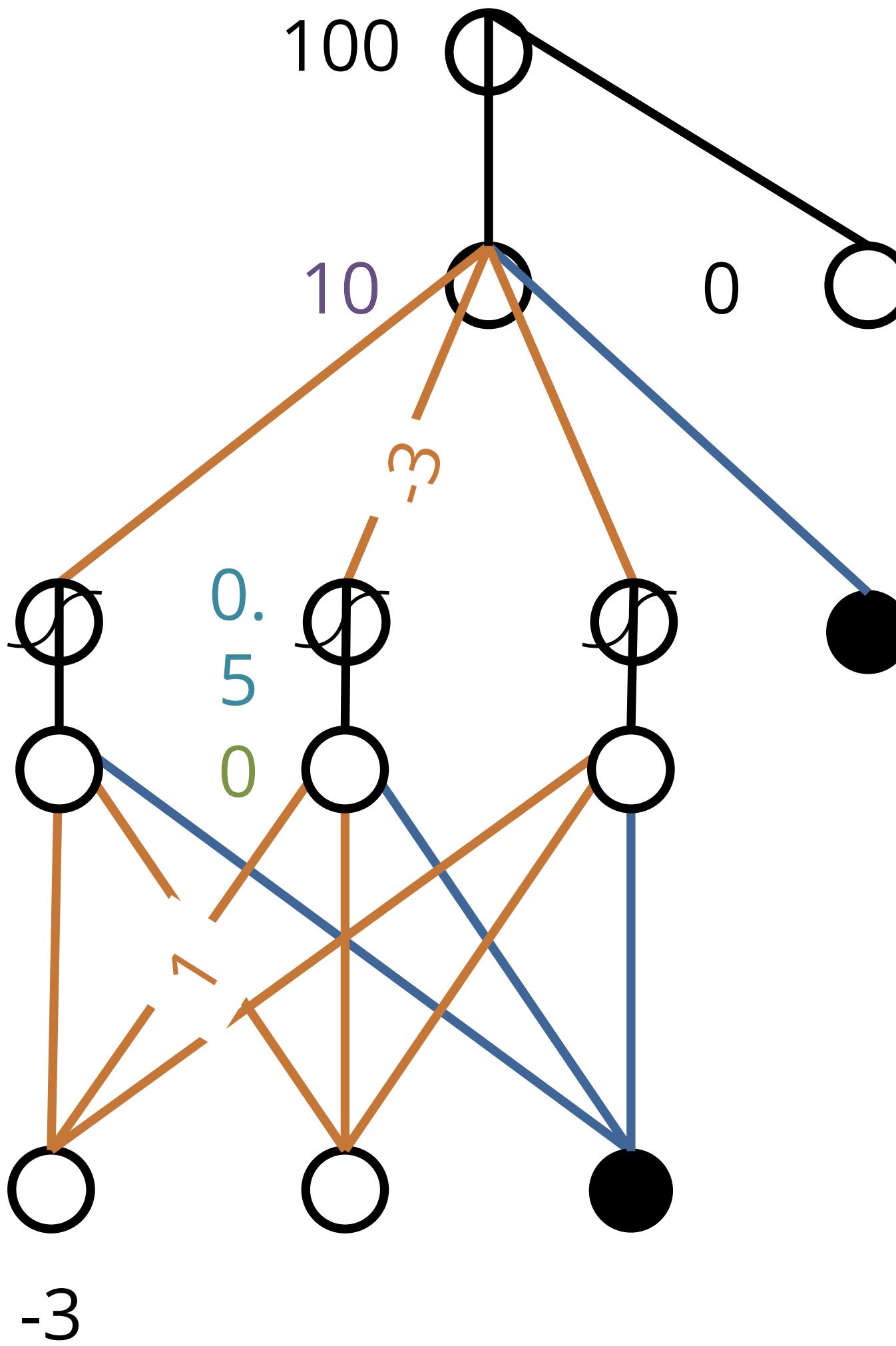
$$\frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} \frac{\partial k_2}{\partial w_{12}} = \frac{\partial l}{\partial w_{12}}$$

BACKPROPAGATION

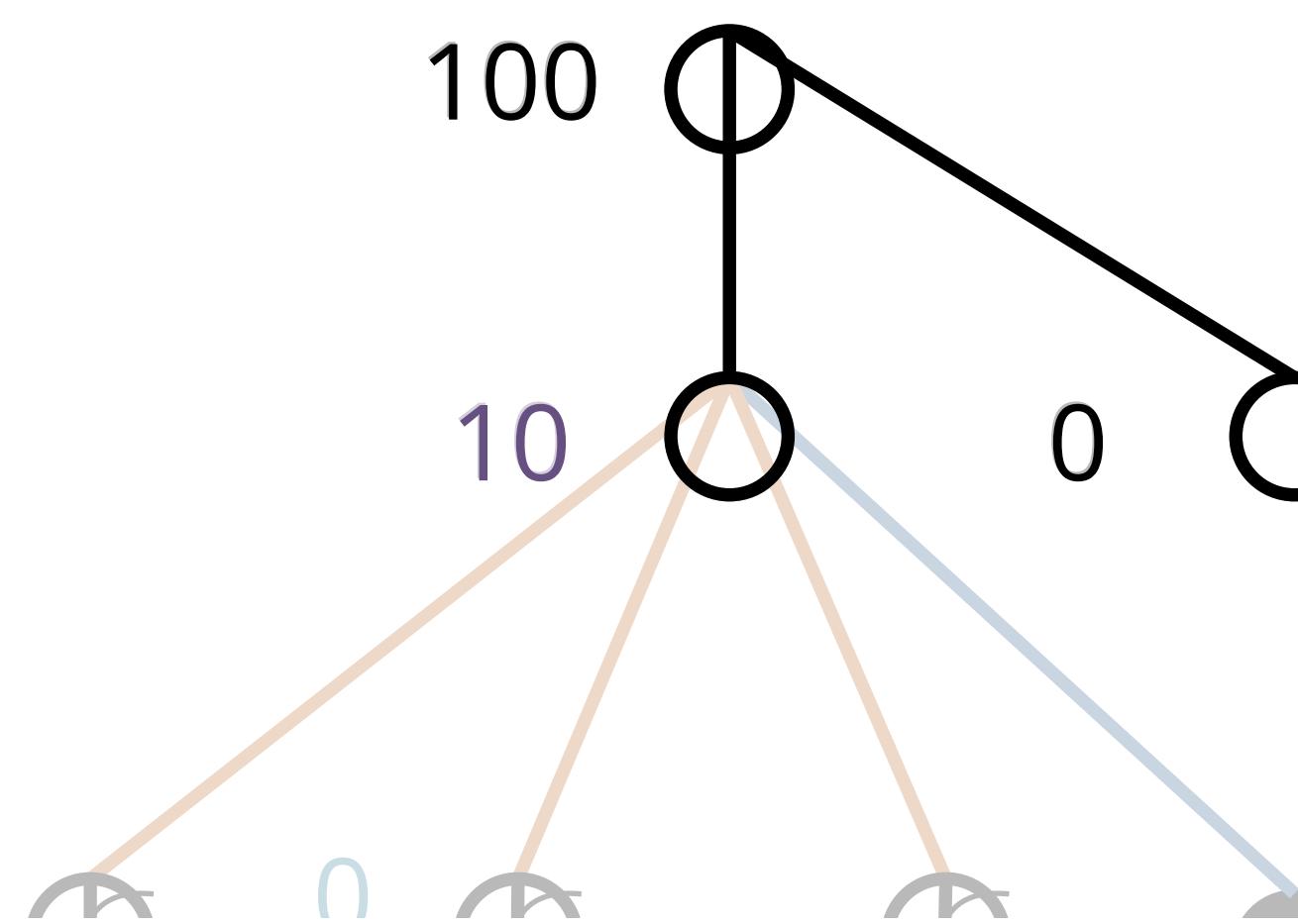
The BACKPROPAGATION algorithm:

- break your computation up into a sequence of operations
what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically
by computing the local derivatives and multiplying them
- Walk backward from the loss, *accumulating* the derivatives.

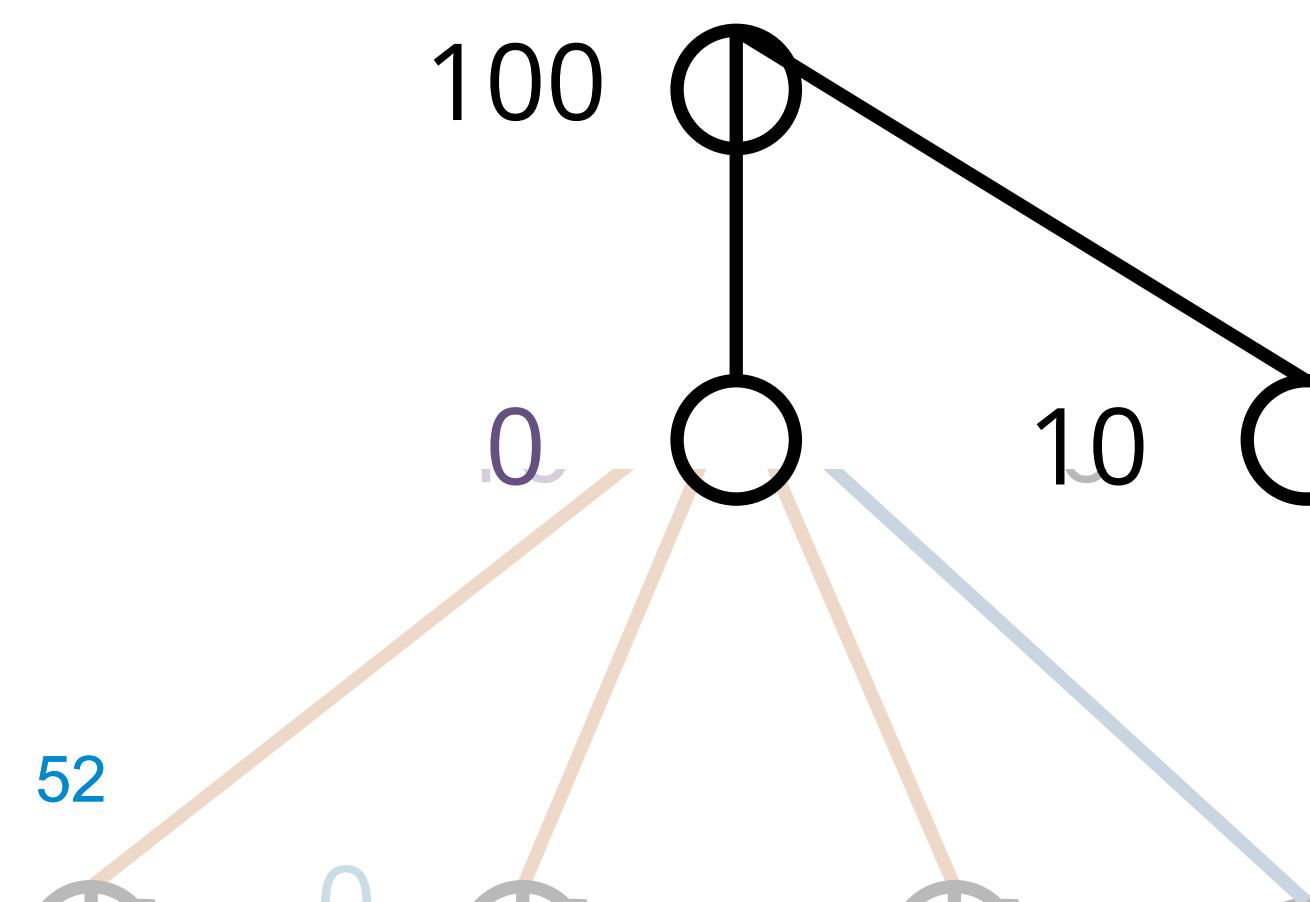
BUILDING SOME INTUITION



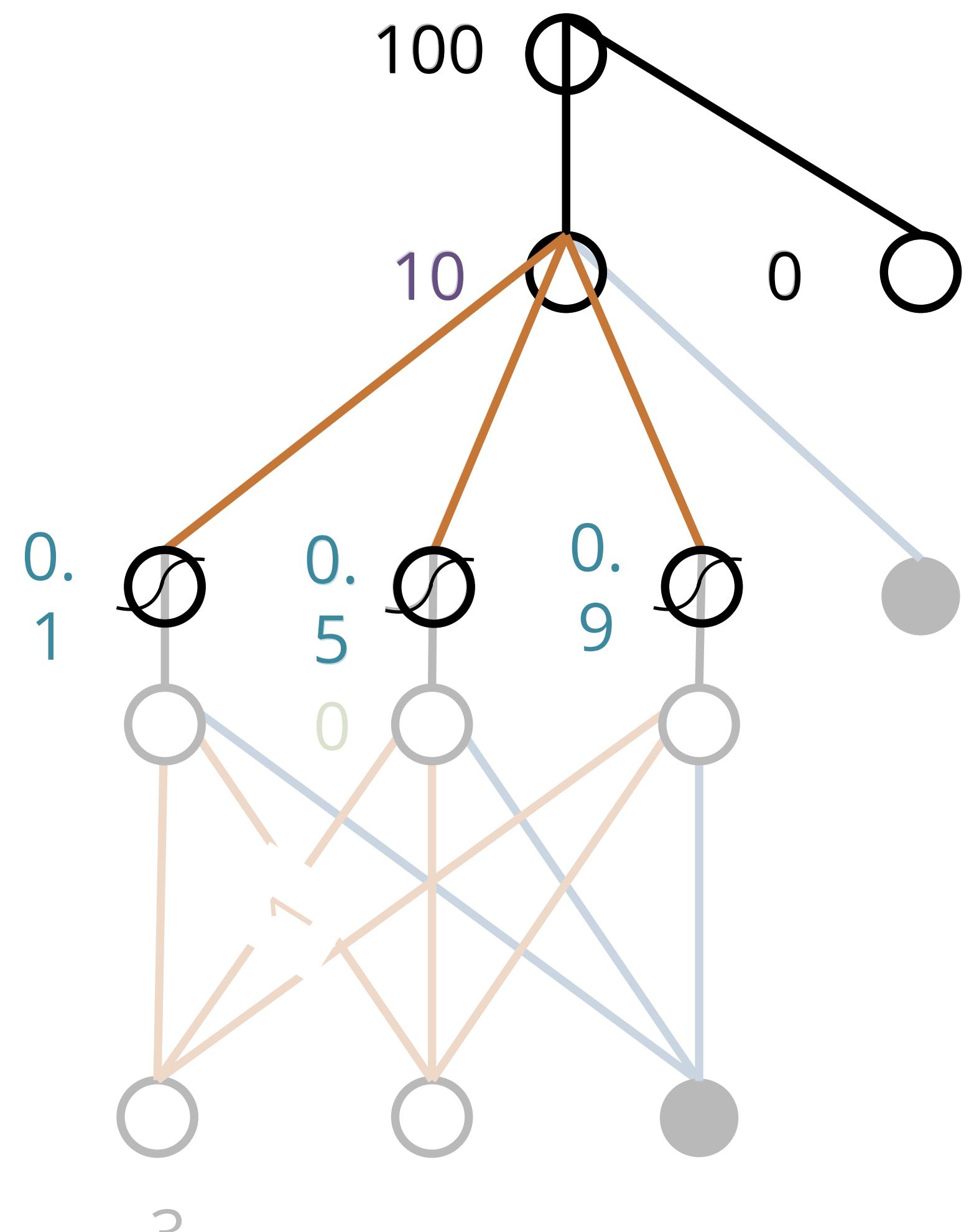
IMAGINE THAT y IS A PARAMETER



$$\begin{aligned}y &\leftarrow y - \alpha \cdot 2(y - t) \\&= y - \alpha \cdot 20\end{aligned}$$



$$\begin{aligned}y &\leftarrow y - \alpha \cdot 2(y - t) \\&= y + \alpha \cdot 20\end{aligned}$$



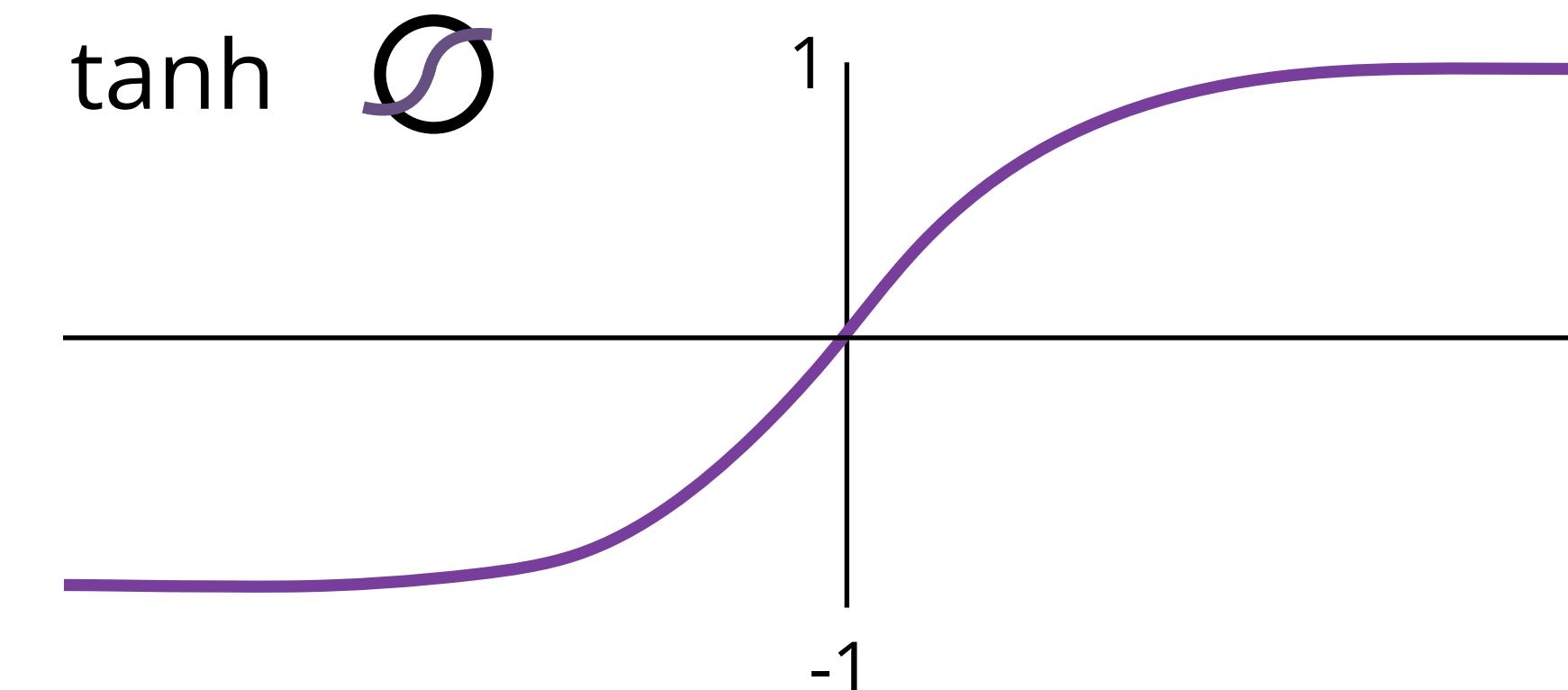
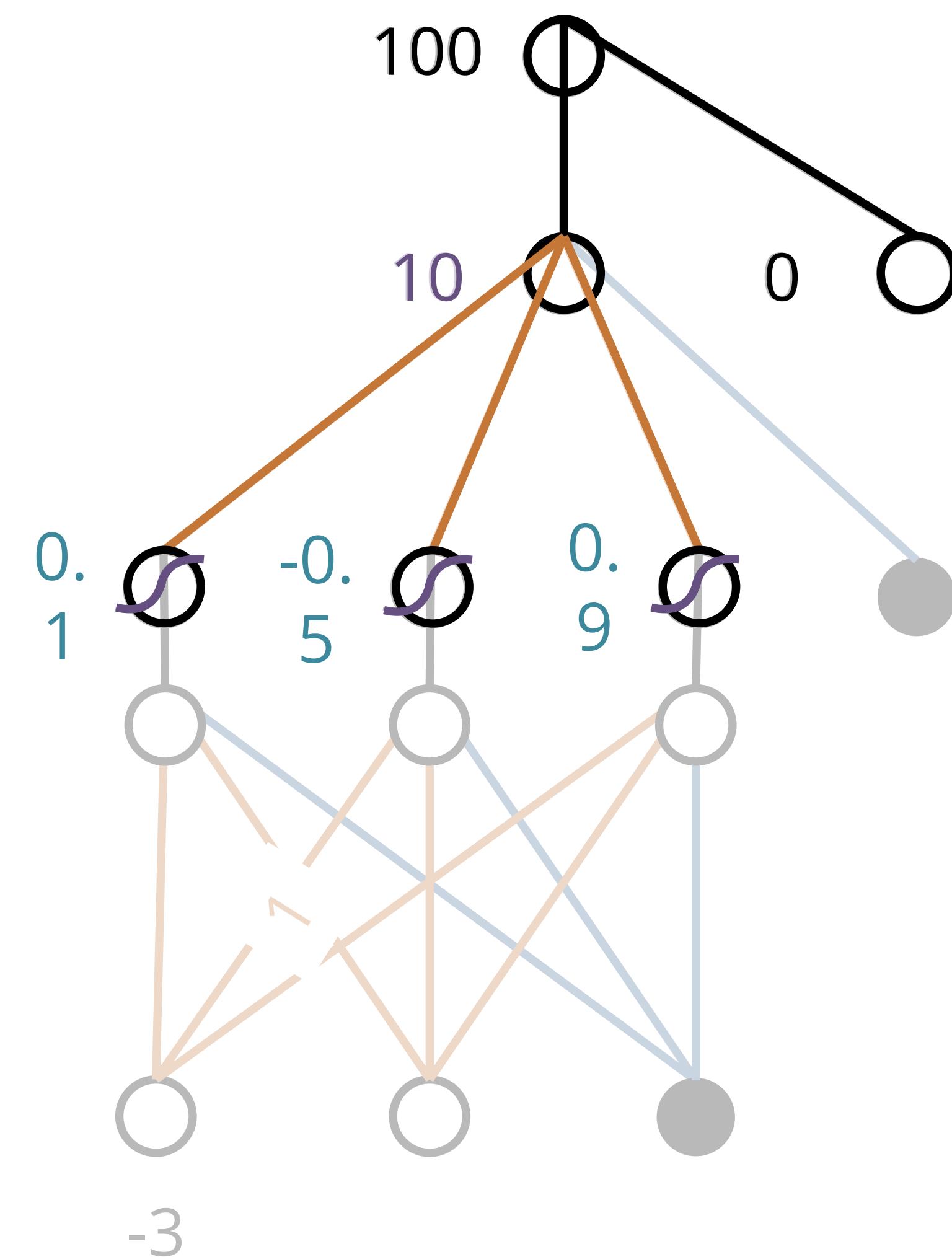
$$v_2 \leftarrow v_2 - \alpha \cdot 2(y - t)h_2$$

$$v_1 \leftarrow v_1 - \alpha \cdot 20 \cdot 0.1$$

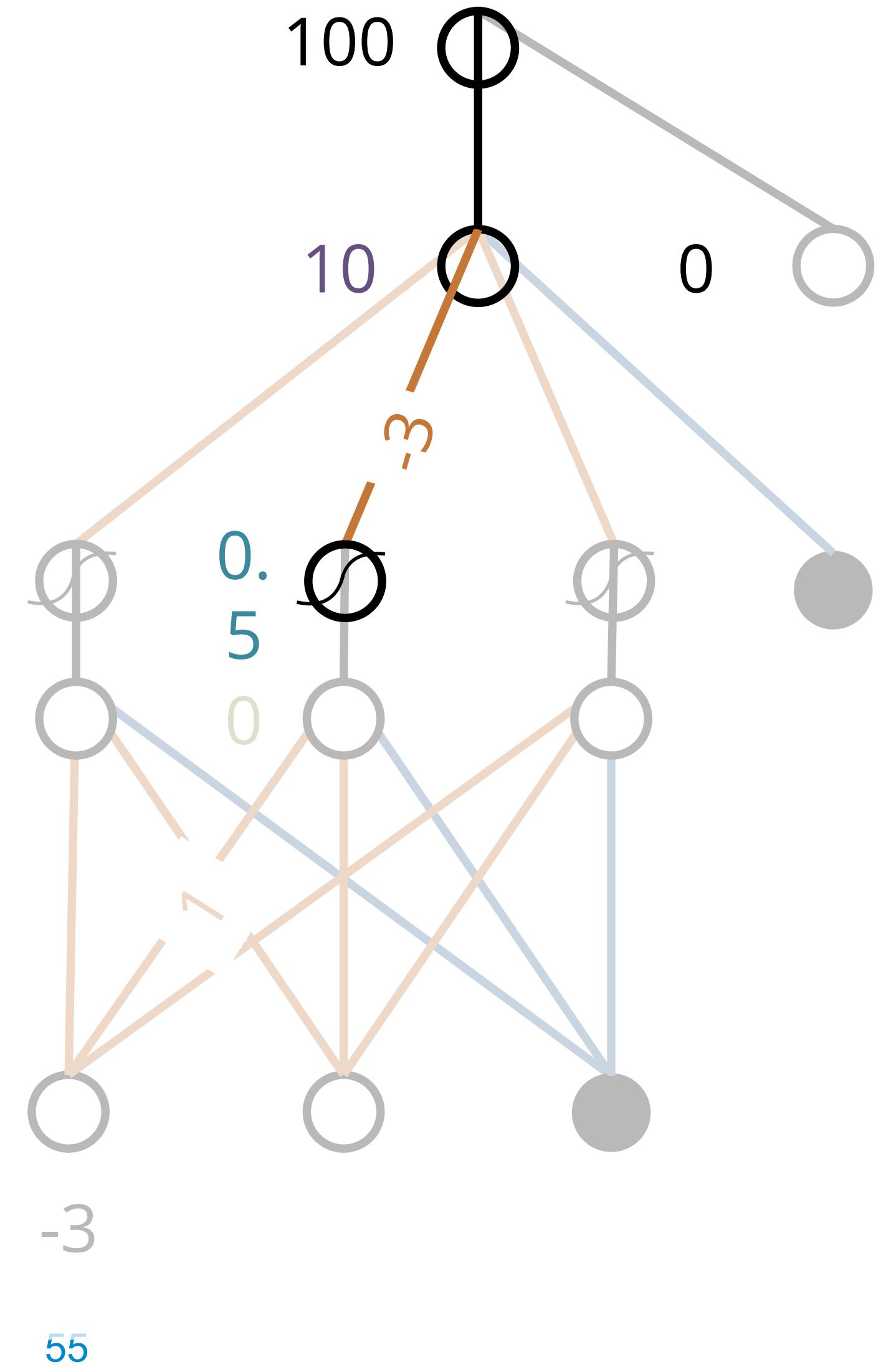
$$v_2 \leftarrow v_2 - \alpha \cdot 20 \cdot 0.5$$

$$v_3 \leftarrow v_3 - \alpha \cdot 20 \cdot 0.9$$

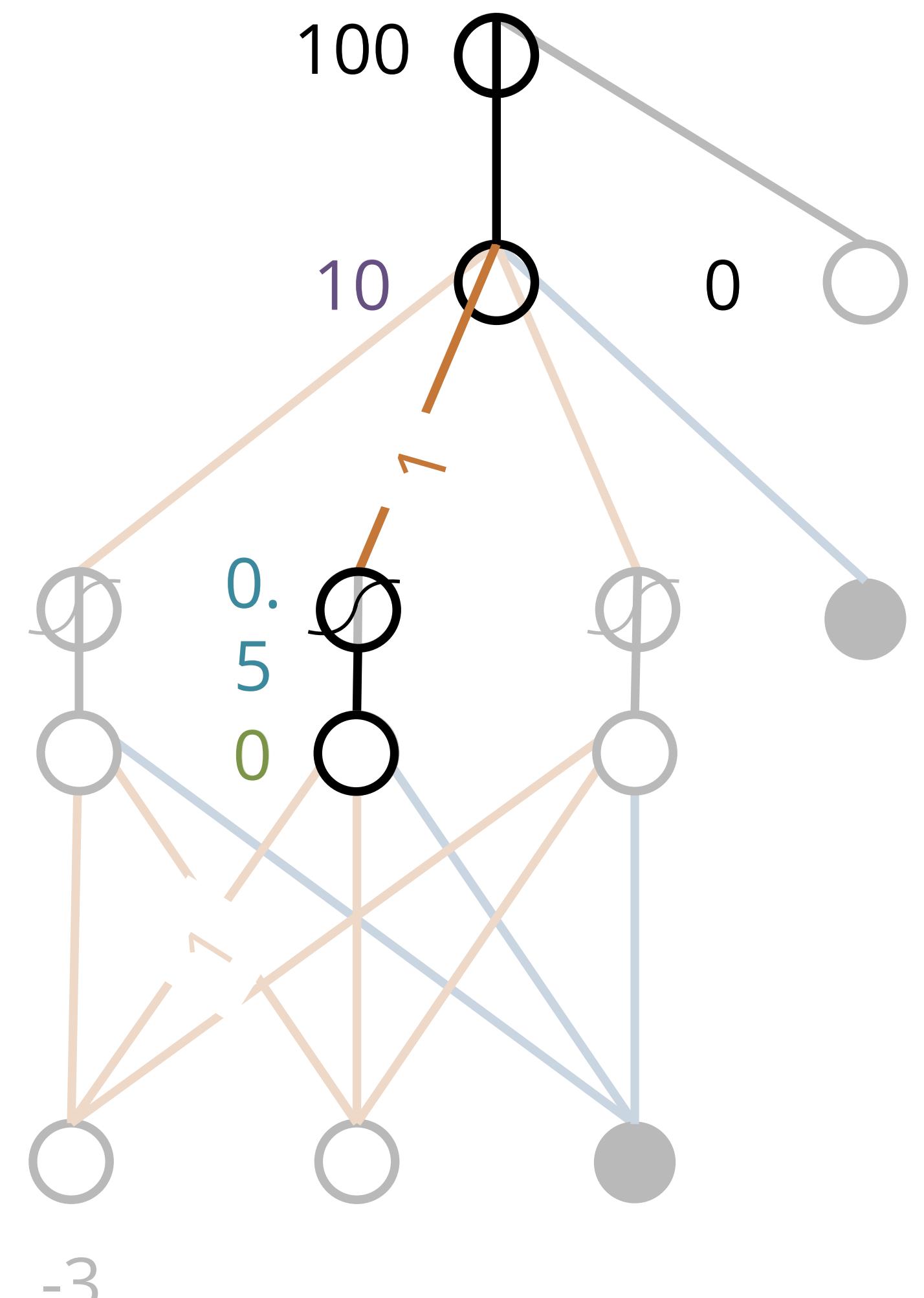
NEGATIVE ACTIVATION



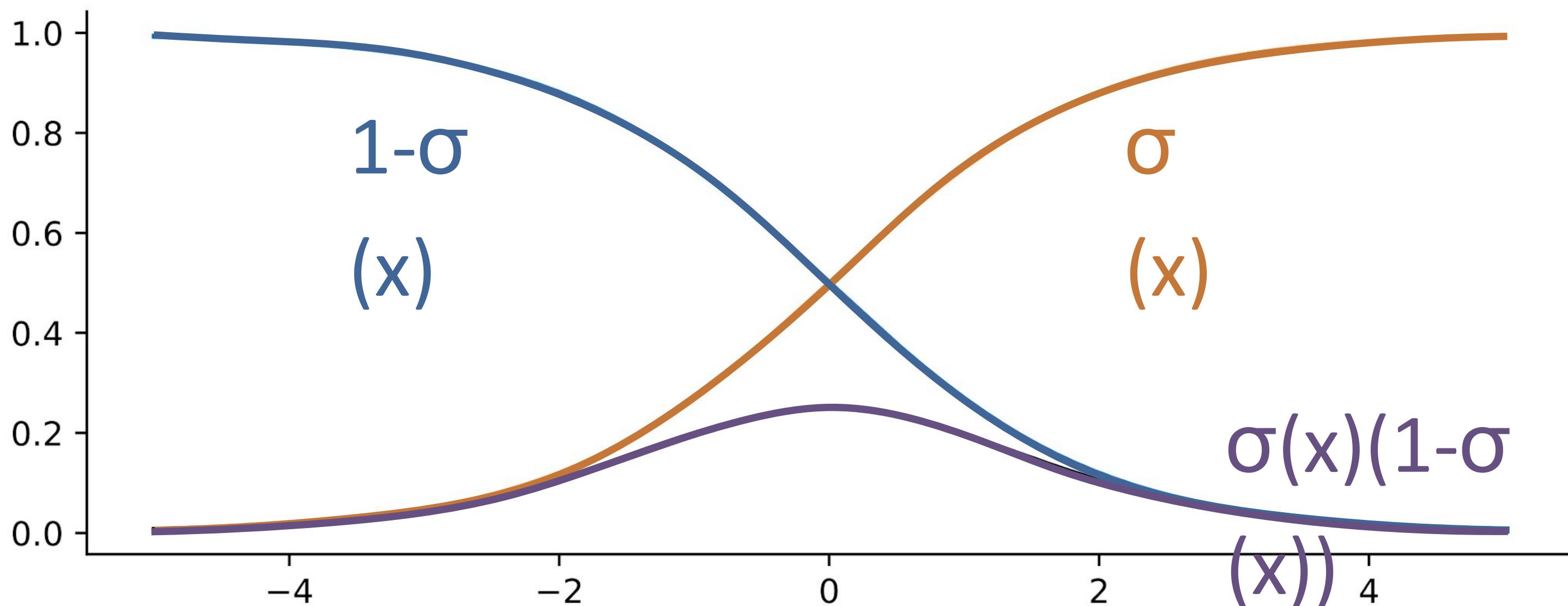
$$\begin{aligned}v_1 &\leftarrow v_1 - \alpha \cdot 20 \cdot 0.1 \\v_2 &\leftarrow v_2 + \alpha \cdot 20 \cdot 0.5 \\v_3 &\leftarrow v_3 - \alpha \cdot 20 \cdot 0.9\end{aligned}$$

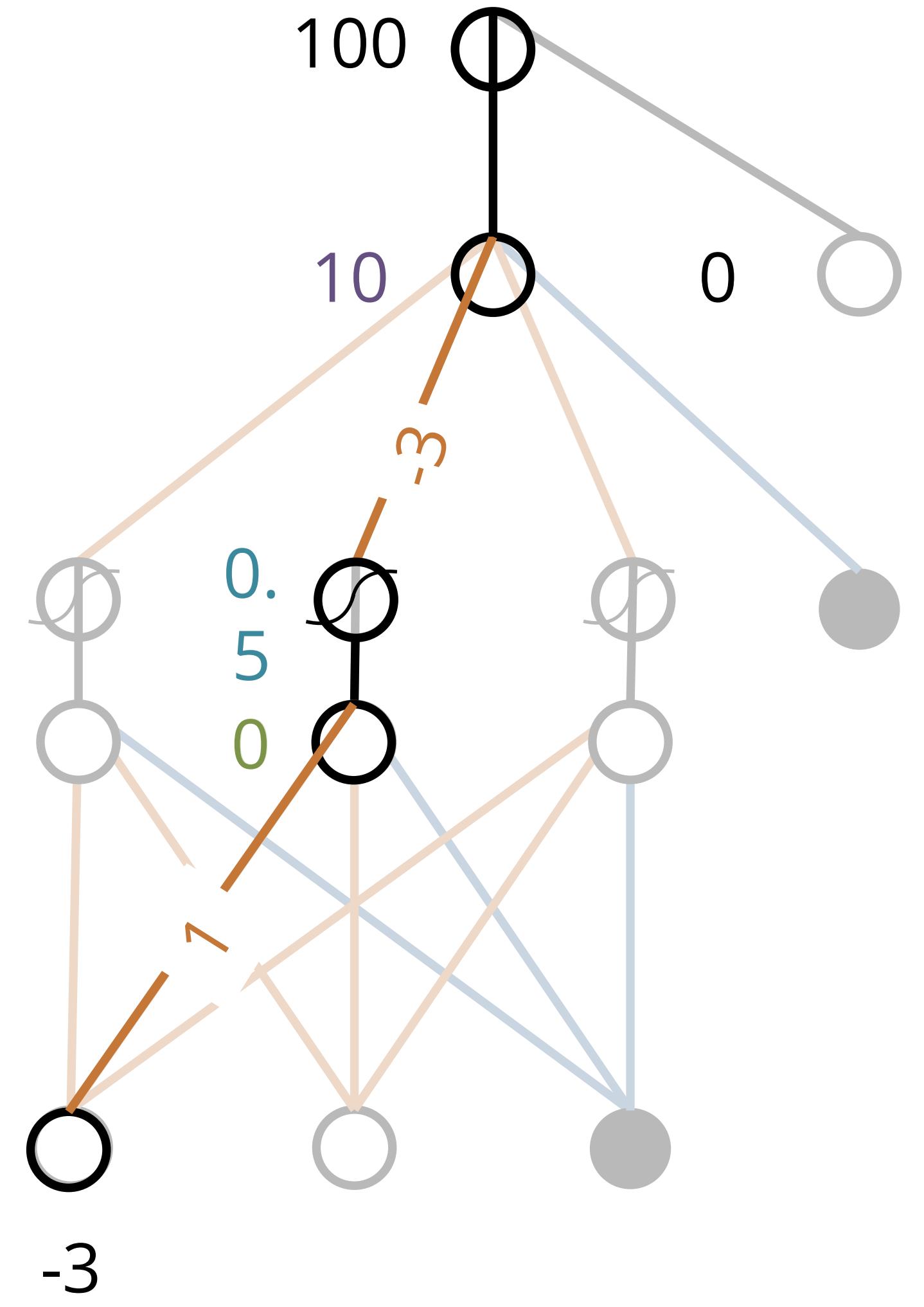


$$\begin{aligned}
 h_2 &\leftarrow h_2 - \alpha \cdot 2(y - t)v_2 \\
 &= h_2 + \alpha \cdot 20 \cdot 3
 \end{aligned}$$



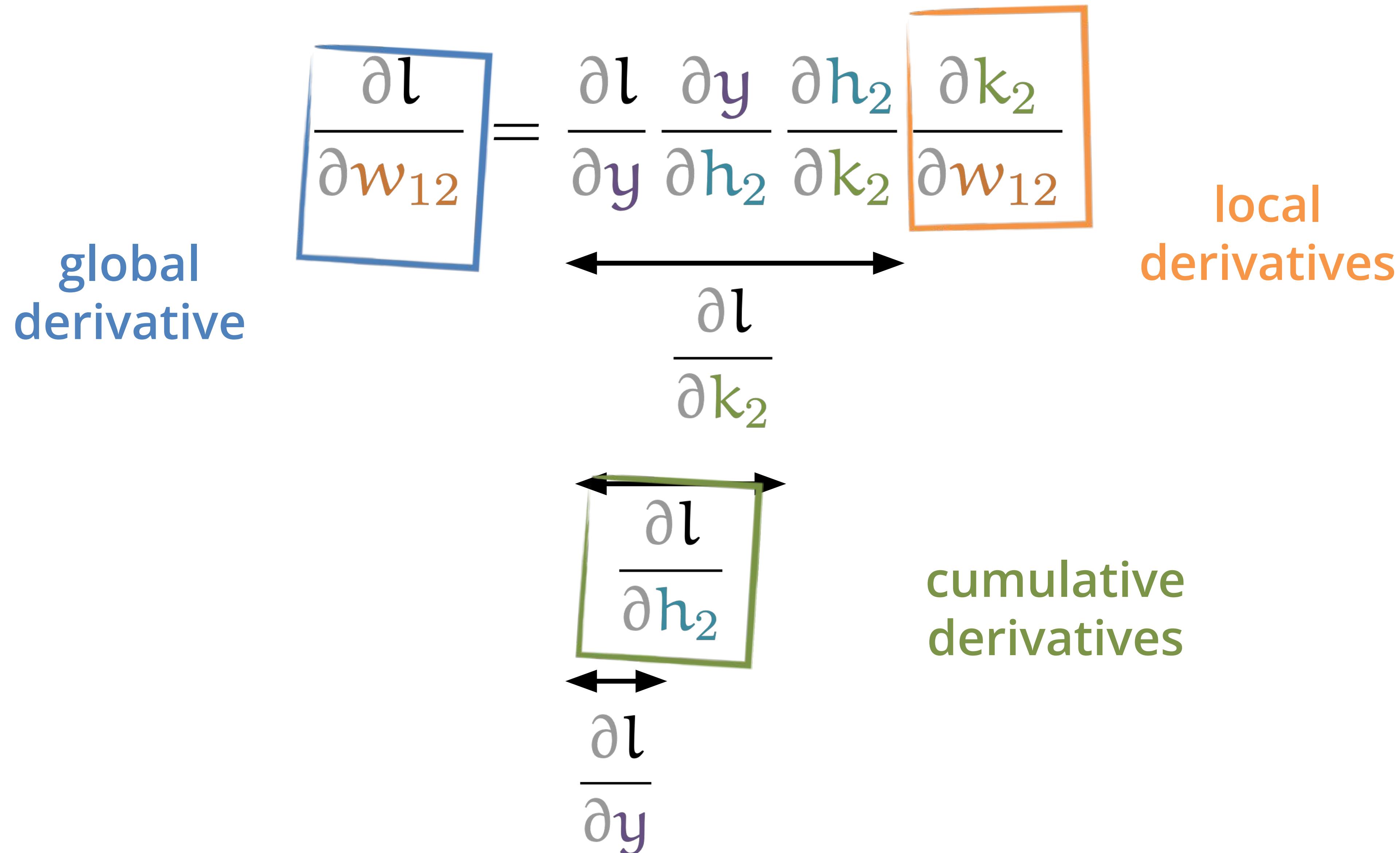
$$k_2 \leftarrow k_2 - \alpha \cdot 2(y - t)v_2 h_2(1 - h_2)$$



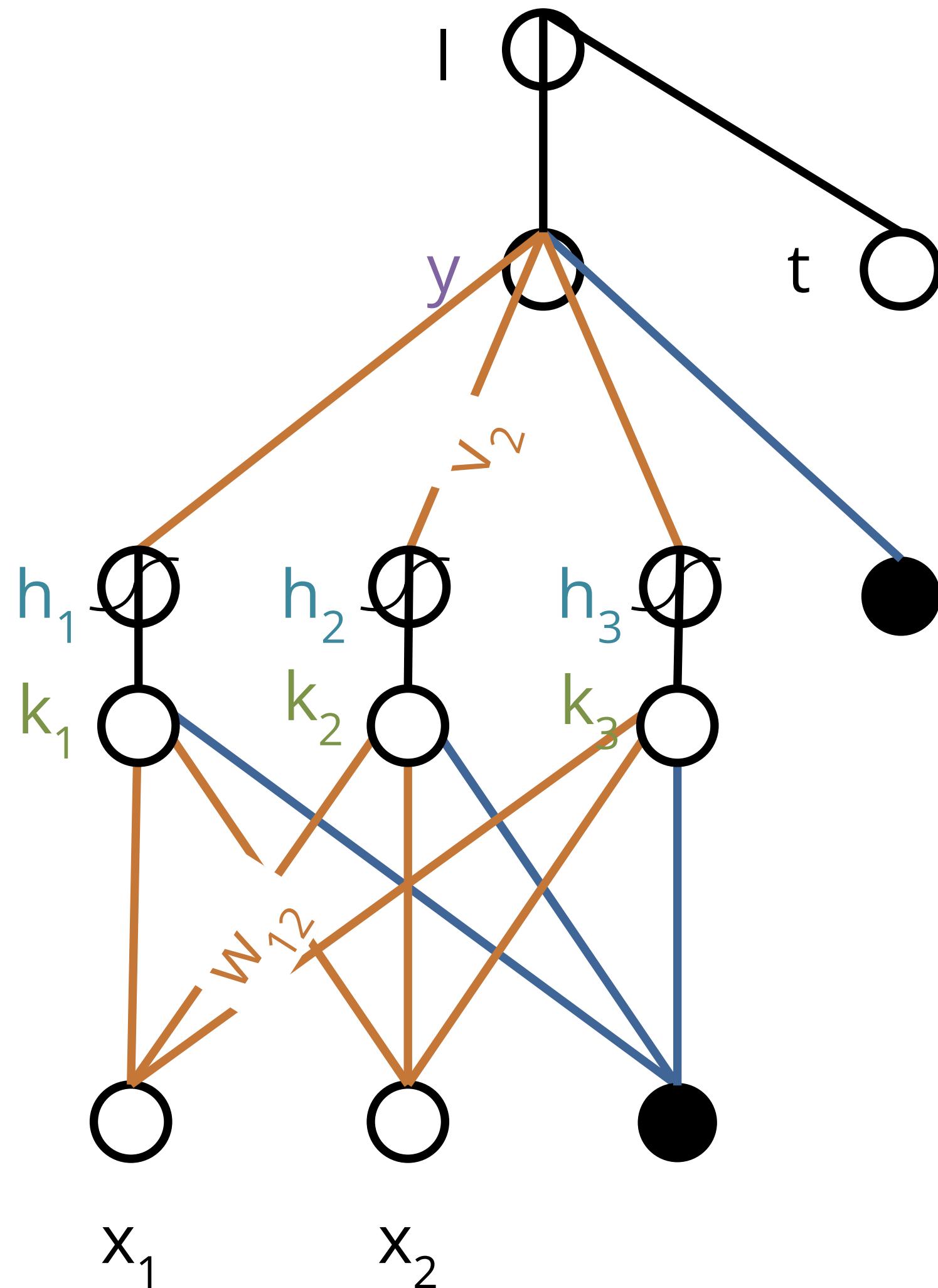


$$\begin{aligned}
 w_{12} &\leftarrow w_{12} - \alpha \cdot 2(y - t) v_2 h_2 (1 - h_2) x_1 \\
 &= w_{12} - \alpha \cdot 20 \cdot -3 \cdot \frac{1}{4} \cdot -3 \\
 &= w_{12} - \alpha \cdot 20 \cdot 3 \cdot \frac{1}{4} \cdot 3
 \end{aligned}$$

GLOBAL AND LOCAL AND ACCUMULATED DERIVATIVES



FORWARD PASS IN PSEUDOCODE



for i **in** $[1 \dots 2]$:

for j **in** $[1 \dots 3]$:

$k[i] += W[i,j] * x[i]$

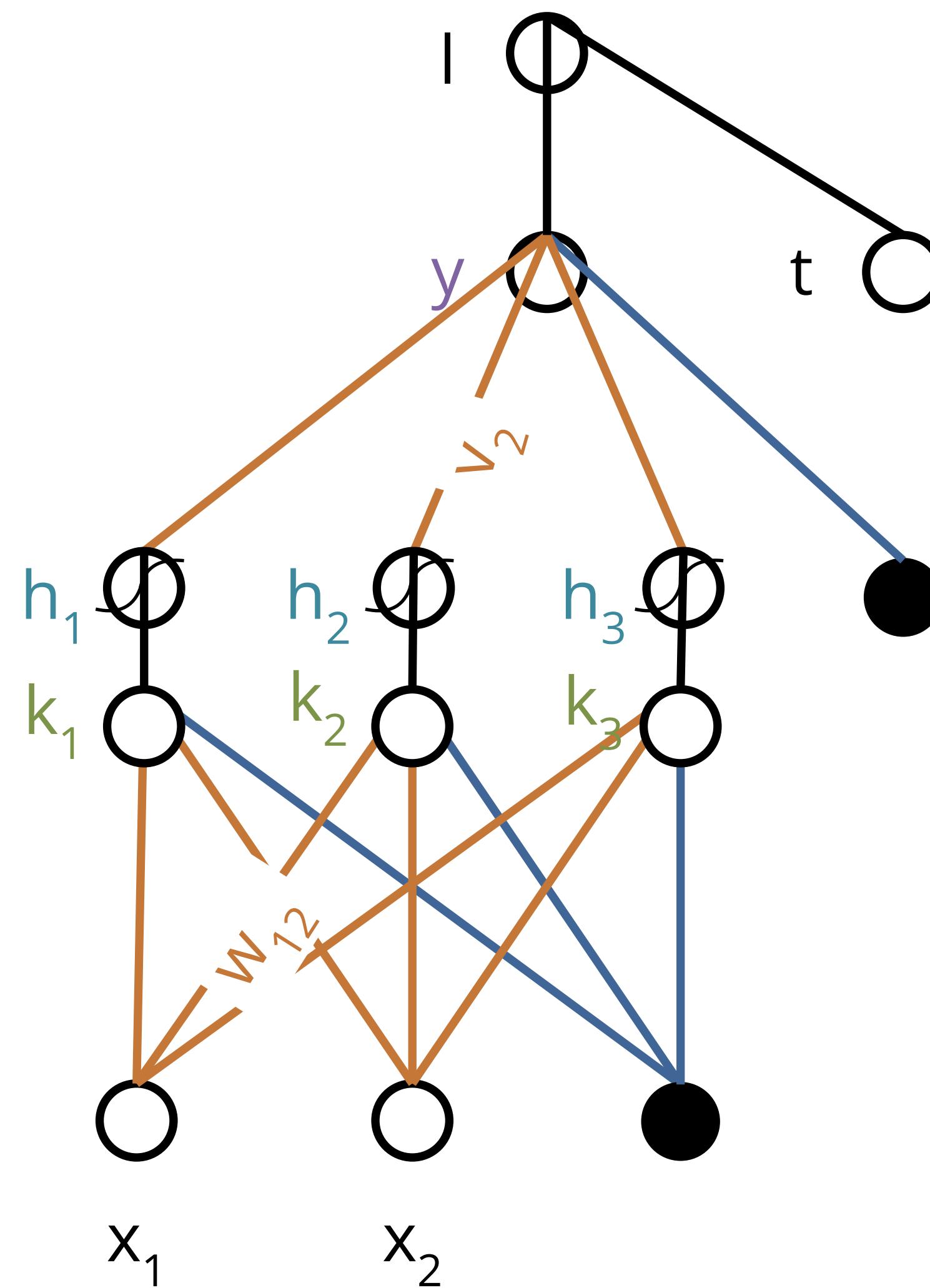
$k[i] += b[i]$

for i **in** $[1 \dots 3]$:

$h[i] = \text{sigmoid}(k[i])$

for i **in** $[1 \dots 3]$:

BACKWARD PASS IN PSEUDOCODE



$y' = 2 * (y - t) \# the error$

for i **in** $[1 \dots 3]$:

$v'[i] = y' * h[i]$

$h'[i] = y' * v[i]$

$c' = y'$

for i **in** $[1 \dots 3]$:

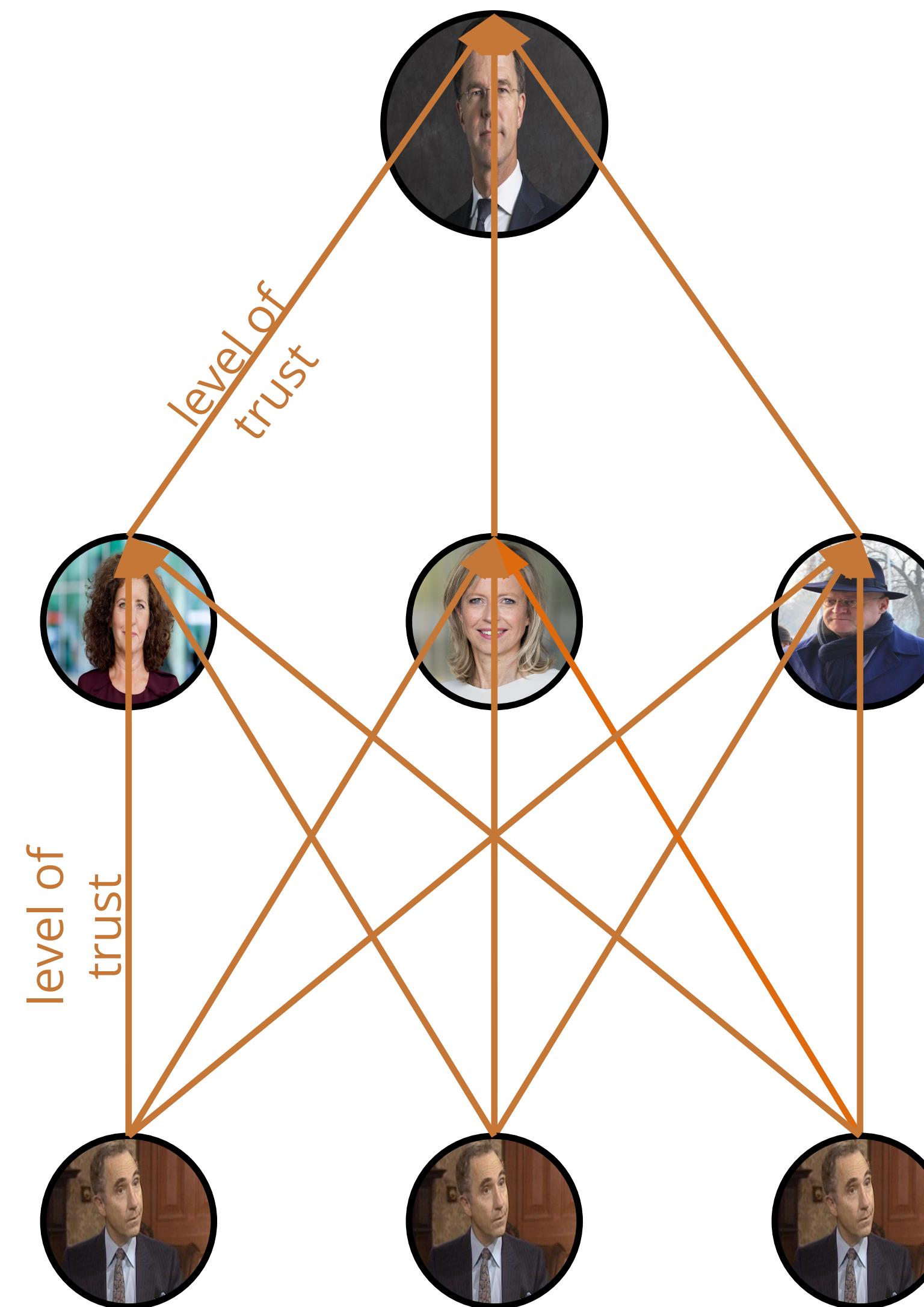
$k'[i] = \text{sig}(h'[i]) * \text{sig}(1-h'[i])$

AN ANALOGY

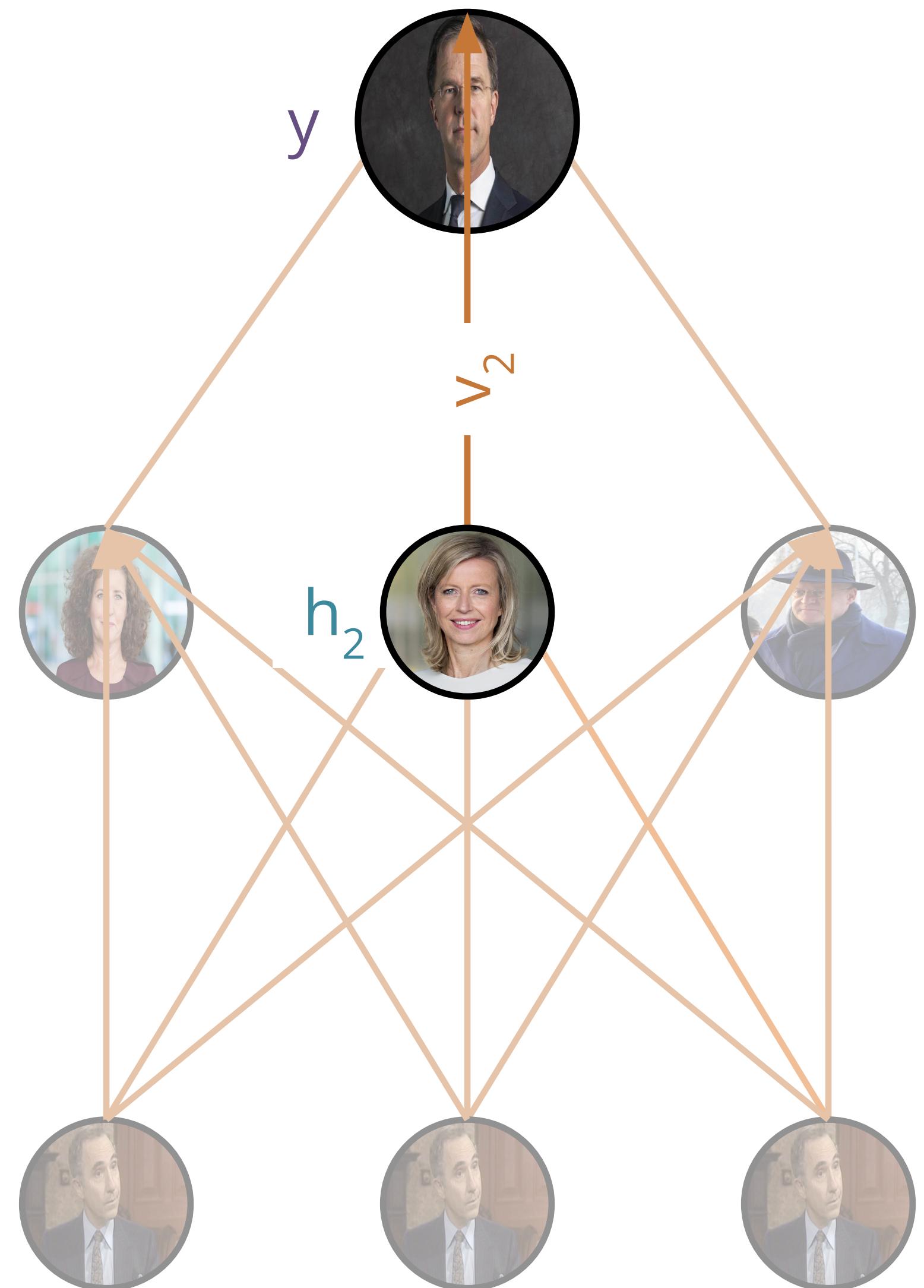
prime
minister

ministe
rs

civil
servants

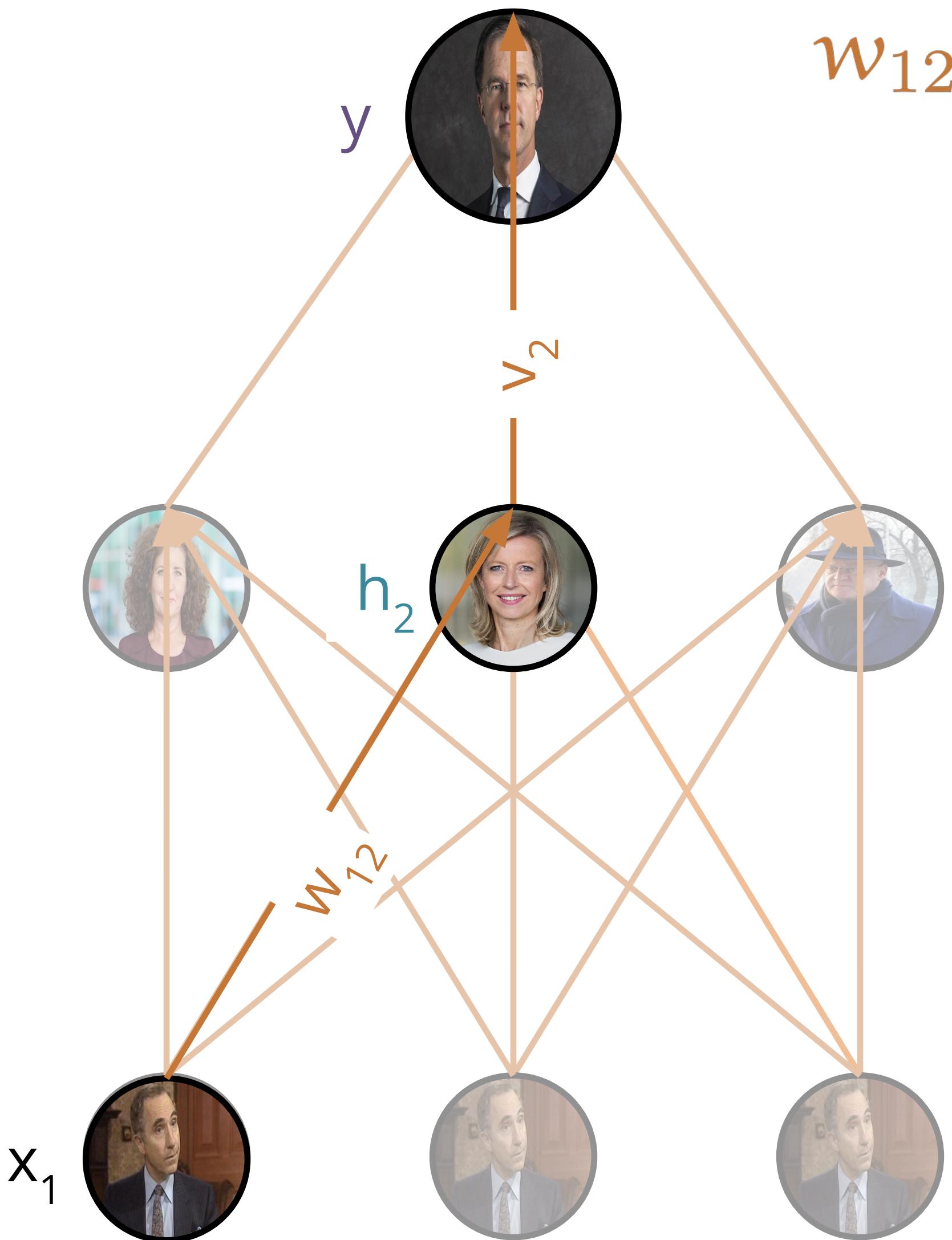


AN ANALOGY



$$v_2 \leftarrow v_2 - \alpha \cdot 2(y - t) \cdot h_2$$

AN ANALOGY



$$w_{12} \leftarrow w_{12} - \alpha \cdot 2(y - t) \cdot v_2 \cdot h_2(1 - h_2) \cdot x_1$$

global error

h_2 's contribution

h_2 's contribution, considering the activation function

x_1 's contribution

Backpropagation: method for computing derivatives.

Combined with gradient descent to train NNs.

Middle ground between symbolic and numeric computation.

- Break computation up into operations.
- Work out local derivatives symbolically.
- Work out global derivative numerically.

Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

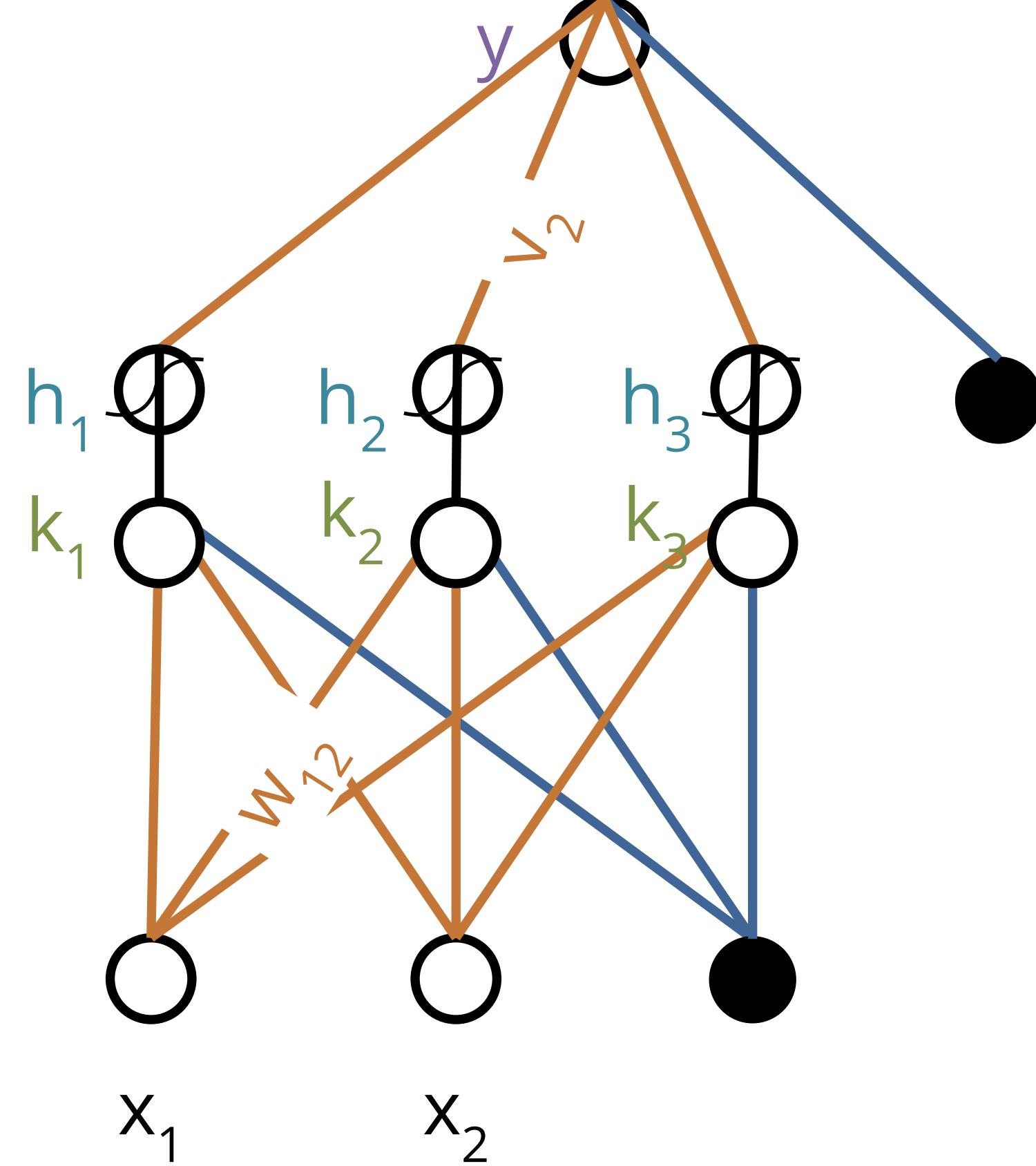
dlvu.github.io



PART THREE: TENSOR BACKPROPAGATION

Expressing backpropagation in vector, matrix and tensor operations.

IT'S ALL JUST LINEAR ALGEBRA

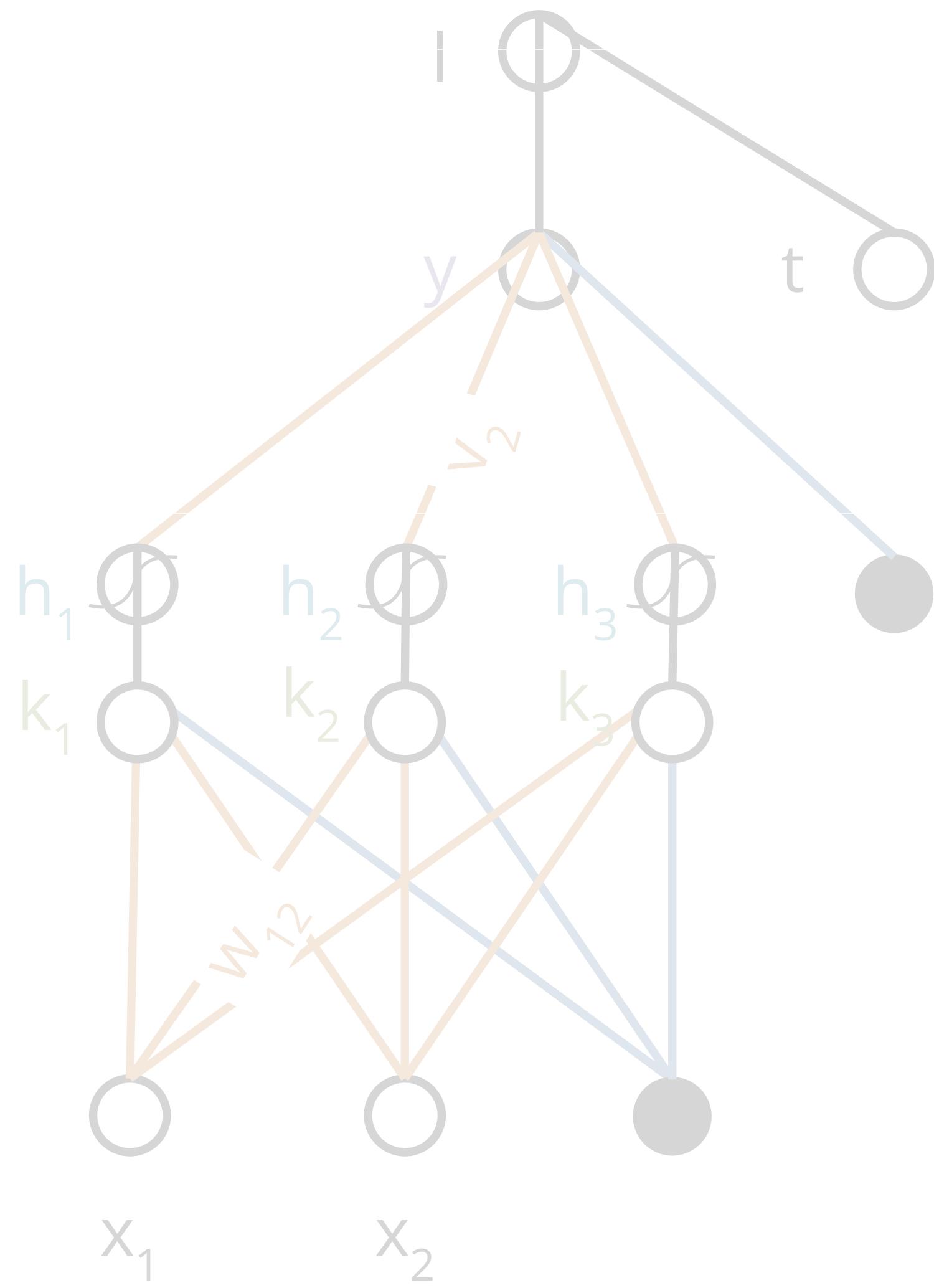


$$\begin{array}{c}
 \text{Matrix Form:} \\
 \begin{matrix}
 & \begin{matrix} x_1 \\ x_2 \end{matrix} \\
 \begin{matrix} w_{11} & w_{21} \\ w_{12} & w_{22} \\ w_{13} & w_{23} \end{matrix} & + & \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \end{matrix} \\
 & \xrightarrow{\sigma} \begin{matrix} h_1 \\ h_2 \\ h_3 \end{matrix}
 \end{matrix}
 \end{array}$$

Diagram illustrating the linear algebra representation of the neural network layer. The inputs x_1 and x_2 are represented as a column vector $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$. The weights w are represented as a matrix $\begin{pmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \\ w_{13} & w_{23} \end{pmatrix}$. The bias b is represented as a column vector $\begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$. The hidden states k are represented as a column vector $\begin{pmatrix} k_1 \\ k_2 \\ k_3 \end{pmatrix}$. The final output h is represented as a column vector $\begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix}$. The sigmoid function σ is applied to the sum of the weighted inputs and bias to produce the final output y .

$$f(\mathbf{x}) = \textcolor{brown}{V} \, \sigma(\textcolor{brown}{W}\mathbf{x} + \mathbf{b}) + \mathbf{c}$$

MATRIX MULTIPLICATION



k, w, k, h, etc are zero'd arrays

for i **in** [1 ... 3]:

for j **in** [1 ... 3]:

$k[i] += W[i,j] * x[i]$

$k[i] += b[i]$

for i **in** [1 ... 3]:

$h[i] = \text{sigmoid}(k[i])$

for i **in** [1 ... 3]:

VECTORIZING

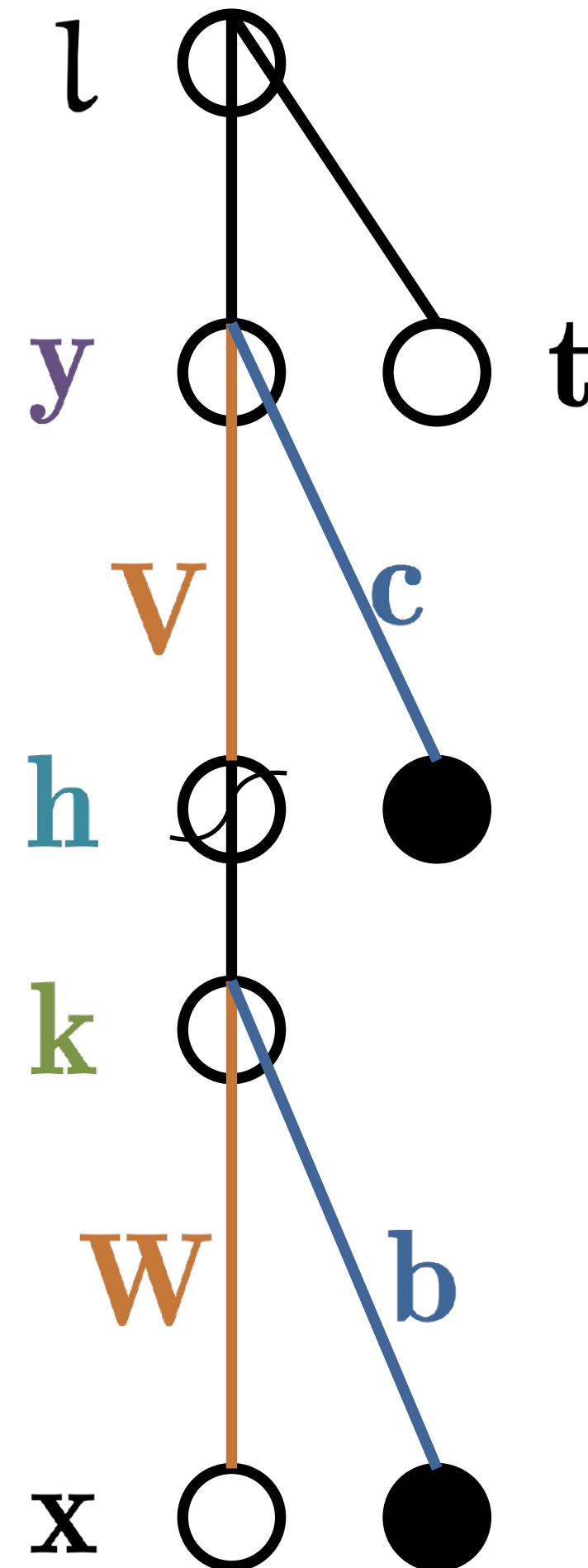
Express everything as operations on vectors, matrices and tensors.

Get rid of all the loops

Makes the notation *simpler*.

Makes the execution *faster*.

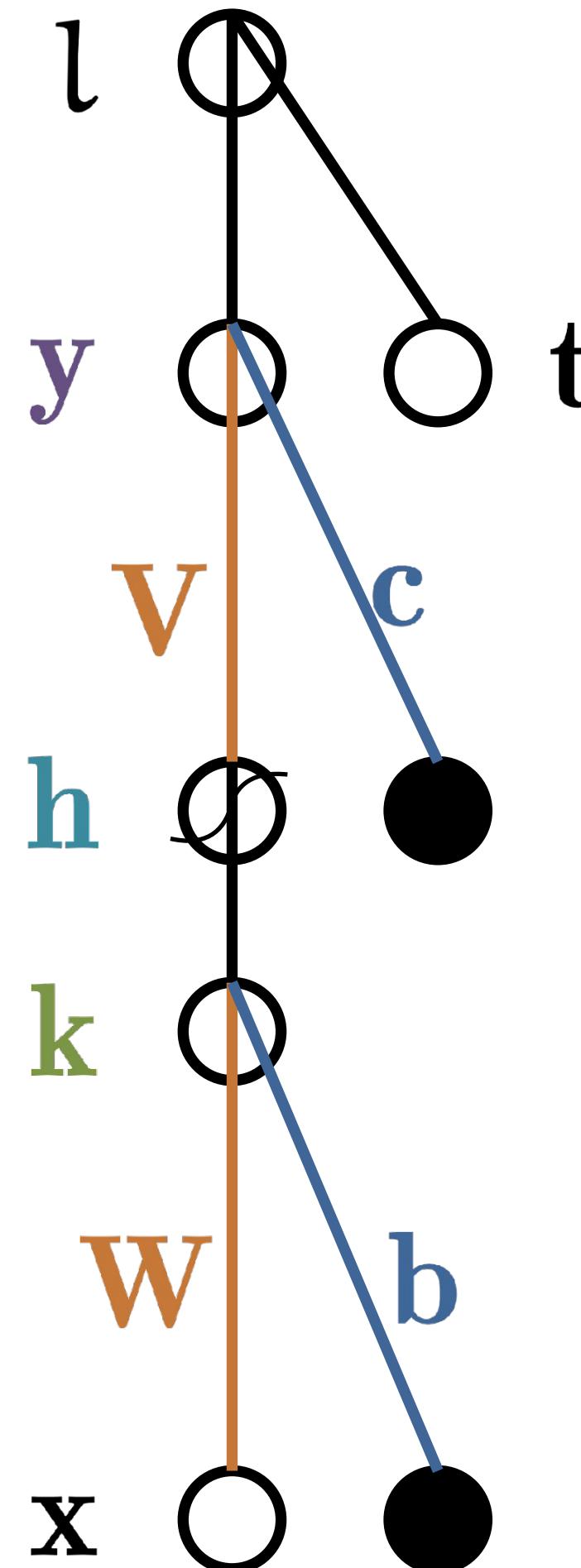
FORWARD



$$l = \sum (y - t)^2$$
$$y = Vh + c$$
$$h = \sigma(k)$$
$$k = Wx + b$$

$$k = \text{W} \cdot \text{dot}(x) + b$$
$$h = \text{sigmoid}(k)$$
$$y = v \cdot \text{dot}(h) + c$$
$$l = (y - t)^2$$

BUT WHAT ABOUT THE BACKWARD PASS?



$$l = \sum (y - t)^2$$

$$y = Vh + c$$

$$h = \sigma(k)$$

$$k = Wx + b$$

Can we do something like this?

$$\frac{\partial l}{\partial W} = ? \quad \frac{\partial l}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial k} \frac{\partial k}{\partial W}$$

GRADIENTS, JACOBIANS, ETC

$f(\mathbf{A}) = \mathbf{B}$, $\frac{\partial \mathbf{B}}{\partial \mathbf{A}}$: derivatives of every element of \mathbf{A} over every element of \mathbf{B}

for instance:

$$f \left(\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \right) = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$\mathbf{J}_f = \begin{pmatrix} \frac{\partial b_1}{\partial a_1} & \frac{\partial b_2}{\partial a_1} \\ \frac{\partial b_1}{\partial a_2} & \frac{\partial b_2}{\partial a_2} \\ \frac{\partial b_1}{\partial a_3} & \frac{\partial b_2}{\partial a_3} \end{pmatrix}$$

input is

function returns a	scalar	vector	matrix
arrange derivatives in a:	scalar	vector	matrix
scalar	scalar	vector	matrix
vector	vector	matrix	?
matrix	matrix	?	?

$$\frac{\partial l}{\partial \mathbf{W}} = ? \quad \frac{\partial l}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{k}} \frac{\partial \mathbf{k}}{\partial \mathbf{W}}$$

uh
oh

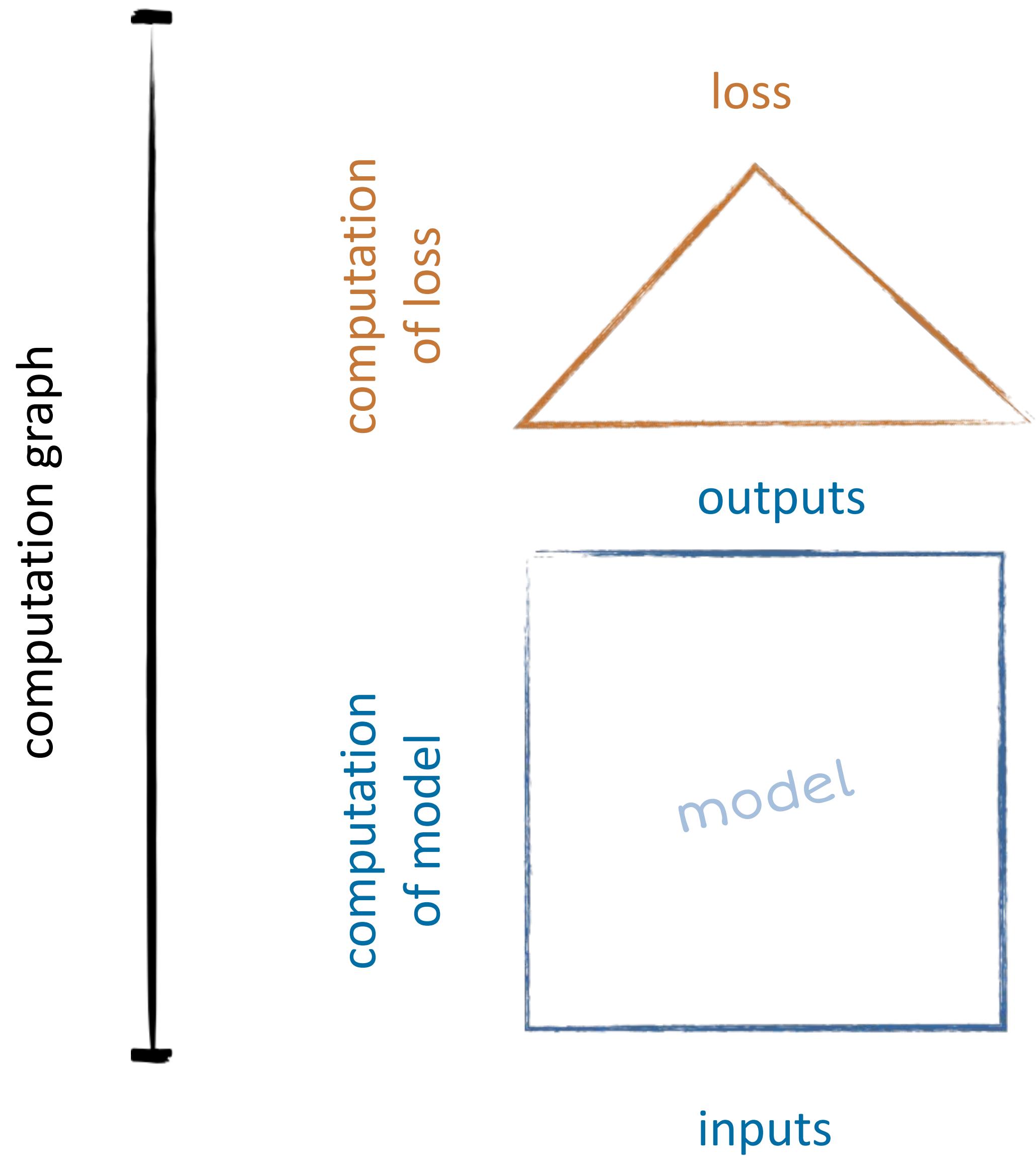
SIMPLIFYING ASSUMPTIONS

- 1) The computation graph *always* has a single, scalar output: l
- 2) We are only ever interested in the derivative of l .

$$\frac{\partial l}{\partial \mathbf{W}}$$

\leftarrow scalar
 \leftarrow tensor of any dimension

	scalar	vector	matrix
scalar	scalar	vector	matrix
vector	vector	matrix	?
matrix	matrix	?	?
3-tensor	3-tensor	?	?



THE GRADIENT

We will call $\frac{\partial l}{\partial \mathbf{W}}$ the gradient of l (with respect to \mathbf{W})

Commonly written as $\nabla_{\mathbf{W}} l$

Nonstandard assumption: $\nabla_{\mathbf{W}} l$ has the same *shape* as \mathbf{W}

$$[\nabla_{\mathbf{W}} l]_{ijk} = \frac{\partial l}{\partial W_{ijk}}$$

NEW NOTATION: THE GRADIENT FOR W

$$\nabla_{\mathbf{W}} l \quad \begin{matrix} \leftarrow \text{always the} \\ \text{same} \\ \leftarrow \text{actual object} \\ \text{of interest} \end{matrix} \quad = \mathbf{W}^\nabla$$

$$\mathbf{b}^\nabla = \nabla_{\mathbf{b}} l = \left[\frac{\partial l}{\partial b_1} \cdots \frac{\partial l}{\partial b_n} \right]^\top$$

$$y^\nabla = \nabla_y l = \frac{\partial l}{\partial y}$$

$$[\mathbf{W}^\nabla]_{ij} = [W_{ij}]^\nabla = \frac{\partial l}{\partial W_{ij}} \\ = w_{ij}^\nabla$$

TENSOR BACKPROPAGATION

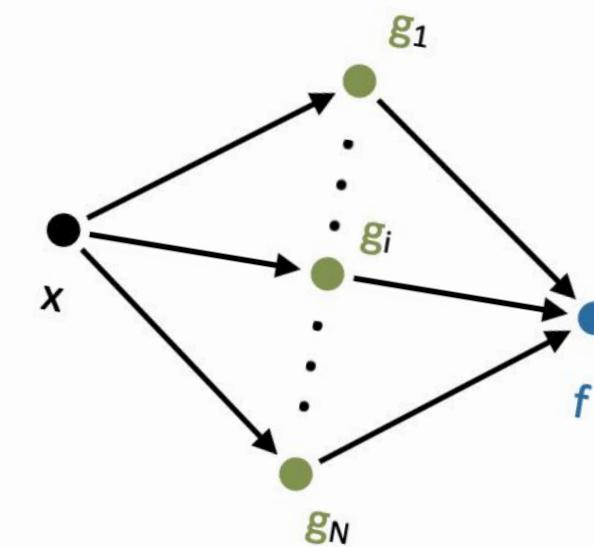
Work out *scalar* derivatives first, then vectorize.

Use the **multivariate chain rule** to derive the scalar derivative.

Apply the chain rule step by step.

Start at the loss and work backward, *accumulating* the gradients.

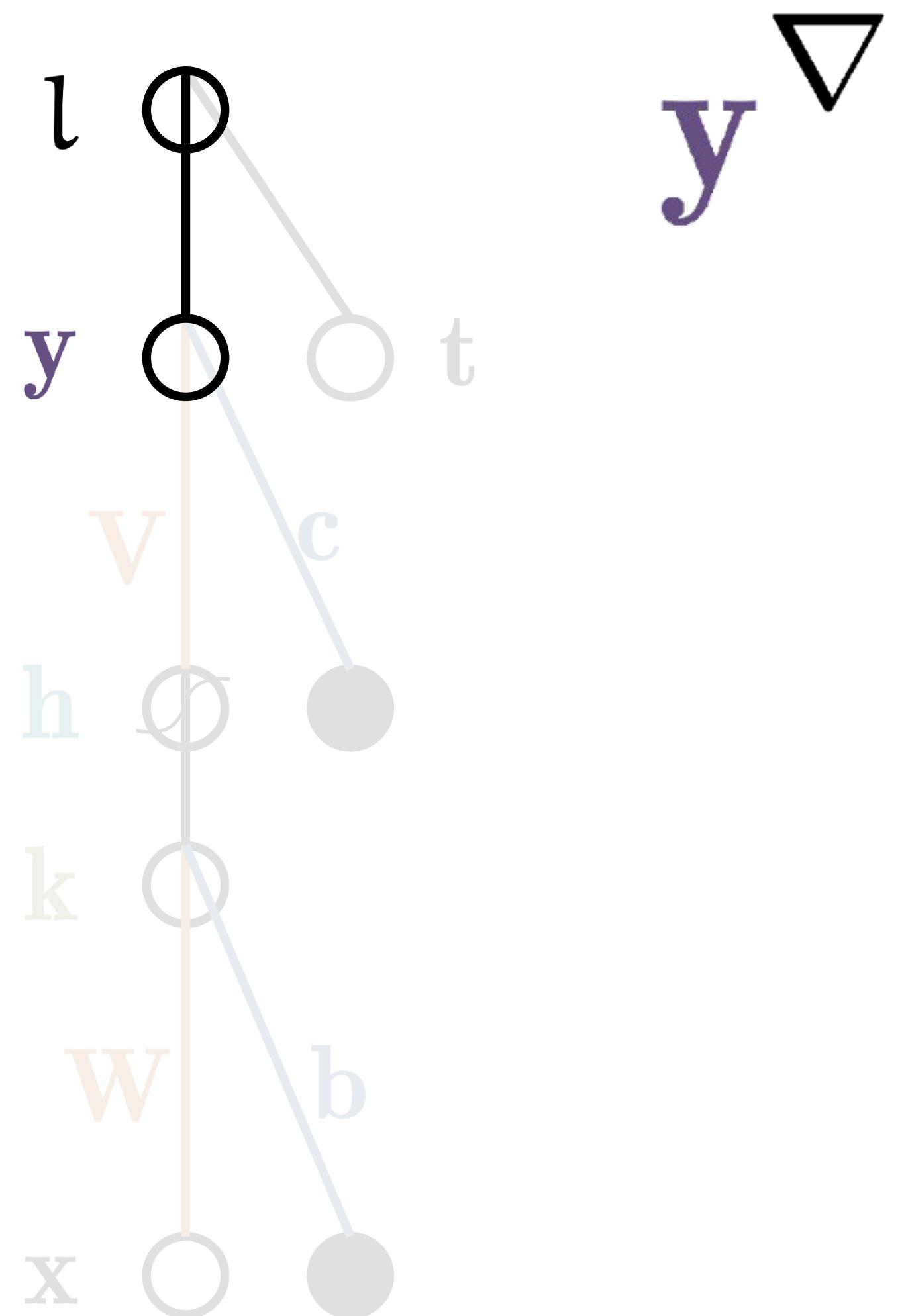
MULTIVARIATE CHAIN RULE



30

$$\frac{\partial \mathbf{f}}{\partial x} = \sum_i \frac{\partial \mathbf{f}}{\partial \mathbf{g}_i} \frac{\partial \mathbf{g}_i}{\partial x}$$

GRADIENT FOR \mathbf{y}

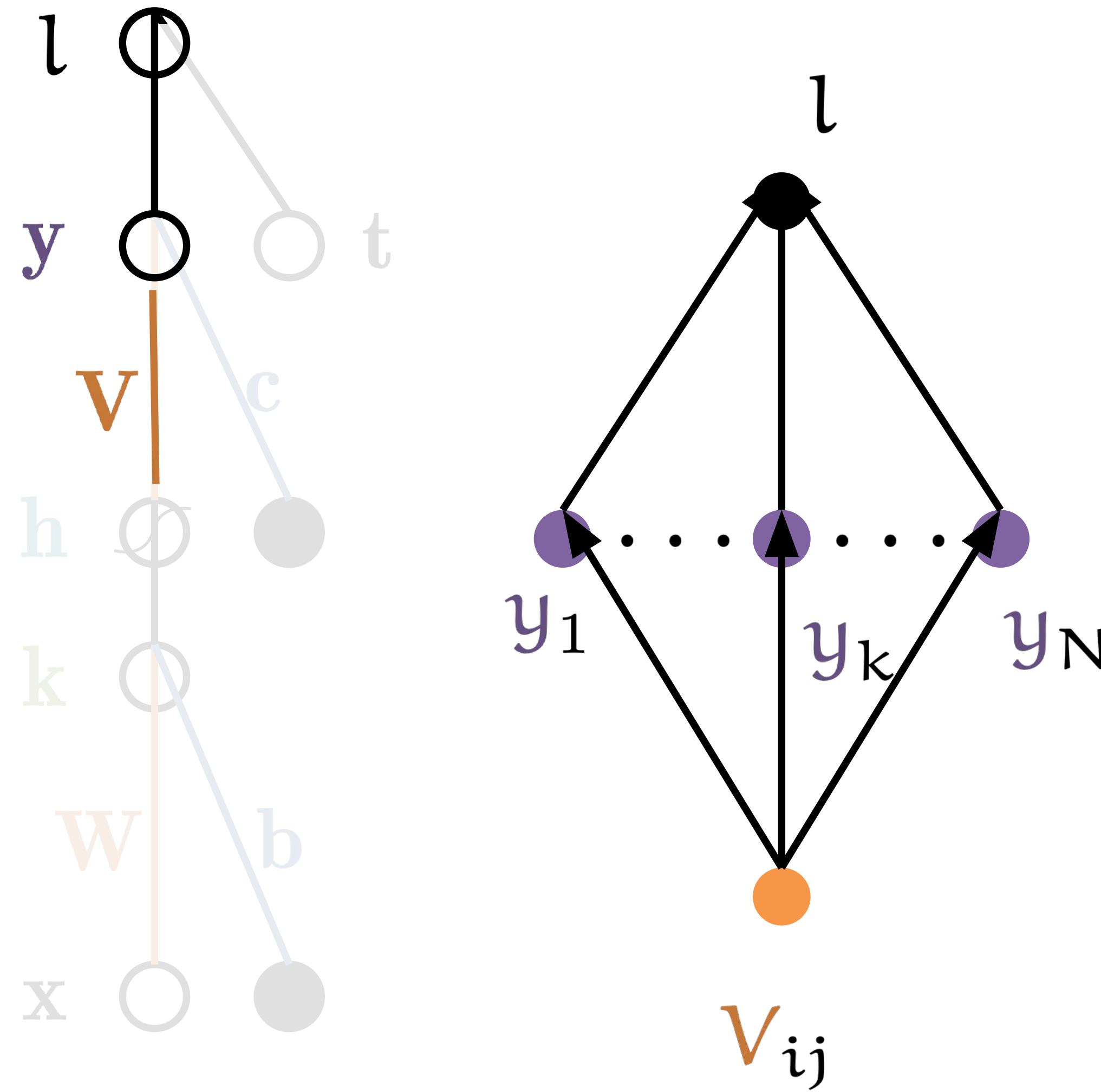


$$\begin{aligned} \mathbf{y}^\nabla &= \frac{\partial l}{\partial \mathbf{y}_i} \\ &= \frac{\partial \sum_k (\mathbf{y}_k - \mathbf{t}_k)^2}{\partial \mathbf{y}_i} = \sum_k \frac{\partial (\mathbf{y}_k - \mathbf{t}_k)^2}{\partial \mathbf{y}_i} \end{aligned}$$

$$= \sum_k 2(\mathbf{y}_k - \mathbf{t}_k) \frac{\partial \mathbf{y}_k - \mathbf{t}_k}{\partial \mathbf{y}_i} = 2(\mathbf{y}_i - \mathbf{t}_i)$$

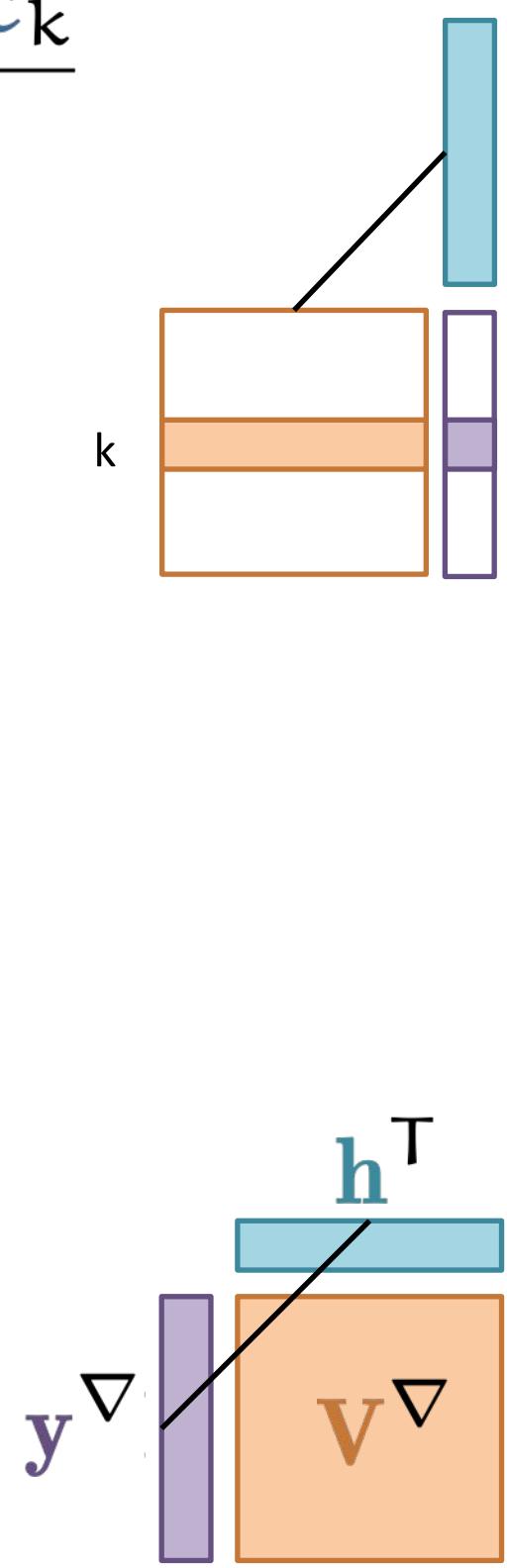
$$\mathbf{y}^\nabla = 2 \begin{pmatrix} \mathbf{y}_1 - \mathbf{t}_1 \\ \vdots \\ \mathbf{y}_n - \mathbf{t}_n \end{pmatrix} = 2 \cdot (\mathbf{y} - \mathbf{t})$$

GRADIENT FOR \mathbf{V}



$$\begin{aligned}
 \mathbf{V}_{ij}^{\nabla} &= \frac{\partial l}{\partial \mathbf{V}_{ij}} \\
 &= \sum_k \frac{\partial l}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{V}_{ij}} = \sum_k y_k^{\nabla} \frac{\partial y_k}{\partial \mathbf{V}_{ij}} \\
 &= \sum_k y_k^{\nabla} \frac{\partial [\mathbf{V}\mathbf{h} + \mathbf{c}]_k}{\partial \mathbf{V}_{ij}} = \sum_k y_k^{\nabla} \frac{\partial [\mathbf{V}_{k:\cdot} \mathbf{h}]_k + c_k}{\partial \mathbf{V}_{ij}} \\
 &= \sum_k y_k^{\nabla} \frac{\partial \sum_l \mathbf{V}_{kl} h_l + c_k}{\partial \mathbf{V}_{ij}} \\
 &= \sum_{kl} y_k^{\nabla} \frac{\partial \mathbf{V}_{kl} h_l}{\partial \mathbf{V}_{ij}} = y_i^{\nabla} \frac{\partial \mathbf{V}_{ij} h_j}{\partial \mathbf{V}_{ij}} \\
 &= y_i^{\nabla} h_j
 \end{aligned}$$

$$\mathbf{V}^{\nabla} = \begin{pmatrix} \dots & y_i^{\nabla} h_j & \dots \\ & \vdots & \\ & \vdots & \end{pmatrix} = \mathbf{y}^{\nabla} \mathbf{h}^{\top}$$

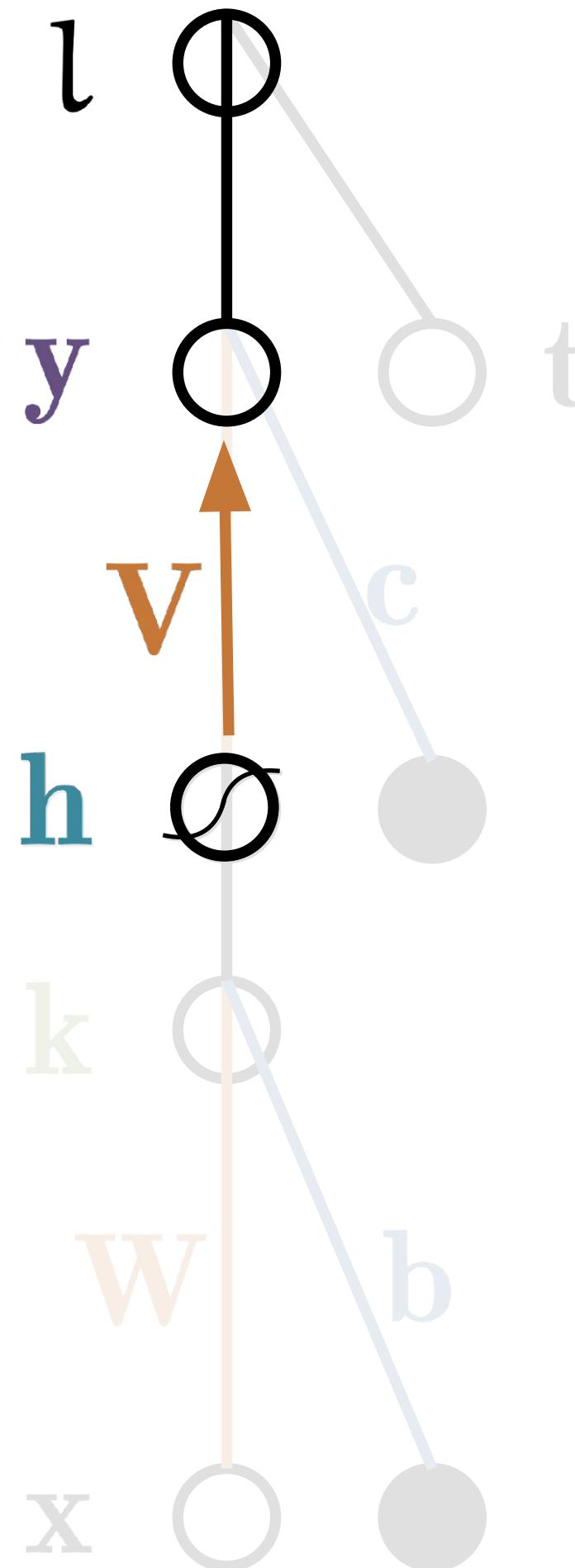


RECIPE

To work out a gradient \mathbf{X}^∇ for some \mathbf{X} :

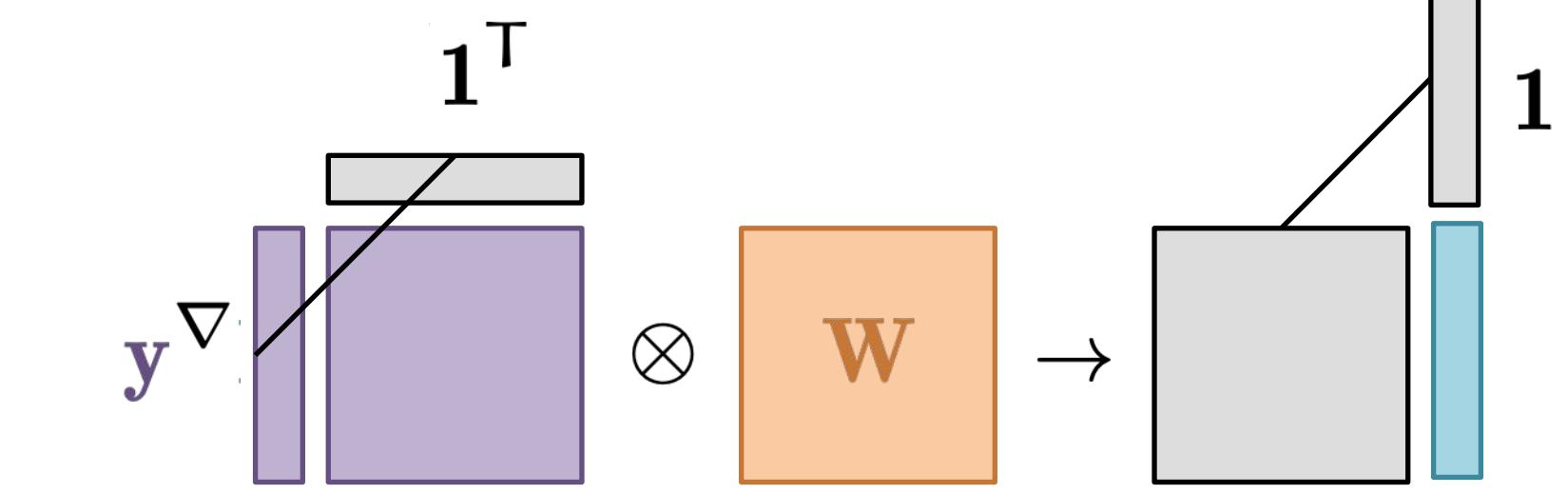
- Write down a **scalar derivative** of the loss over *one element* of \mathbf{X} .
 - Use the **multivariate chain rule** to sum over all outputs.
 - Work out the scalar derivative.
- ❖ **Vectorize** the computation of \mathbf{X}^∇ in terms of the original inputs.

GRADIENT FOR \mathbf{h}



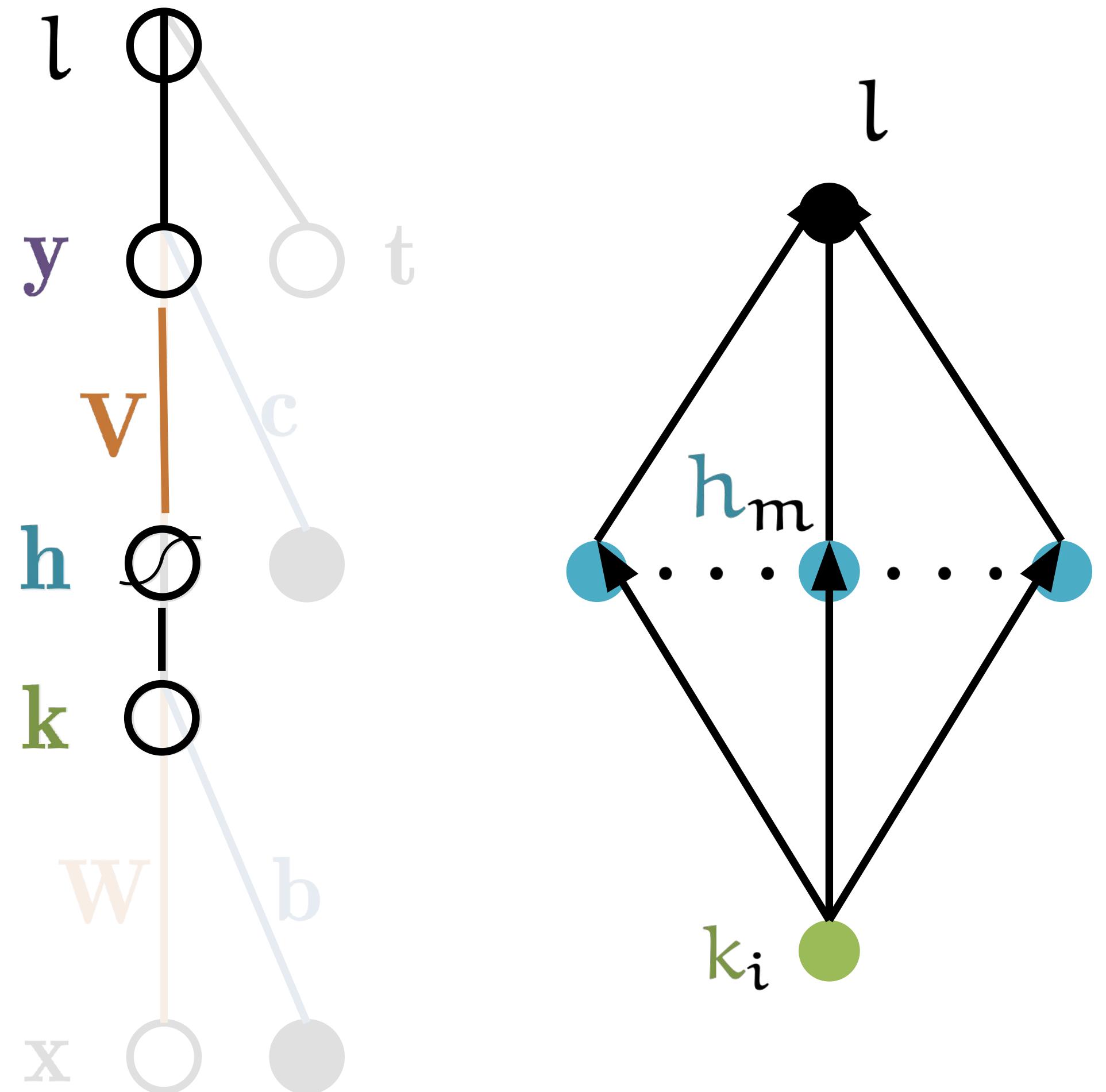
$$\begin{aligned}
 \mathbf{h}_i^\nabla &= \frac{\partial l}{\partial \mathbf{h}_i} \\
 &= \sum_k \frac{\partial l}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{h}_i} = \sum_k y_k^\nabla \frac{\partial y_k}{\partial \mathbf{h}_i} \\
 &= \sum_k y_k^\nabla \frac{\partial [\mathbf{W}\mathbf{h} + \mathbf{c}]_k}{\partial \mathbf{h}_i} = \sum_k y_k^\nabla \frac{\partial \sum_l \mathbf{W}_{kl} \mathbf{h}_l + c_k}{\partial \mathbf{h}_i} \\
 &= \sum_{kl} y_k^\nabla \frac{\partial \mathbf{W}_{kl} \mathbf{h}_l + c_k}{\partial \mathbf{h}_i} = \sum_k y_k^\nabla \frac{\partial \mathbf{W}_{ki} \mathbf{h}_i}{\partial \mathbf{h}_i} \\
 &= \sum_k y_k^\nabla \mathbf{W}_{ki}
 \end{aligned}$$

$$\mathbf{h}^\nabla = \begin{pmatrix} \vdots \\ \sum_k y_k^\nabla \mathbf{W}_{ki} \\ \vdots \end{pmatrix} = (\mathbf{y}^\nabla \mathbf{1}^T \otimes \mathbf{W}) \mathbf{1}$$



or in numpy:
 $(\mathbf{g}_y[\text{None}, :] * \mathbf{W}).\text{sum}(\text{axis}=1)$

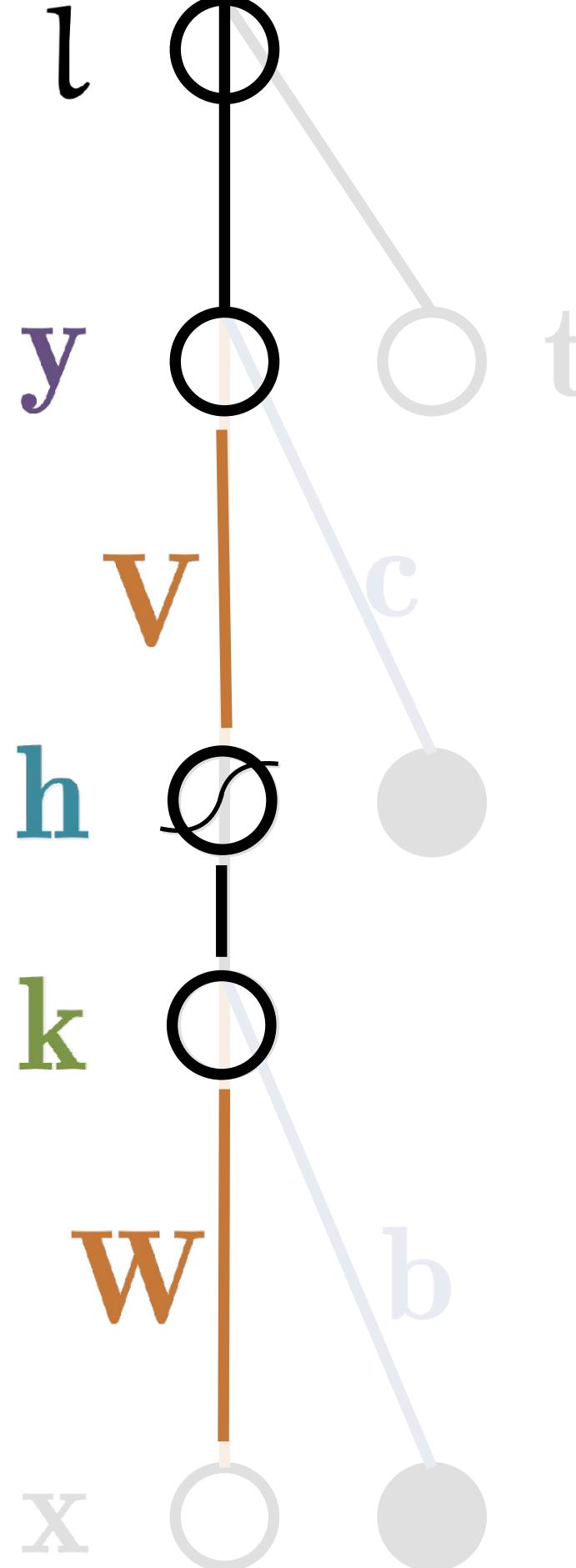
GRADIENT FOR \mathbf{k}



$$\begin{aligned}
 \mathbf{k}_i^\nabla &= \frac{\partial l}{\partial \mathbf{k}_i} \\
 &= \sum_m \mathbf{h}_m^\nabla \frac{\partial \mathbf{h}_m}{\partial \mathbf{k}_i} = \sum_m \mathbf{h}_m^\nabla \frac{\partial \sigma(\mathbf{k}_m)}{\partial \mathbf{k}_i} \\
 &= \mathbf{h}_i^\nabla \frac{\partial \sigma(\mathbf{k}_i)}{\partial \mathbf{k}_i} \\
 &= \mathbf{h}_i^\nabla \sigma'(\mathbf{k}_i) = \mathbf{h}_i^\nabla \sigma(\mathbf{k}_i)(1 - \sigma(\mathbf{k}_i)) \\
 &= \mathbf{h}_i^\nabla \mathbf{h}_i(1 - \mathbf{h}_i)
 \end{aligned}$$

$$\mathbf{k}^\nabla = \mathbf{h}^\nabla \otimes \mathbf{h} \otimes (1 - \mathbf{h})$$

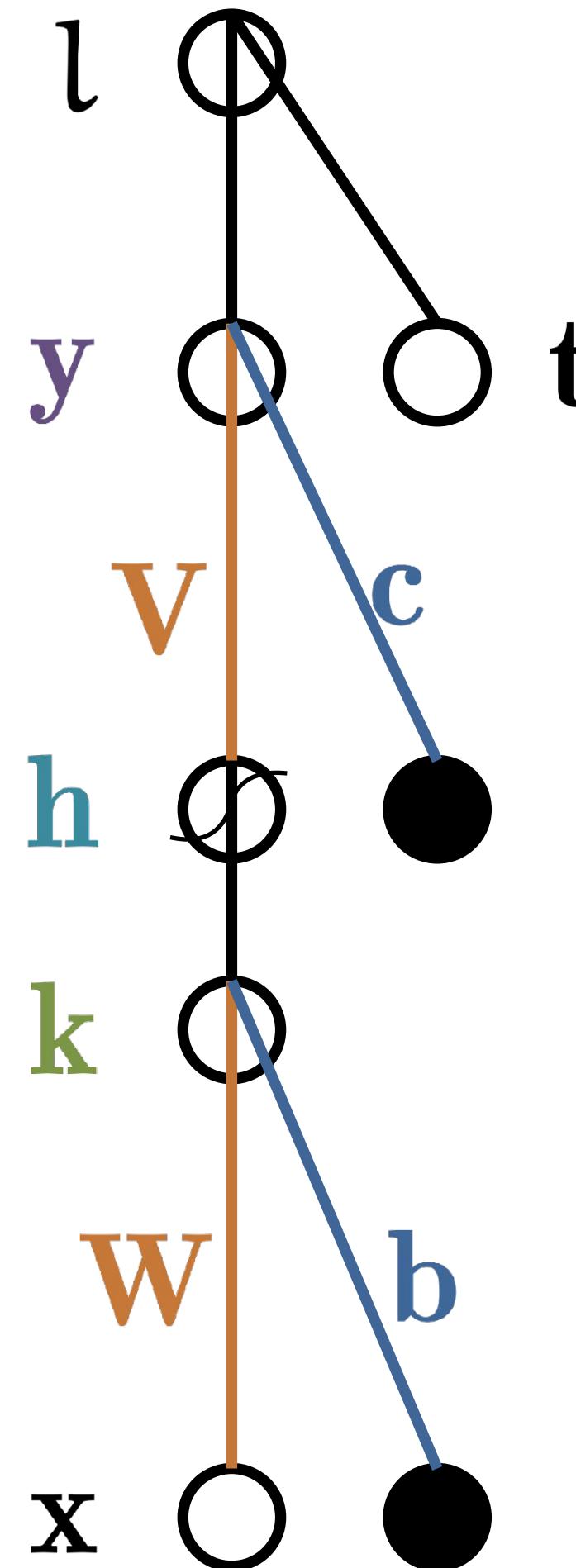
GRADIENT FOR \mathbf{W}



$$\begin{aligned}
 \mathbf{w}_{ij}^{\nabla} &= \frac{\partial l}{\partial \mathbf{w}_{ij}} \\
 &= \sum_m \frac{\partial l}{\partial k_m} \frac{\partial k_m}{\partial \mathbf{w}_{ij}} = \sum_m k_m^{\nabla} \frac{\partial k_m}{\partial \mathbf{w}_{ij}} \\
 &= \sum_m k_m^{\nabla} \frac{\partial [\mathbf{W}\mathbf{x} + \mathbf{b}]_m}{\partial \mathbf{w}_{ij}} = \sum_m k_m^{\nabla} \frac{\partial [\mathbf{W}_{m:\cdot}\mathbf{x}]_m + b_m}{\partial \mathbf{w}_{ij}} \\
 &= \sum_m k_m^{\nabla} \frac{\partial \sum_l \mathbf{w}_{ml} x_l + b_m}{\partial \mathbf{w}_{ij}} \\
 &= \sum_{ml} k_m^{\nabla} \frac{\partial \mathbf{w}_{ml} x_l}{\partial \mathbf{w}_{ij}} = k_i^{\nabla} \frac{\partial \mathbf{w}_{ij} x_j}{\partial \mathbf{w}_{ij}} \\
 &= k_i^{\nabla} x_j
 \end{aligned}$$

$$\mathbf{w}^{\nabla} = \begin{pmatrix} & \vdots & \\ \dots & k_i^{\nabla} x_j & \dots \\ & \vdots & \end{pmatrix} = \mathbf{k}^{\nabla} \mathbf{x}^T$$

PSEUDOCODE: BACKWARD



$$y' = 2 * (\textcolor{violet}{y} - \textcolor{black}{t})$$

$$\textcolor{brown}{v}' = \textcolor{violet}{y}'.\text{mm}(\textcolor{teal}{h}.T)$$

$$\textcolor{teal}{h}' = \textcolor{violet}{y}'[\text{None}, :] * \textcolor{brown}{W}.\text{sum}(\text{axis}=1)$$

$$\textcolor{brown}{k}' = \text{sigmoid}(\textcolor{brown}{k}) * \text{sigmoid}(1 - \textcolor{brown}{k})$$

$$\textcolor{brown}{w}' = \textcolor{brown}{k}'.\text{mm}(\textcolor{black}{x}.T)$$

Vectorizing makes notation simpler, and computation faster.

Vectorizing the **forward pass** is usually easy.

Backward pass: work out the scalar derivatives, accumulate, then vectorize.

Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

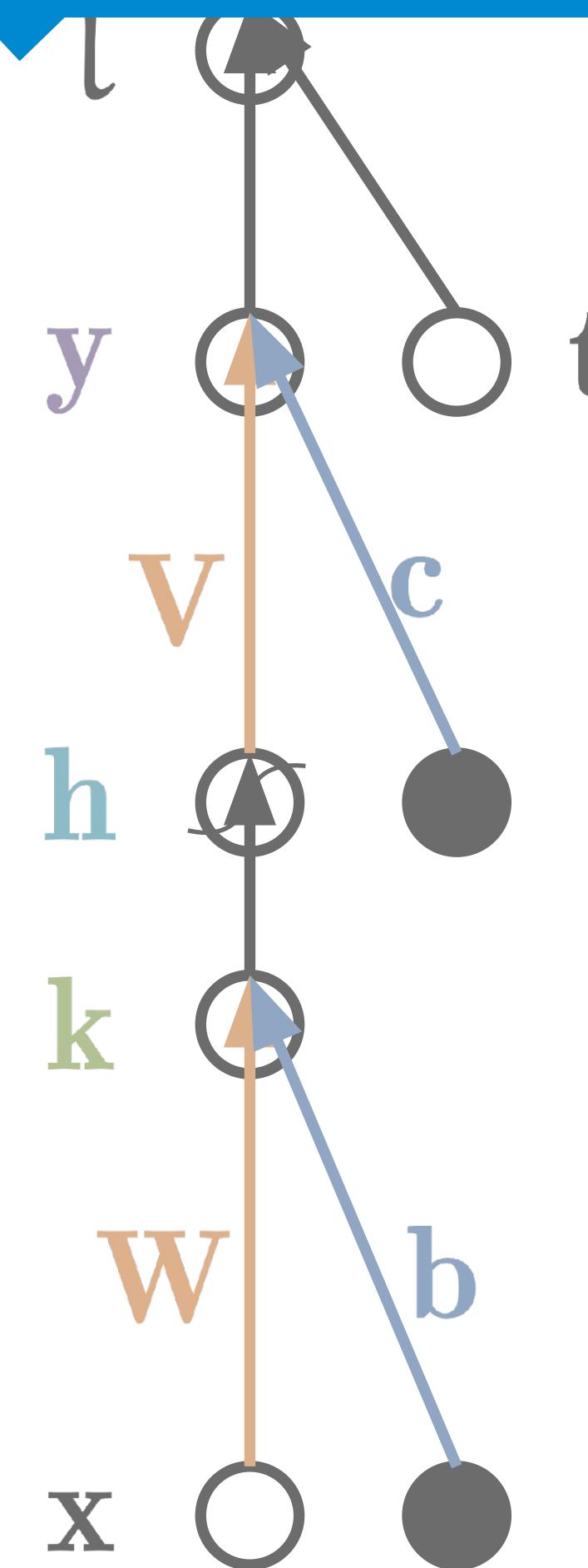
dlvu.github.io



PART FOUR: AUTOMATIC DIFFERENTIATION

Letting the computer do all the work.

THE STORY SO FAR



pen and paper

$$v' = y'.mm(h.T)$$

$$h' = y'[None, :] * W.sum(axis=1)$$

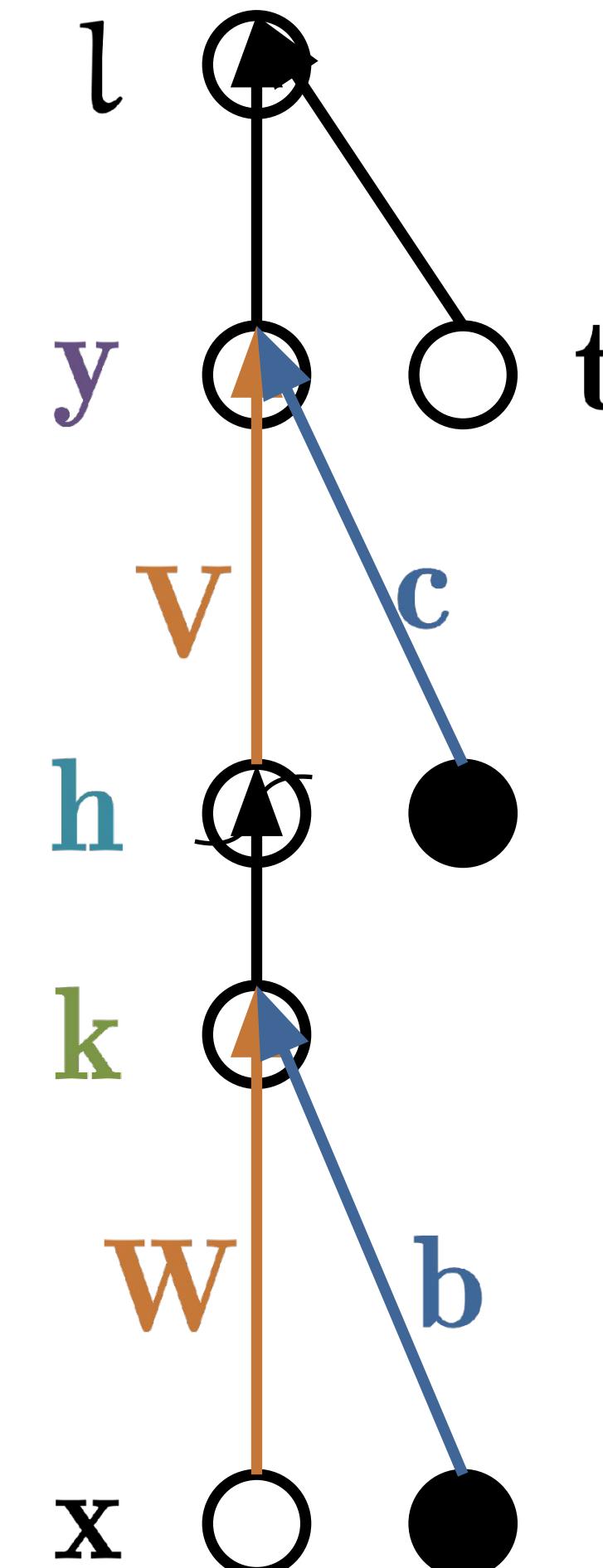
$$k' = \text{sigmoid}(k) * \text{sigmoid}(1 - k)$$

$$w' = k'.mm(x.T)$$

in the computer

THE IDEAL

○
○ +
○ x



$$k = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

$$\mathbf{h} = \text{sigmoid}(\mathbf{k})$$

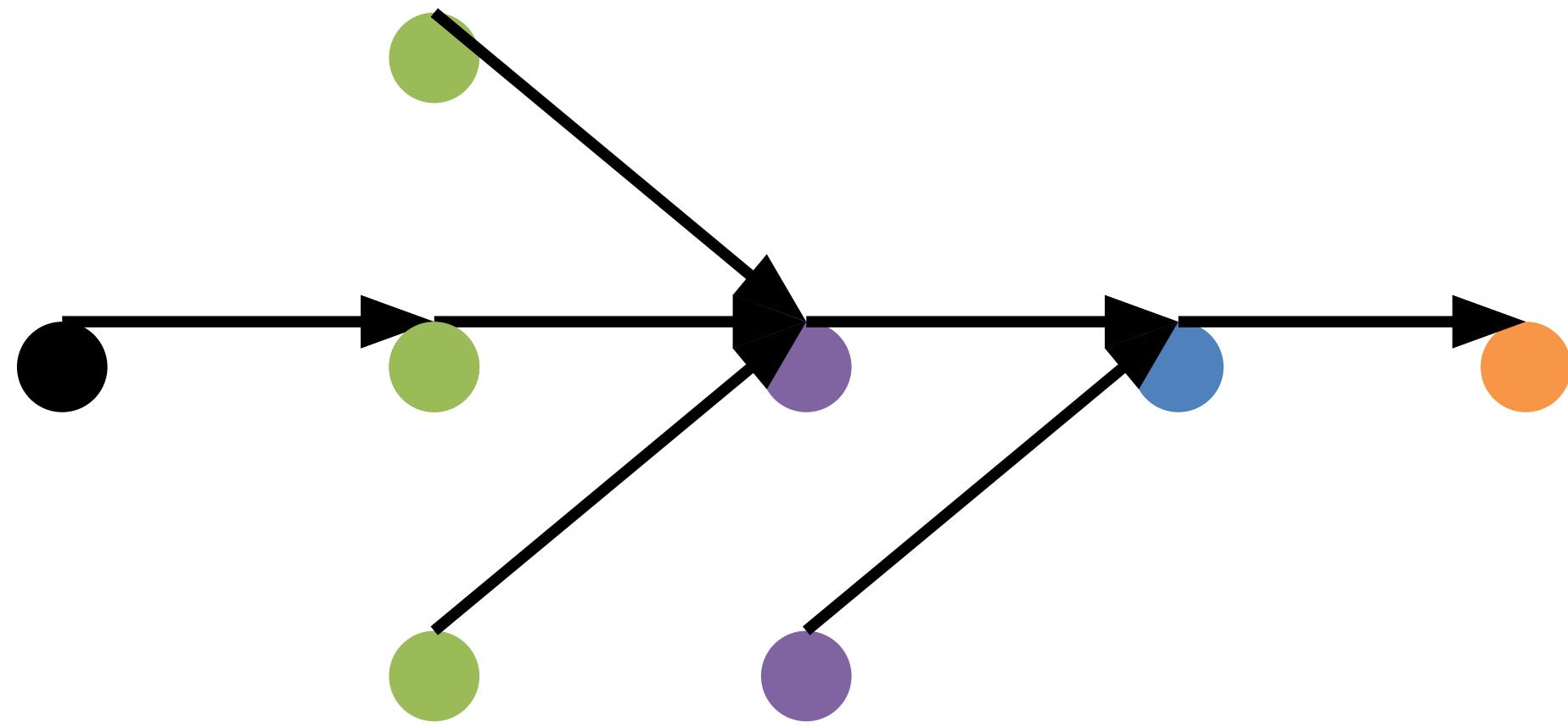
$$y = \mathbf{v} \cdot \mathbf{h} + c$$

$$l = (y - t)^2$$

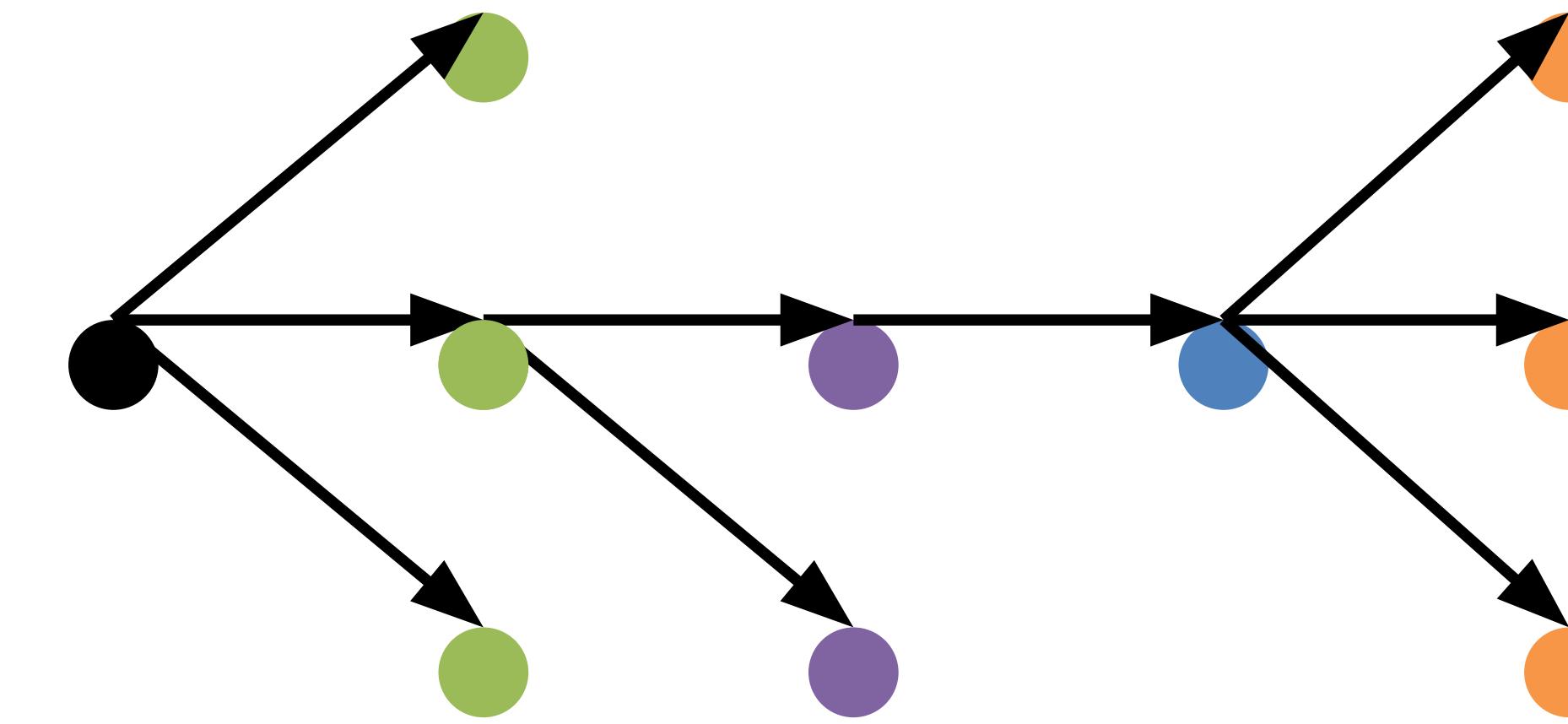
pen and paper

in the computer
`l.backward() # start backprop`

AUTOMATIC DIFFERENTIATION



many inputs, few **outputs**, work *backward*



few inputs, many outputs, work *forward*

<- *deep learning*

THE PLAN

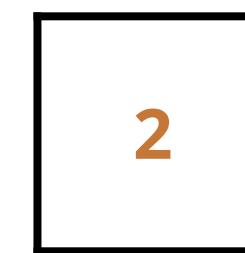
Tensors

Building computation graphs

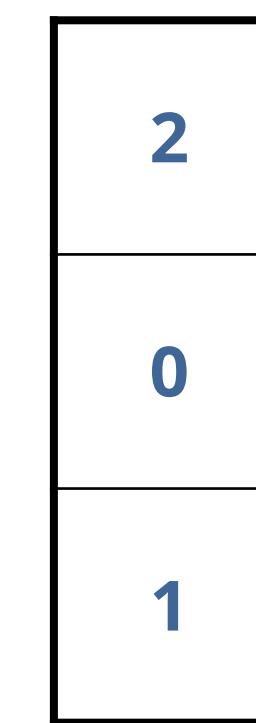
Working out backward functions

TENSORS (AKA MULTIDIMENSIONAL ARRAYS)

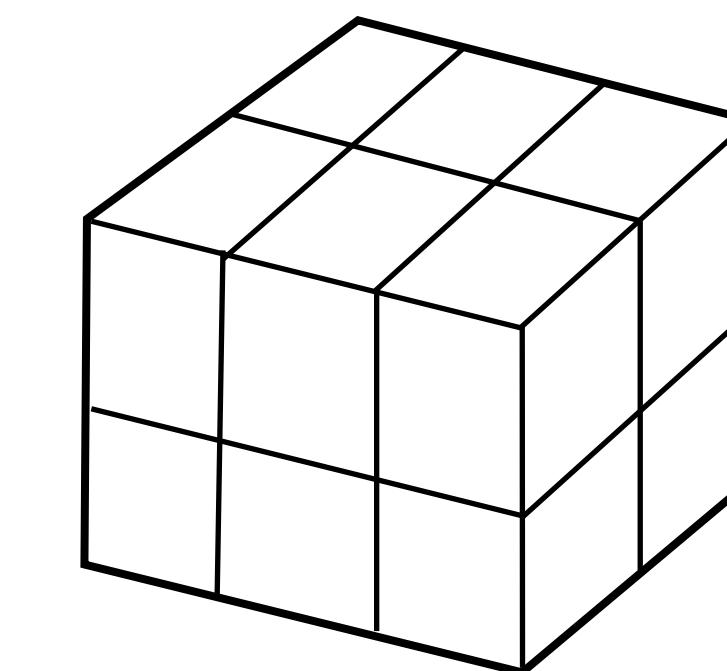
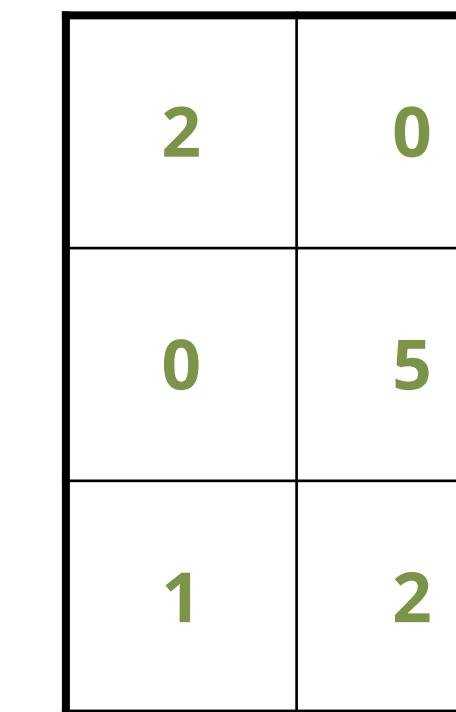
scalar



vector



matrix



0-tensor

1-tensor

shape: 3

2-tensor
shape:

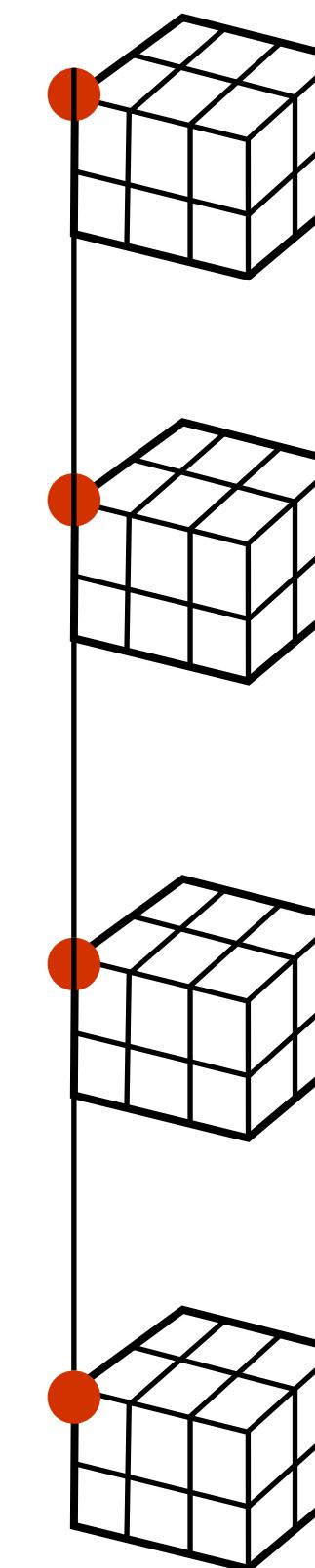
3x2

3-tensor
shape:

2x3x2

4-tensor
shape:

2x3x2x4



A CLASSIFICATION TASK AS TWO TENSORS



181	46	male
181	50	male
166	44	female
171	38	female
152	36	female
156	40	female
167	40	female
170	45	male
178	50	male
191	50	male

X =

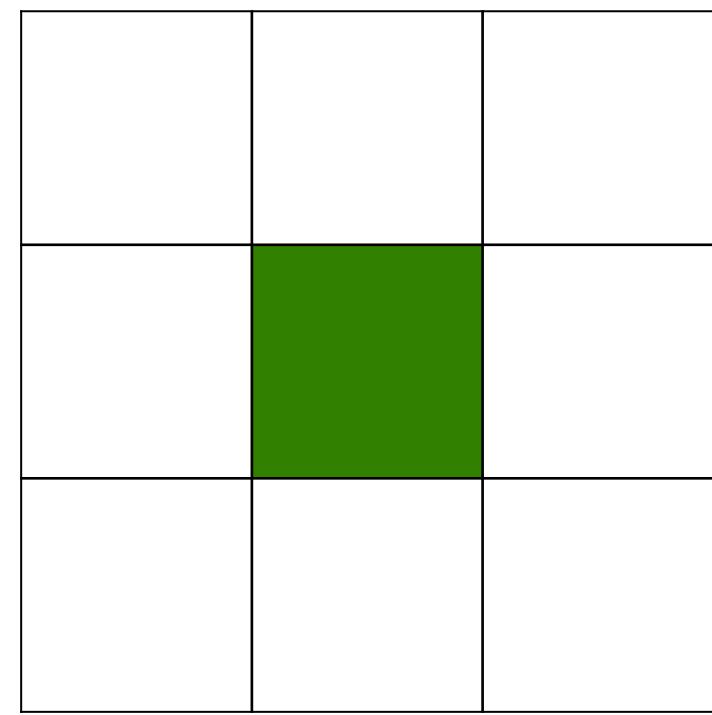
181	46
181	50
166	44
171	38
152	36
156	40
167	40
170	45
178	50
191	50
166	38

y =

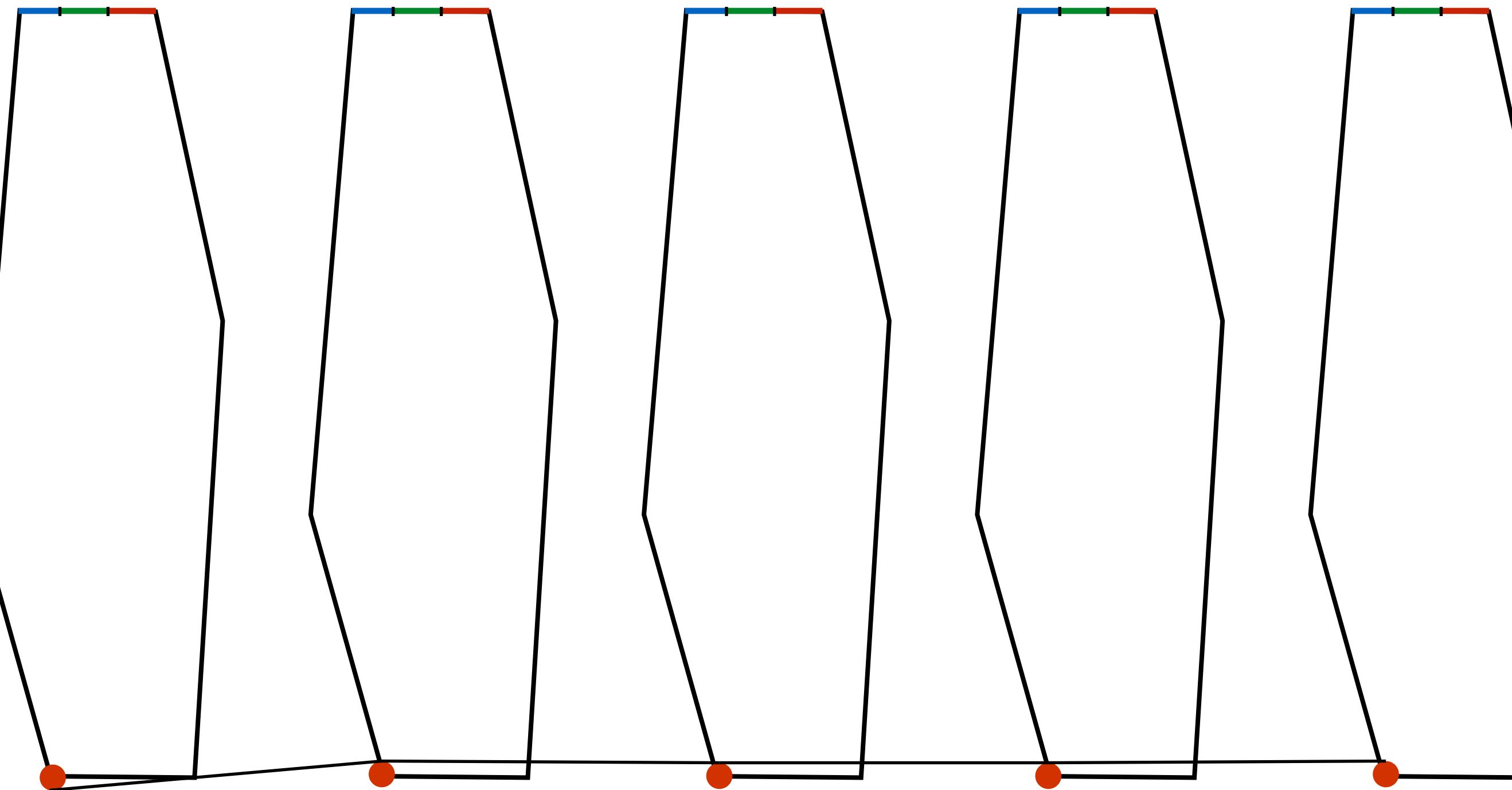
0
0
1
1
1
1
1
0
0
0
1

AN IMAGE AS A 3-TENSOR

pixel



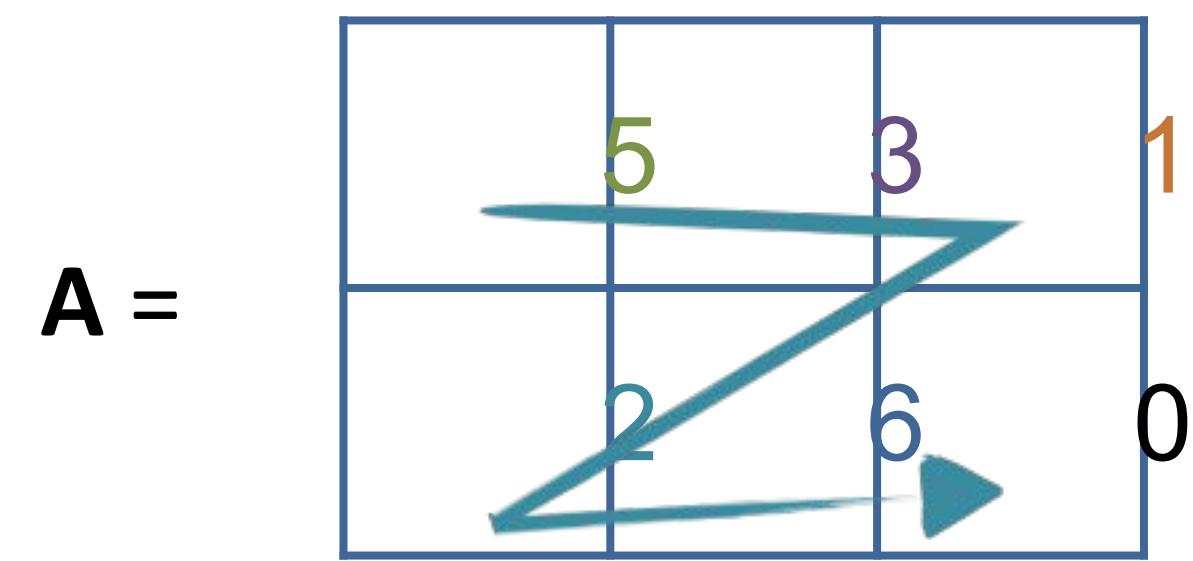
A DATASET OF IMAGES



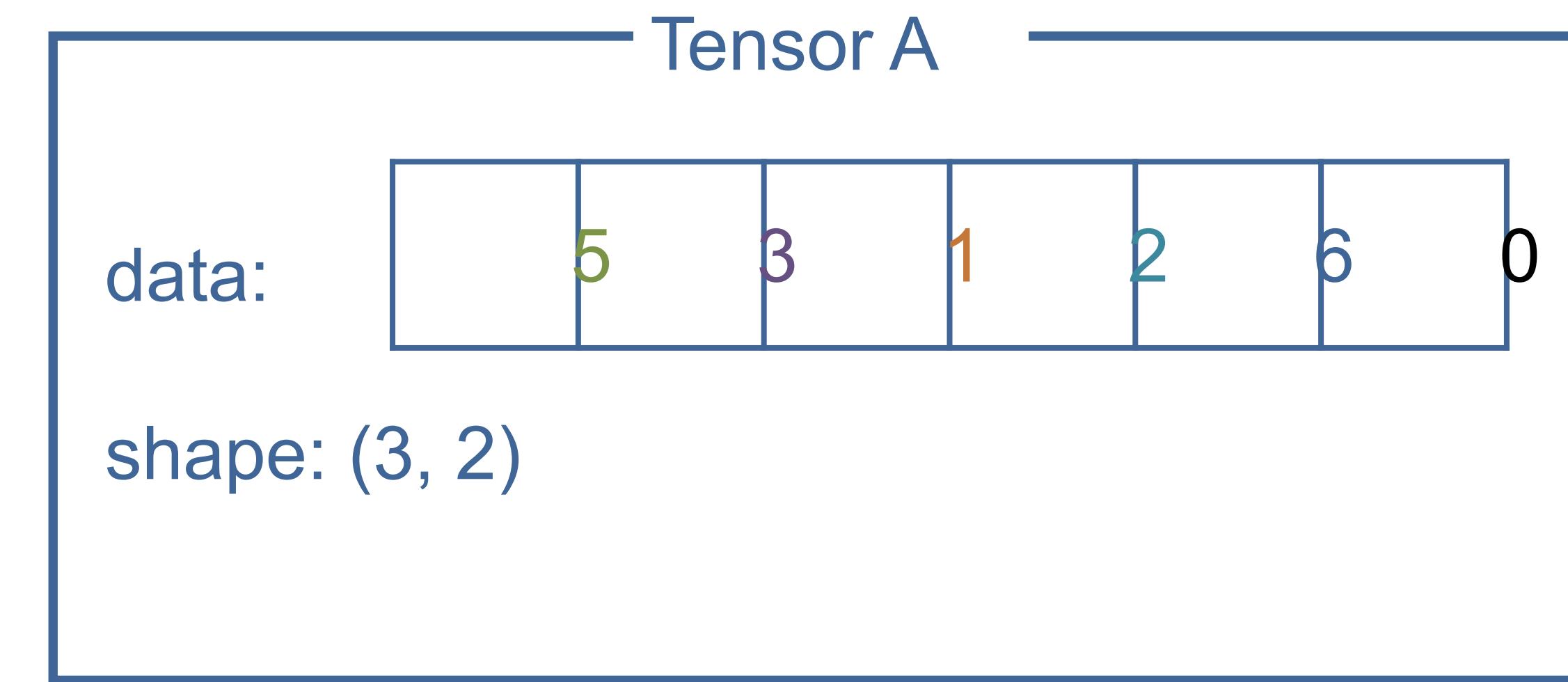
```
In [5]: from keras.datasets import cifar10  
  
(x_train, y_train), (x_test, y_test) = cifar10.load_data()  
  
x_train.shape
```

Out[5]: (50000, 32, 32, 3)

TENSORS IN MEMORY: ROW MAJOR ORDERING



row-major ordering



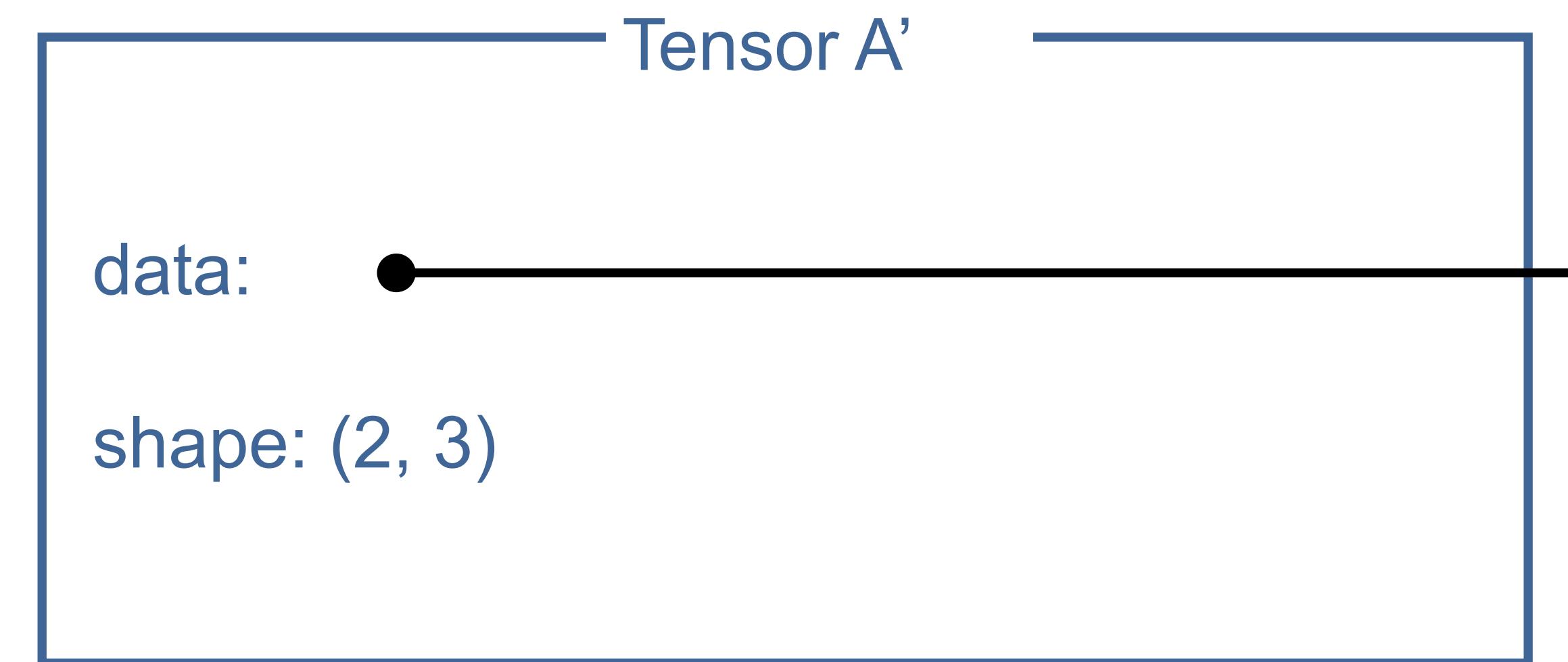
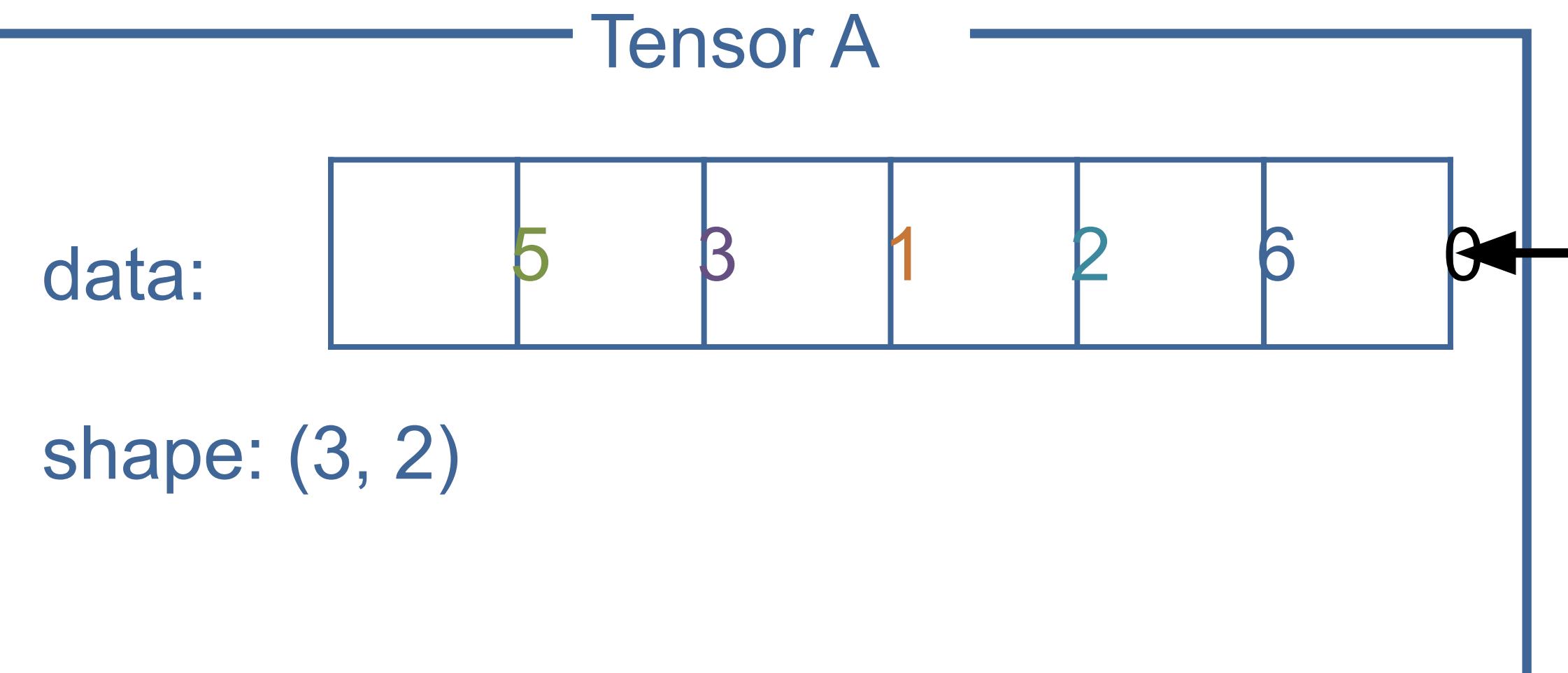
RESHAPE/VIEW

$A =$

	5	3	1
2		6	0

$A' =$

	5	3
1		2
6		0



RESHAPE/VIEW

$A =$

	5	3	1
	2	6	0

$a' =$

5	3	1	2	6	0
---	---	---	---	---	---

Tensor A

data:

5	3	1	2	6	0
---	---	---	---	---	---

shape: (3, 2)

Tensor a'

data:



shape: (6,)

EXAMPLE: NORMALIZE COLORS

```
images = load_images(...)
```

```
n, h, w, c = images.size()
```

```
images = images.reshape(n*h*w, c)
```

```
images = images / images.max(dim=0)  
# this makes no sense, but it is allowed.
```

```
images = images.reshape(n*h*w*c).reshape(n, h, w, c)
```

```
images = images.reshape(h, n, c, w)
```

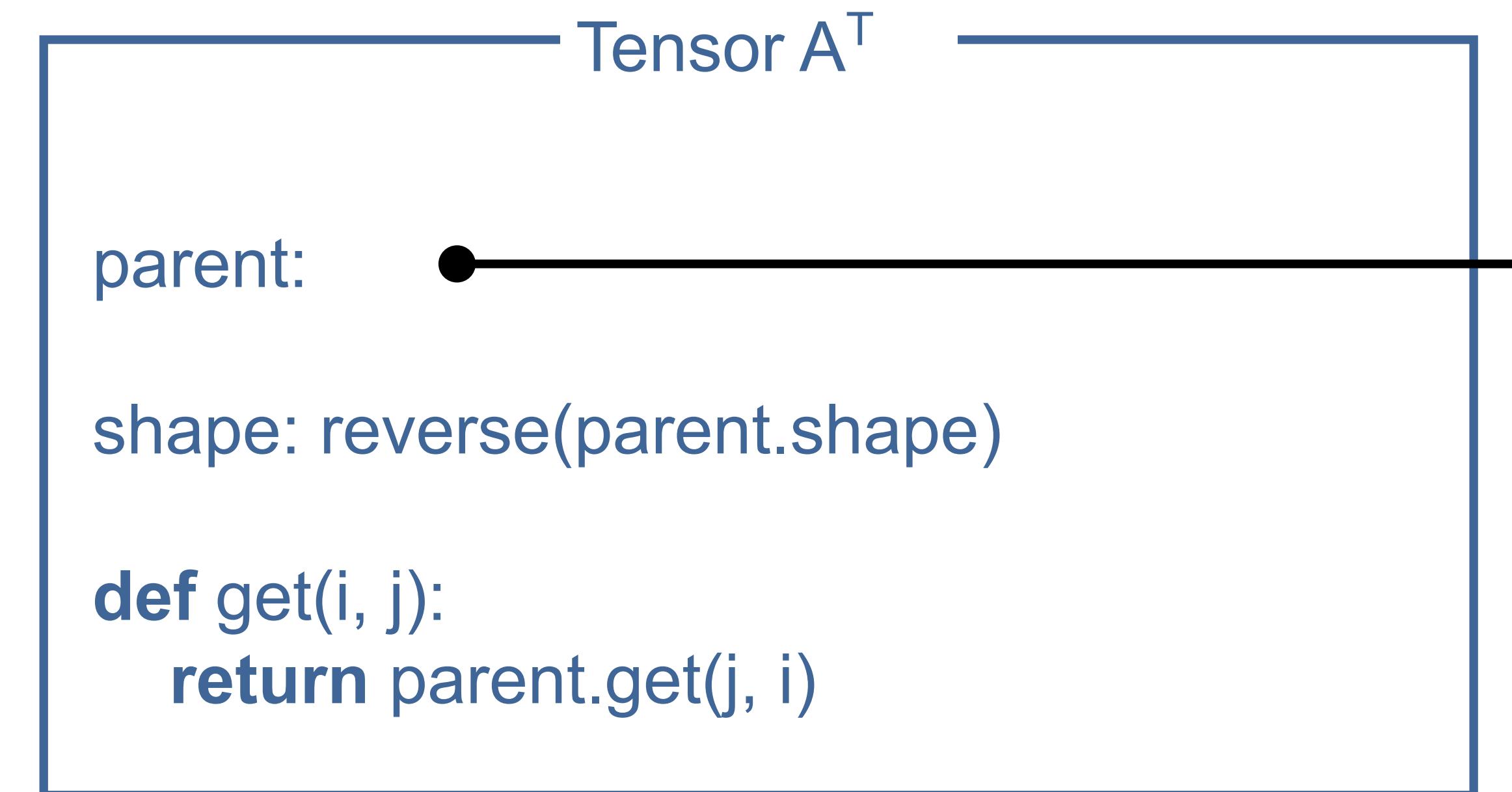
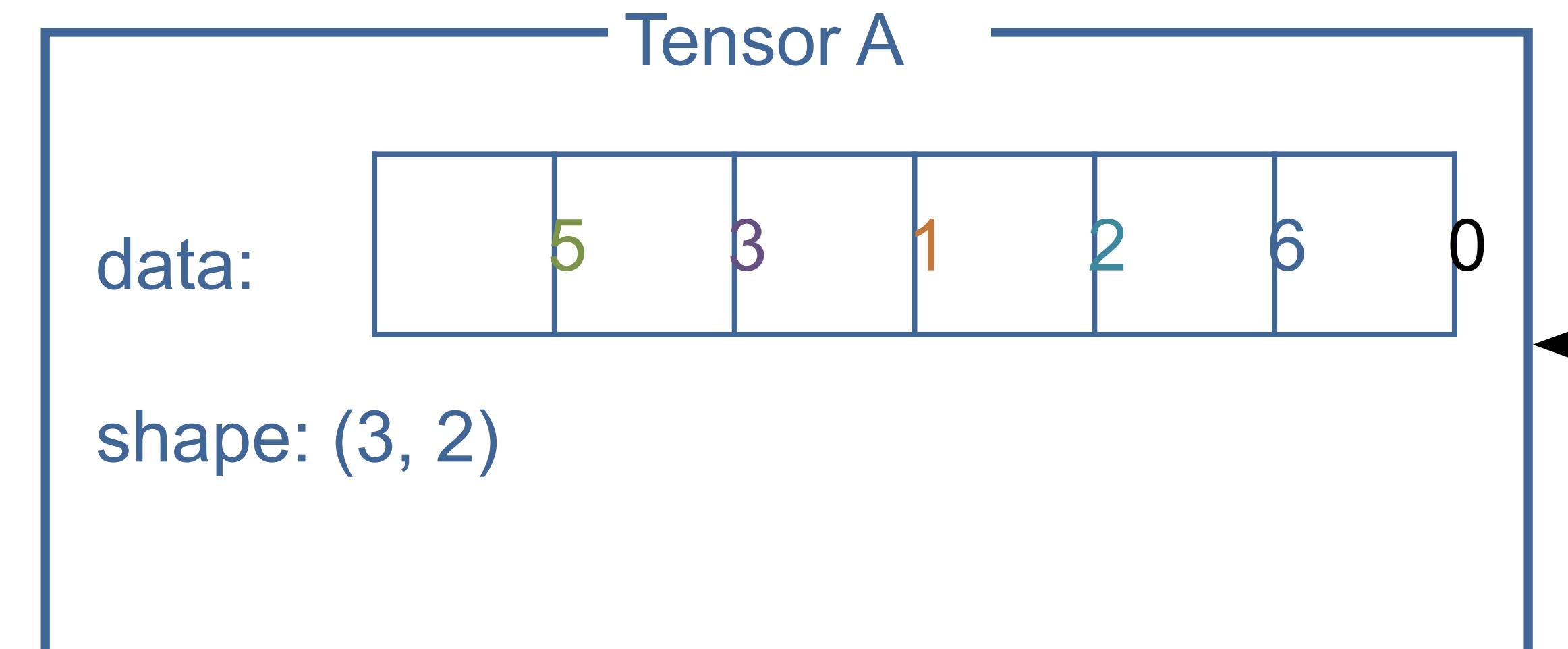
TRANSPOSE

$A =$

	5	3	1
2		6	0
			0

$A^T =$

	5	2
3		6
1		0



SLICE

A =

	5	3	1
	2	6	0

A[:, :2]=

	5	3
	2	6

data:

	5	3	1	2	6	0
--	---	---	---	---	---	---

shape: (3, 2)

Tensor A^T

parent:



shape: (2, 2)

```
def get(i, j):  
    return parent(i, j)
```

NON-CONTIGUOUS

Contiguous tensor: the data are directly laid out in row-major ordering for the shape with no gaps or shuffling required.

```
x = torch.randn(2, 3)  
x = x[:, 1:]  
x.view(6)
```

Traceback (most recent call last):

...

RuntimeError: view size is not compatible with input tensor's size and stride (at least one dimension spans across two contiguous subspaces). Use .reshape(...) instead.

x.contiguous(): make contiguous (by copying the data).

x.view(): reshape if possible without making contiguous.

x.reshape(): reshape, call contiguous() if necessary.

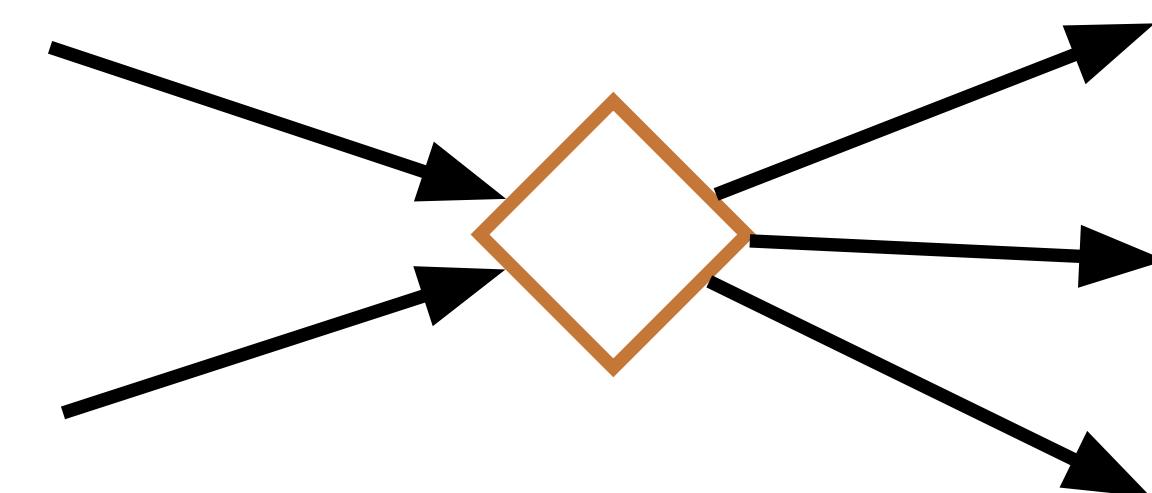
~~Tensors~~

Building computation graphs

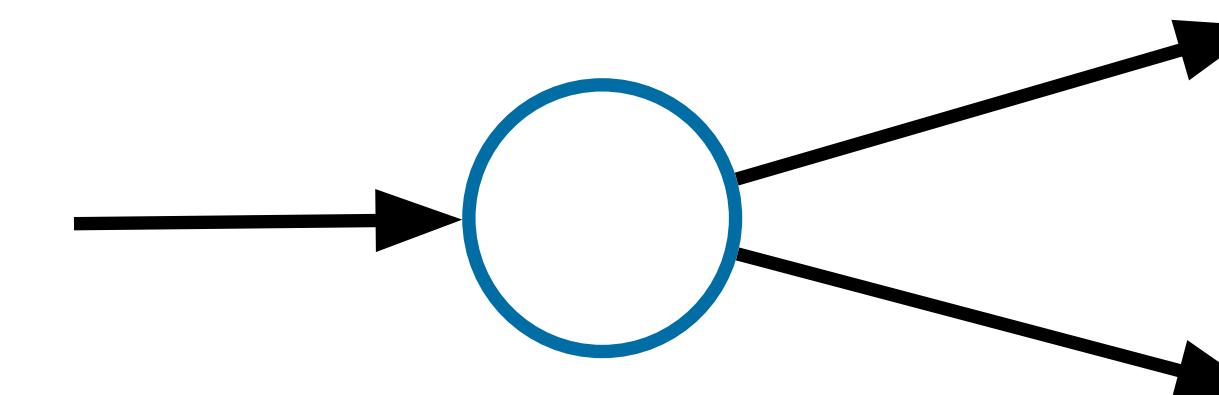
Working out backward functions

COMPUTATION GRAPHS

Operation nodes



Tensor nodes

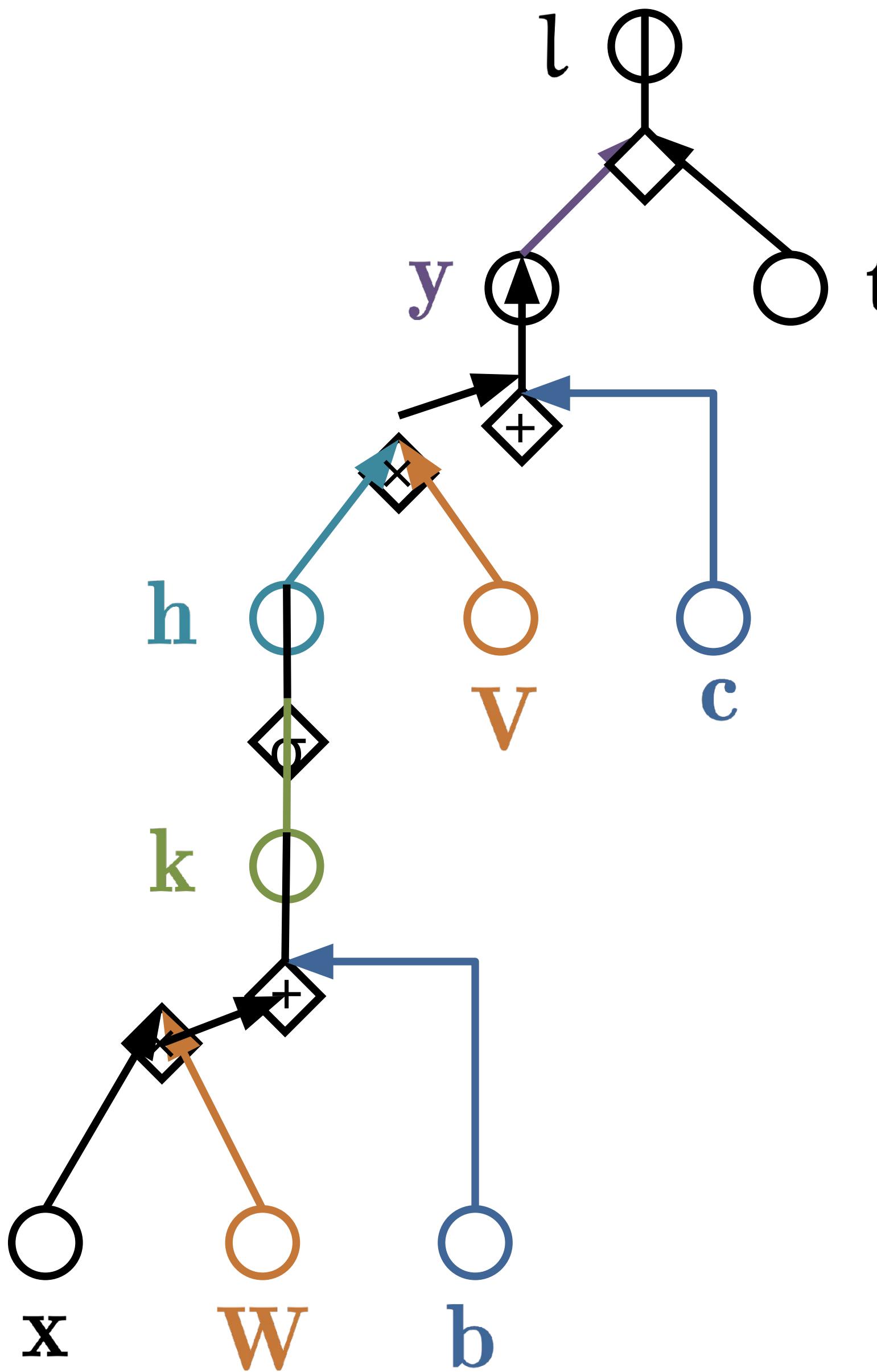


bipartite: only **op-to-tensor** or **tensor-to-op** edges

tensor nodes: one input, multiple outputs

operation nodes: multiple inputs and outputs

FEEDFORWARD NETWORK



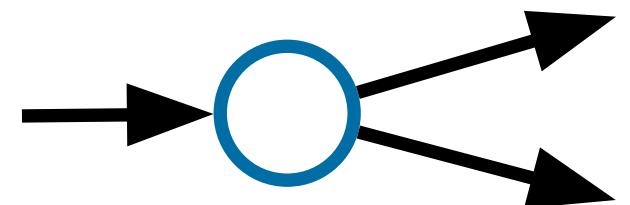
RESTATING OUR ASSUMPTIONS

One output \mathbf{l} node with a scalar value.

All gradients are of \mathbf{l} , with respect to the tensors on the tensor nodes.

IMPLEMENTATIONS

TensorNode:

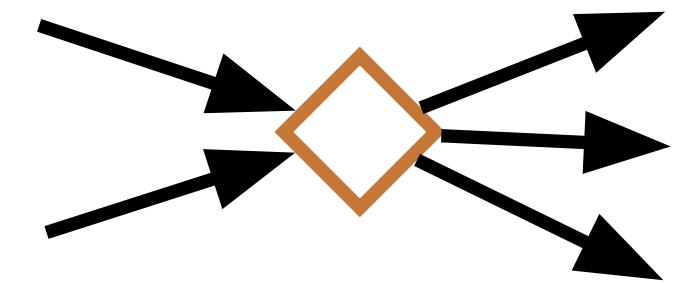


value: <Tensor>

gradient: <Tensor>

source: <OpNode>

OpNode:



inputs: List<TensorNode>

outputs: List<TensorNode>

op: Op

DEFINING AN OPERATION

Op:

<- store anything we need for the backward pass

$$f(\mathbf{A}) = \mathbf{B}$$

def forward(context, inputs):

$$\text{forward}_f : \mathbf{A} \mapsto \mathbf{B}$$

given the inputs, compute the outputs

...

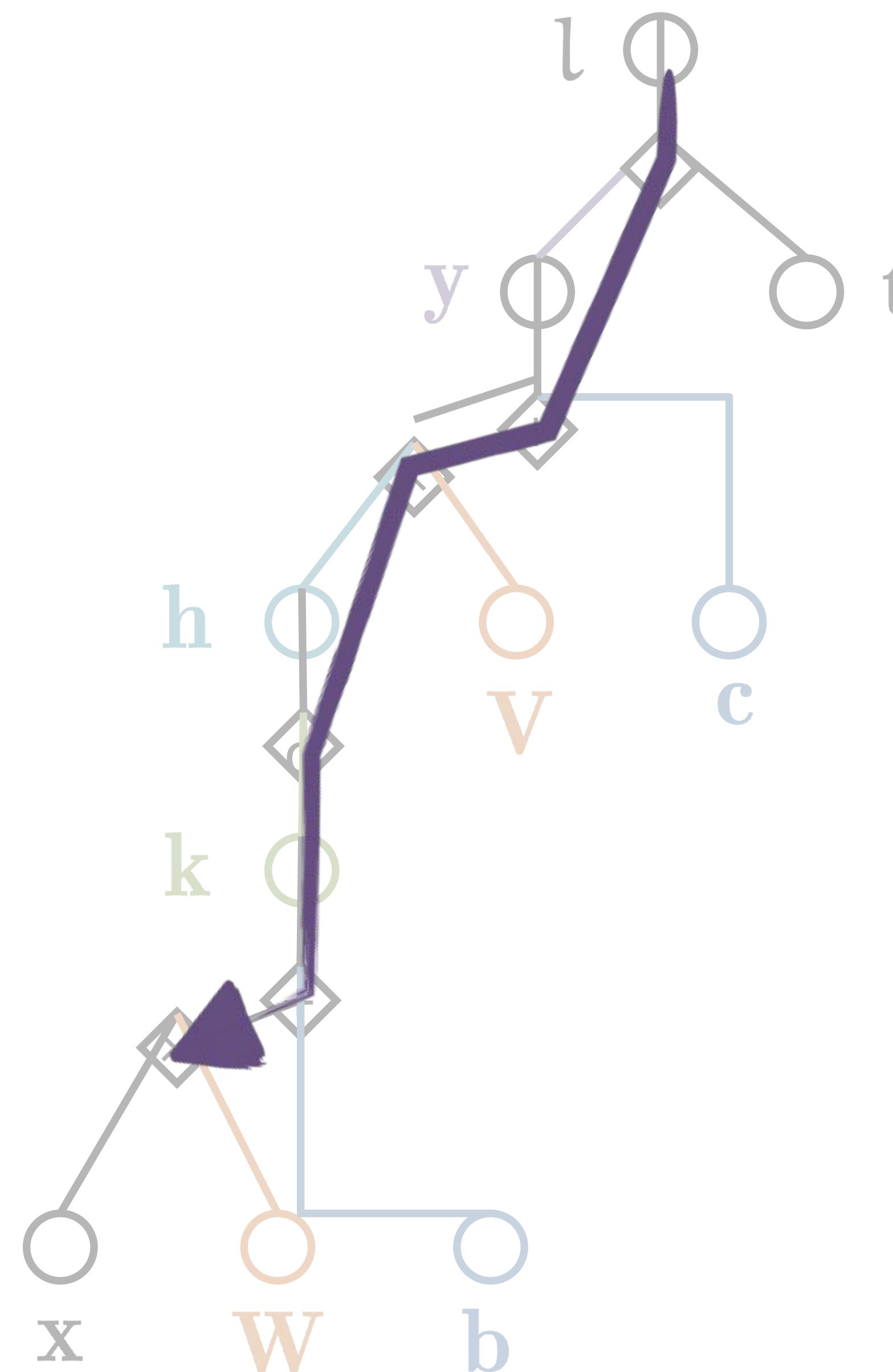
def backward(context, outputs_gradient):

$$\text{backward}_f : \mathbf{B}^\nabla \mapsto \mathbf{A}^\nabla$$

given the gradient of the loss wrt to the outputs

compute the gradient of the loss wrt to the inputs

BACKPROPAGATION



build a computation graph, perform forward pass

traverse the tree backward breadth-first

call backward for each OpNode

breadth-first ensures that the output gradients are known

add the computed gradients to the tensornodes

BUILDING COMPUTATION GRAPHS IN CODE

Lazy execution: build your graph, compile it, feed data through.

Eager execution: perform forward pass, keep track of computation graph.

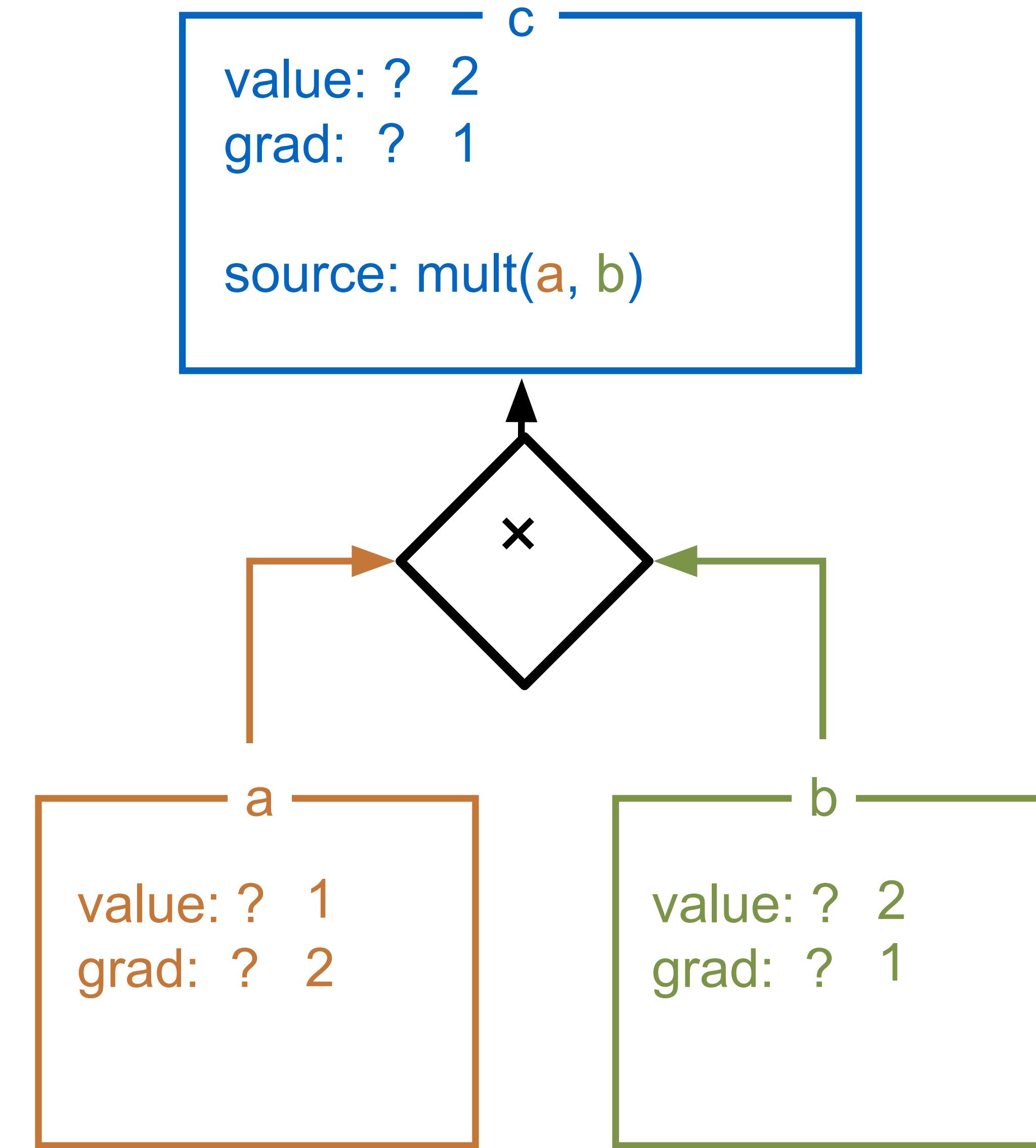
EXAMPLE: LAZY EXECUTION

```
a = TensorNode()  
b = TensorNode()
```

```
c = a * b
```

```
m = Model(  
    in=(a, b),  
    loss=c)
```

```
m.train((1, 2))
```



BUILDING THE COMPUTATION GRAPH: LAZY EXECUTION

Tensorflow 1.x default, Keras default

Define the computation graph.

- Compile it.
- Iterate backward/forward over the data

Fast. Many possibilities for optimization. Easy to serialise models. Easy to make training parallel.

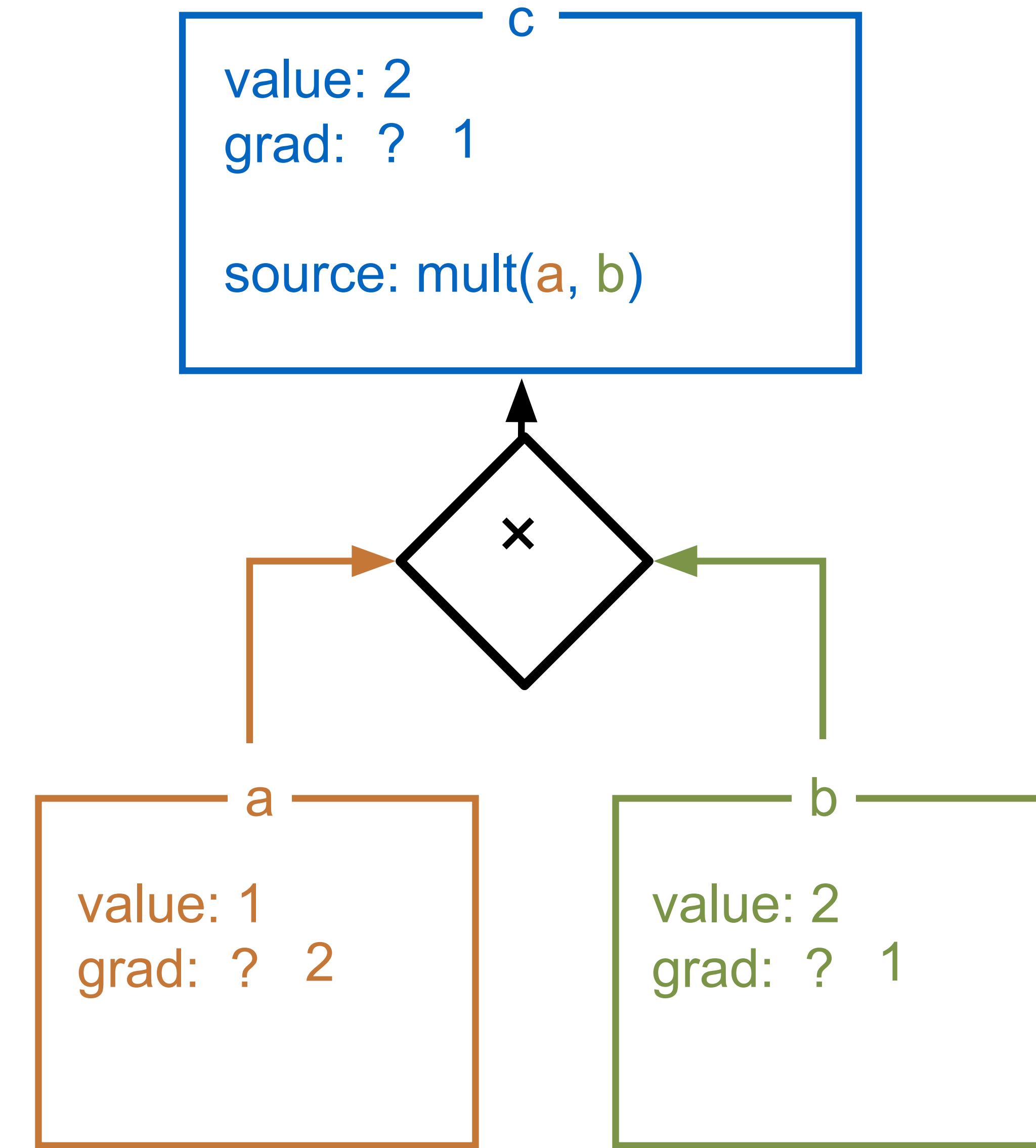
Difficult to debug. Model must remain static during training.

EXAMPLE: EAGER EXECUTION

```
a = TensorNode(value=1)  
b = TensorNode(value=2)
```

```
c = a * b
```

```
c.backward()
```



BUILDING THE COMPUTATION GRAPH: EAGER EXECUTION

PyTorch, Tensorflow 2.0 default, Keras option

- Build the computation graph on the fly during the forward pass.

Easy to debug, problems in the model occur as the module is executing.
Flexible: the model can be entirely different from one forward to the next.

More difficult to optimize. A little more difficult to serialize.

WORKING OUT THE BACKWARD FUNCTION

```
class Plus(Op):
```

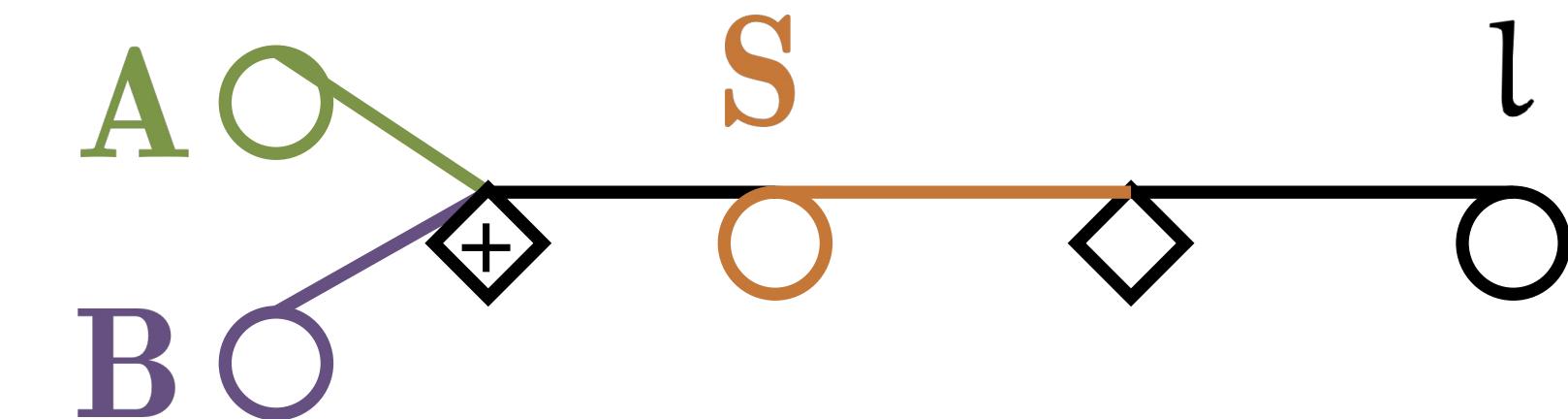
```
def forward(context, a, b):
```

a, b are matrices of the same size

```
return a + b
```

```
return goutput, goutput
```

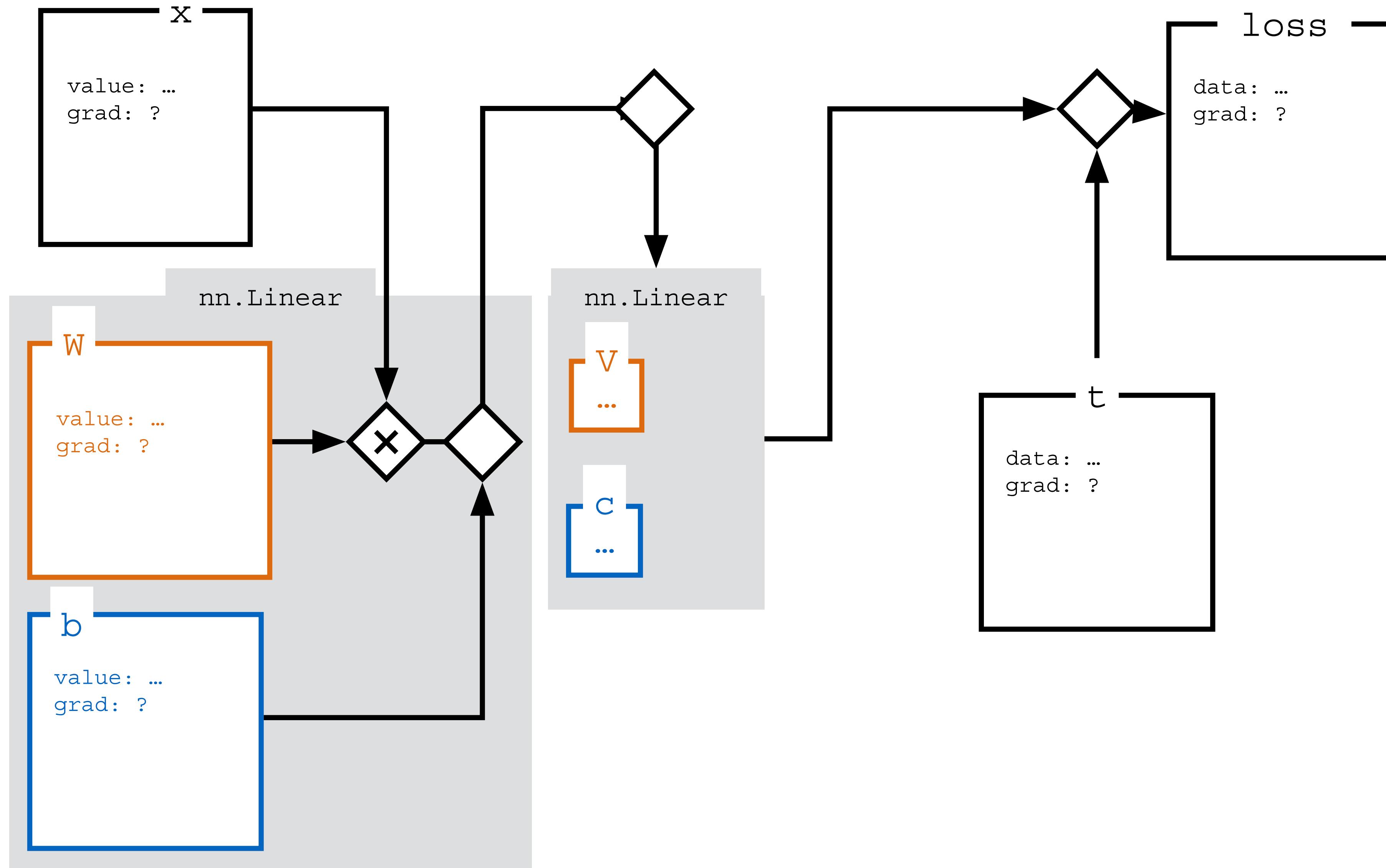
```
def backward(context, goutput):
```



$$\begin{aligned} \mathbf{A}_{ij}^\nabla &= \frac{\partial l}{\partial \mathbf{A}_{ij}} \\ &= \sum_{kl} \frac{\partial l}{\partial \mathbf{S}_{kl}} \frac{\partial \mathbf{S}_{kl}}{\partial \mathbf{A}_{ij}} = \sum_{kl} \mathbf{S}_{kl}^\nabla \frac{\partial \mathbf{S}_{kl}}{\partial \mathbf{A}_{ij}} \\ &= \sum_{kl} \mathbf{S}_{kl}^\nabla \frac{\partial [\mathbf{A} + \mathbf{B}]_{kl}}{\partial \mathbf{A}_{ij}} = \sum_{kl} \mathbf{S}_{kl}^\nabla \frac{\partial \mathbf{A}_{kl} + \mathbf{B}_{kl}}{\partial \mathbf{A}_{ij}} \\ &= \mathbf{S}_{ij}^\nabla \frac{\partial \mathbf{A}_{ij}}{\partial \mathbf{A}_{ij}} = \mathbf{S}_{ij}^\nabla \end{aligned}$$

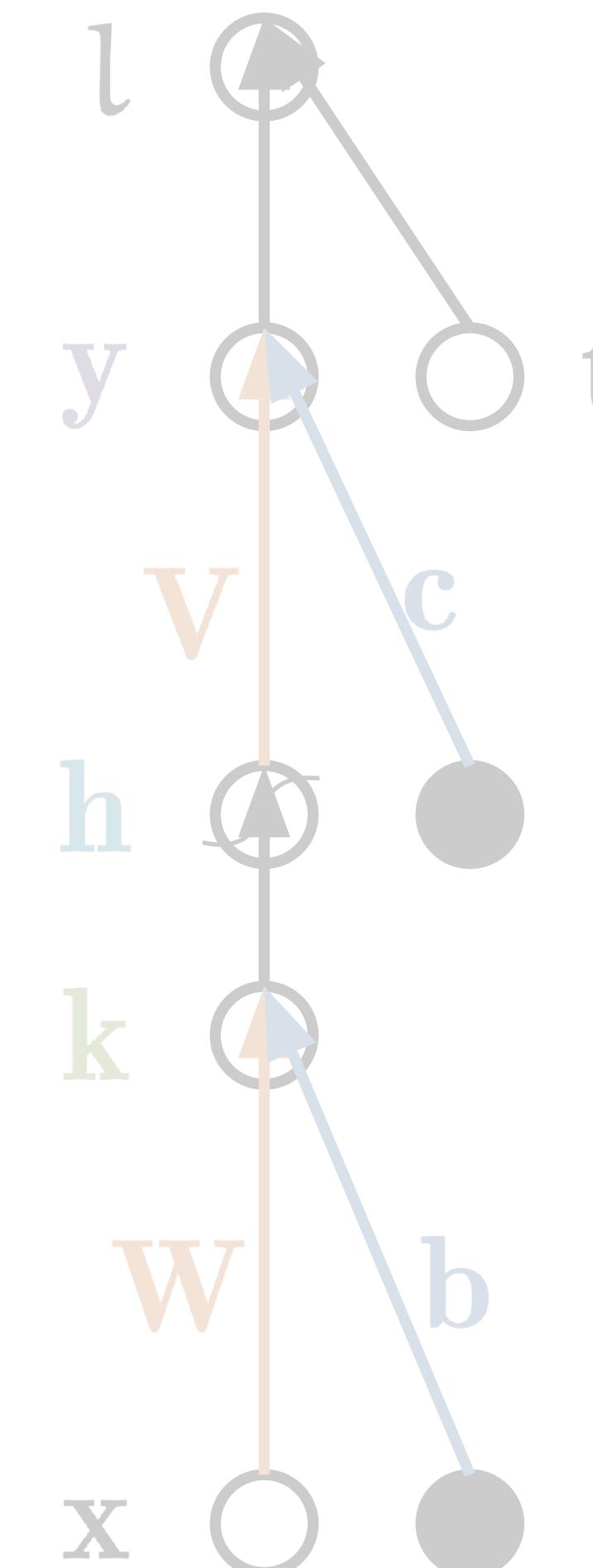
$$\mathbf{A}^\nabla = \mathbf{S}^\nabla \quad \mathbf{B}^\nabla = \mathbf{S}^\nabla$$

MODULES OR LAYERS



THE IDEAL

∅
+ x



pen and paper

$$k = \mathbf{w} \cdot \mathbf{x} + b$$

$$h = \text{sigmoid}(k)$$

$$y = \mathbf{v} \cdot \mathbf{h} + c$$

$$l = (y - t)^2$$

in the computer
`l.backward() # start backprop`

compute forward pass
and
build computation graph

Tensors

Building computation graphs

Working out backward functions

BACKWARD: A FEW MORE EXAMPLES

Sigmoid

Row-wise sum

Expand

SIGMOID

```
class Sigmoid(Op):
```

```
def forward(context, X):
```

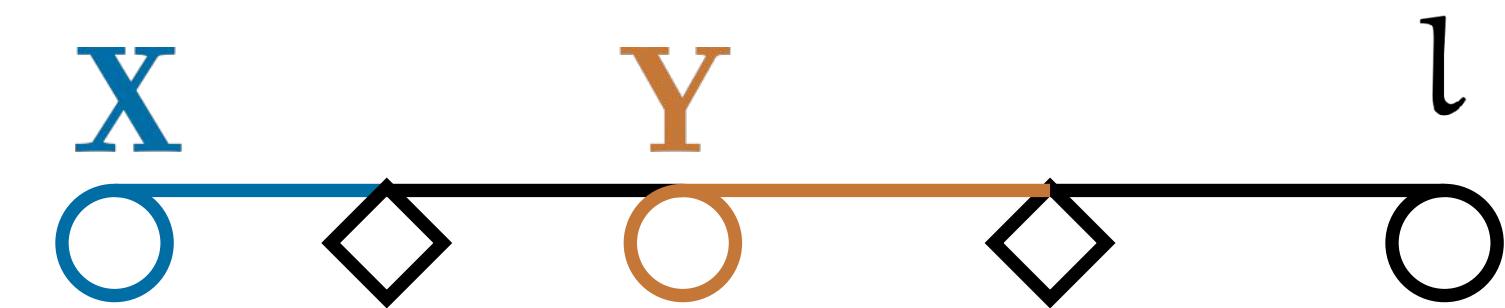
X is a tensor of any shape

```
sigx = 1 / (1 + (-X).exp())
```

```
    sigx = context['sigx']
```

```
r
```

```
return goutput * sigx * (1 - sigx)
```



$$\begin{aligned} \mathbf{X}_{ijk}^\nabla &= \sum_{abc} \mathbf{Y}_{abc}^\nabla \frac{\partial \mathbf{Y}_{abc}}{\partial \mathbf{X}_{ijk}} \\ &= \sum_{abc} \mathbf{Y}_{abc}^\nabla \frac{\partial \sigma(\mathbf{X}_{abc})}{\partial \mathbf{X}_{ijk}} \\ &= \mathbf{Y}_{ijk}^\nabla \frac{\partial \sigma(\mathbf{X}_{ijk})}{\partial \mathbf{X}_{ijk}} \\ &= \mathbf{Y}_{ijk}^\nabla \sigma(\mathbf{X}_{ijk})(1 - \sigma(\mathbf{X}_{ijk})) \end{aligned}$$

$$\mathbf{X}^\nabla = \mathbf{Y}^\nabla \otimes \sigma(\mathbf{X}) \otimes (1 - \sigma(\mathbf{X}))$$

ROW-WISE SUM

```
class RowSum(Op):
```

```
    def forward(context, X):
```

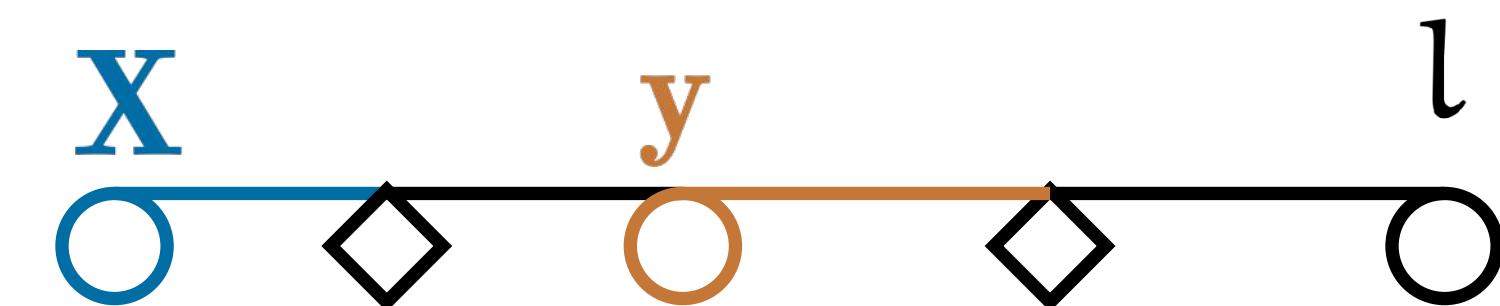
```
# x is a matrix
```

```
sumd = x.sum(axis=1)
```

```
n, m = gy.shape[0], context['m']
```

```
return gy[None, :].expand(n, m)
```

```
def backward(context, gy):
```



$$\begin{aligned} \mathbf{x}_{ij}^\nabla &= \sum_k \mathbf{y}_k^\nabla \frac{\partial \mathbf{y}_k}{\partial \mathbf{x}_{ij}} = \sum_k \mathbf{y}_k^\nabla \frac{\partial \sum_l \mathbf{x}_{kl}}{\partial \mathbf{x}_{ij}} \\ &= \sum_{kl} \mathbf{y}_k^\nabla \frac{\partial \mathbf{x}_{kl}}{\partial \mathbf{x}_{ij}} = \mathbf{y}_i^\nabla \frac{\partial \mathbf{x}_{ij}}{\partial \mathbf{x}_{ij}} \\ &= \mathbf{y}_i^\nabla \end{aligned}$$

$$\mathbf{x}^\nabla = \mathbf{y}^\nabla \mathbf{1}^\top$$

EXPAND

```
class Expand(Op):
```

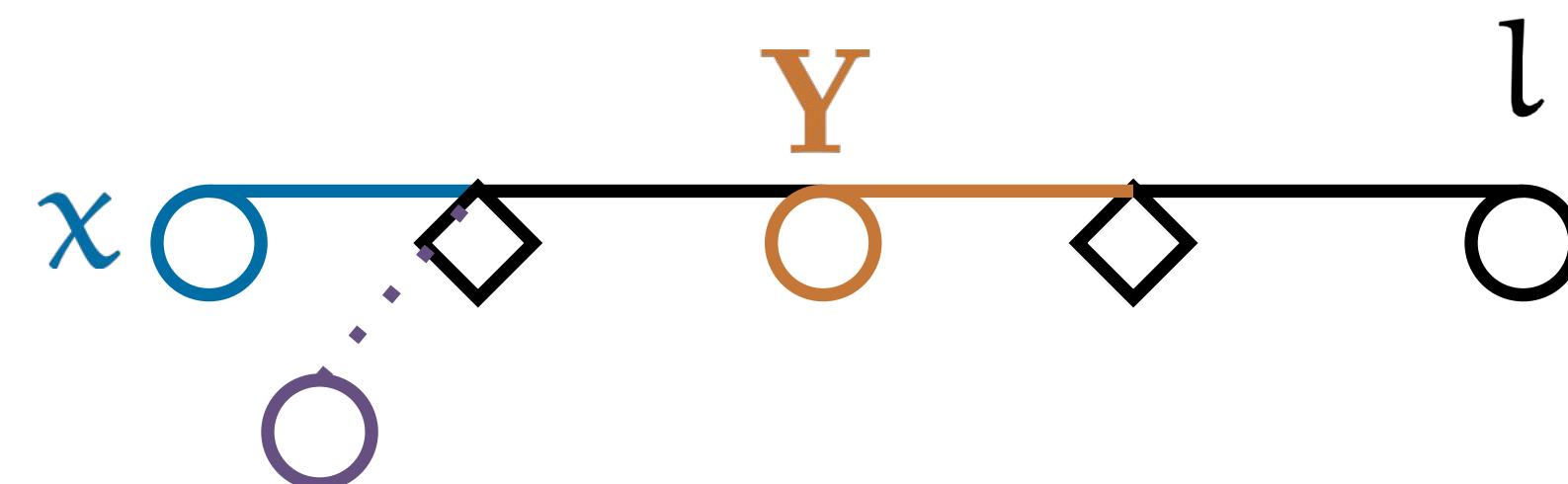
```
def forward(context, x, size):
```

x is a scalar

```
return np.full(x.size=size)
return gy.sum(), None
```

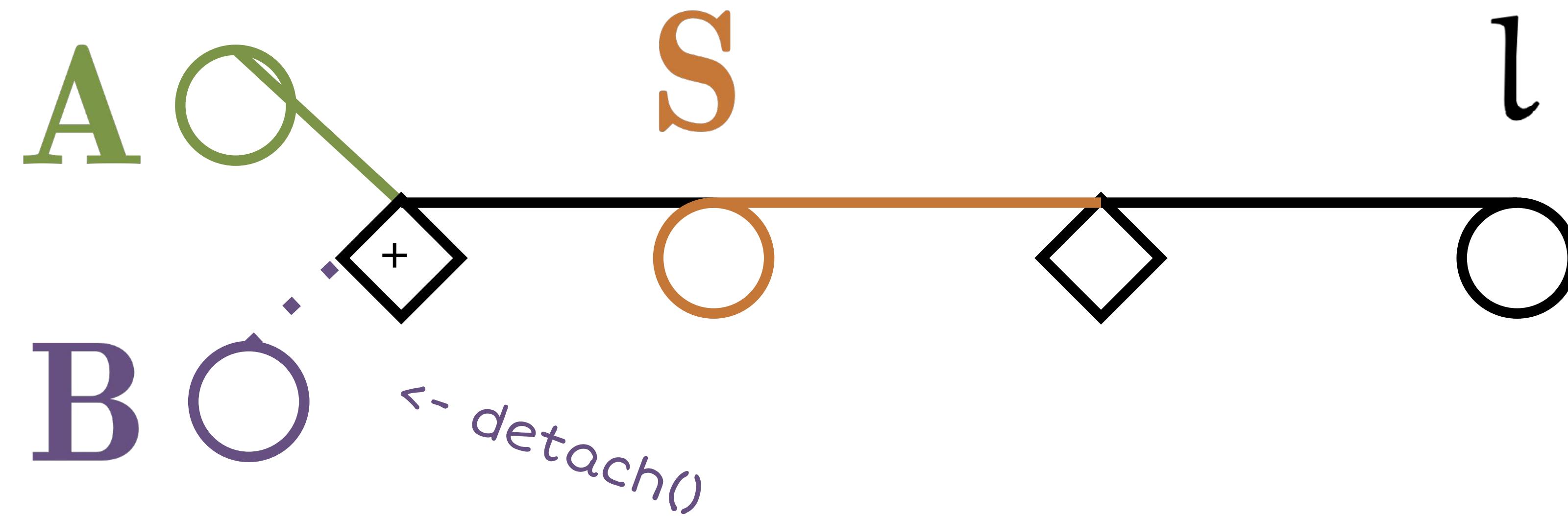
```
def backward(context, gy):
```

???



$$\begin{aligned}x^\nabla &= \sum_{ab} Y_{ab}^\nabla \frac{\partial Y_{ab}}{\partial x} \\&= \sum_{ab} Y_{ab}^\nabla \frac{\partial x}{\partial x} \\&= \sum_{ab} Y_{ab}^\nabla\end{aligned}$$

DETACH()



Tensors

Building computation graphs

Working out backward functions

neural networks: historically inspired by brain structure, but really just a pile of differentiable linear algebra.

backpropagation: work out local gradients symbolically, global gradients numerically. Work from back to front, applying the chain rule.

tensor backpropagation: assume scalar loss l , always compute the gradient of the loss.

automatic differentiation: Build computation graph on the fly.

¹²₈ Backpropagation as a *breadth-first* walk down from the loss.

THANK YOU FOR YOUR ATTENTION

dlvu@peterbloem.nl

target class

1 0

softmax activation

sigmoid activation

