



CS 540 Introduction to Artificial Intelligence

Neural Networks: Backpropagation

University of Wisconsin—Madison
Fall 2025, Section 3
October 13, 2025

Announcements

- HW5 due Friday 10/6 at 11:59 pm
- Midterm exam:

Thursday 10/23
7:30 pm to 9:00 pm
Humanities Building, Room 3650

ML: Unsupervised Learning

ML: Linear Regression

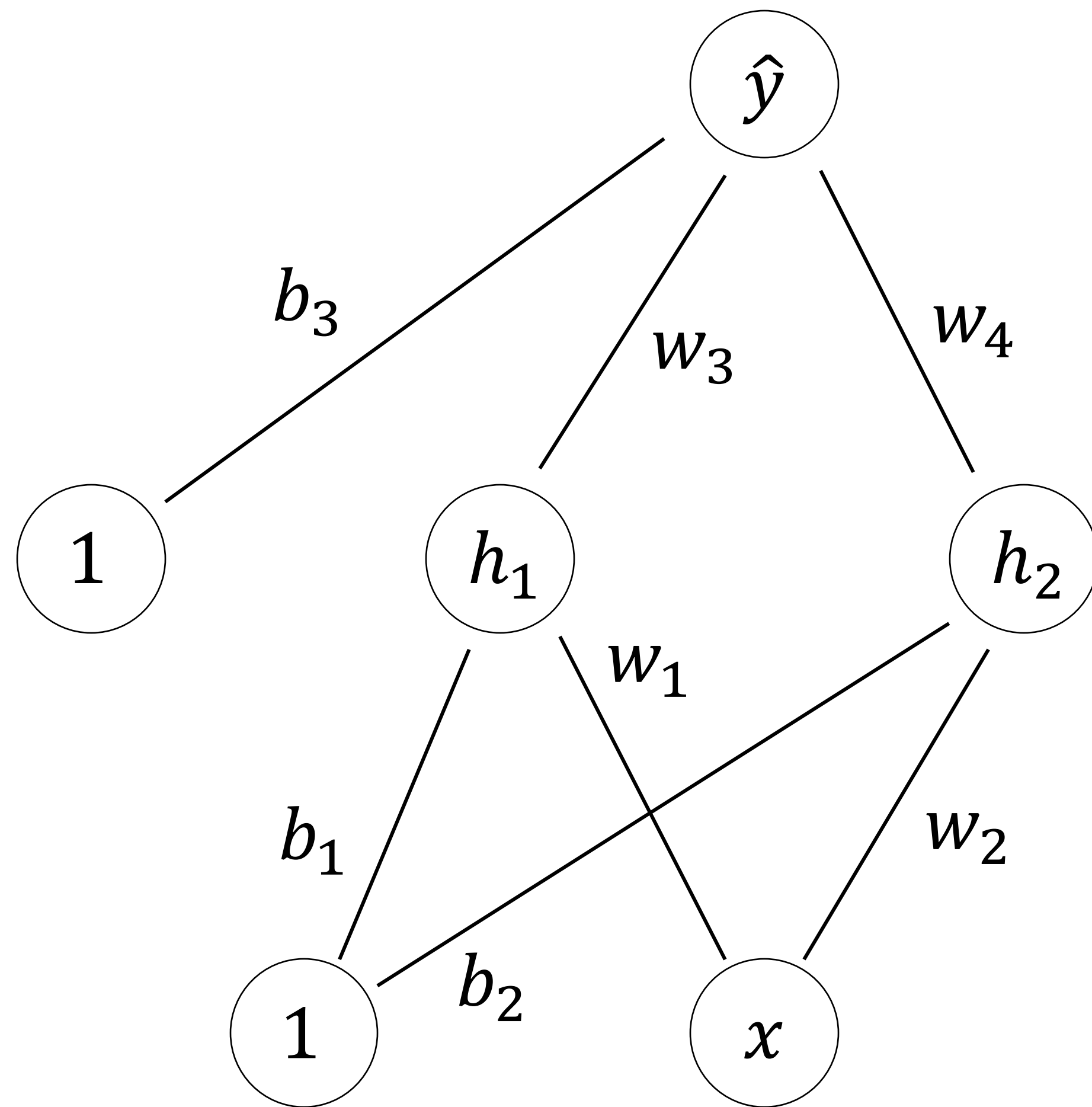
ML: K - Nearest Neighbors & Naive Bayes

ML: Neural Networks I (Perceptron)

ML: Neural Networks II

ML: Neural Networks III

A Simple Multilayer Perceptron



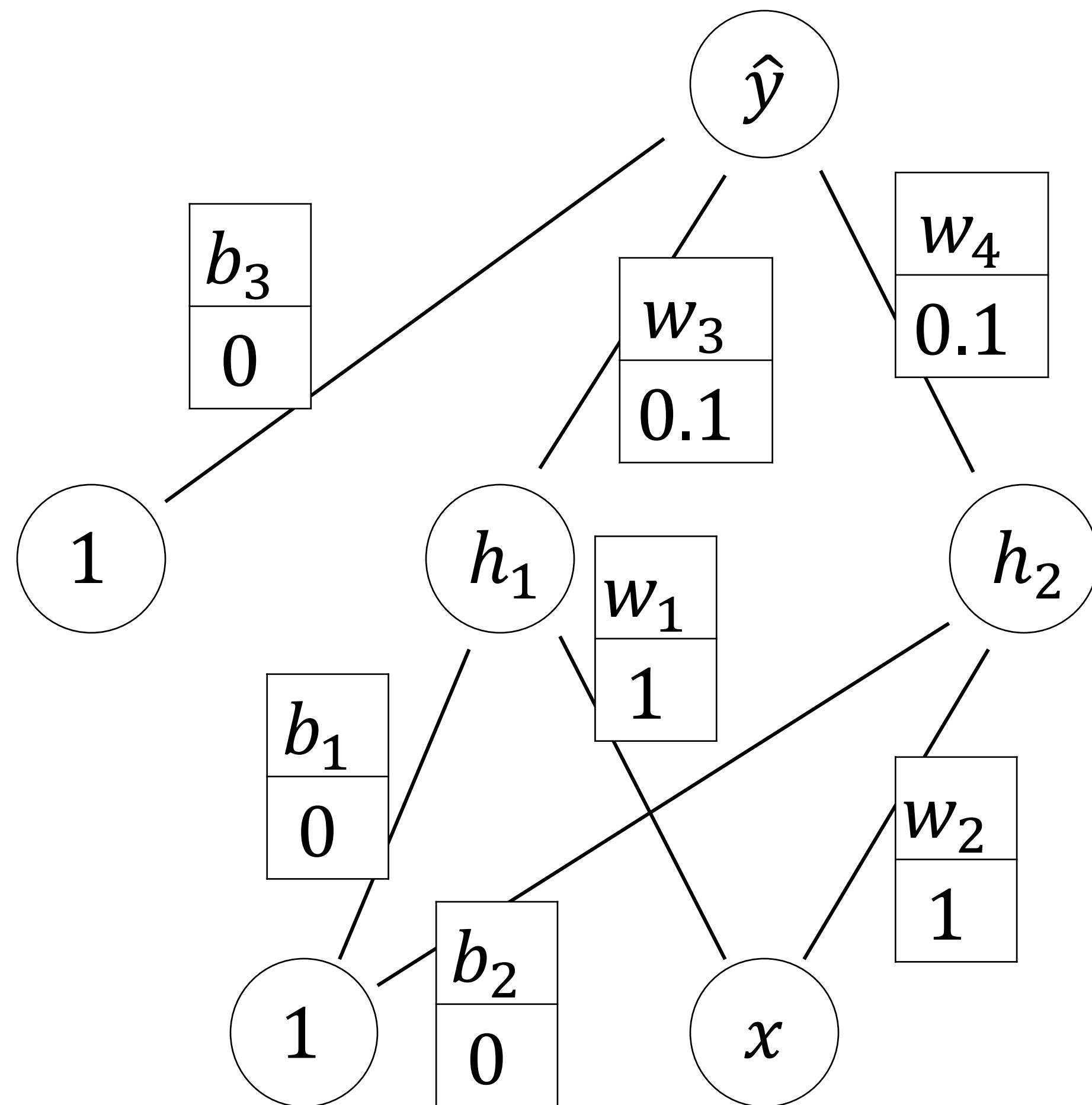
Hidden node: ReLU activations

$$h_1 = \max\{0, w_1x + b_1\}$$

Output node: no activation

$$\hat{y} = w_3h_1 + w_4h_2 + b_3$$

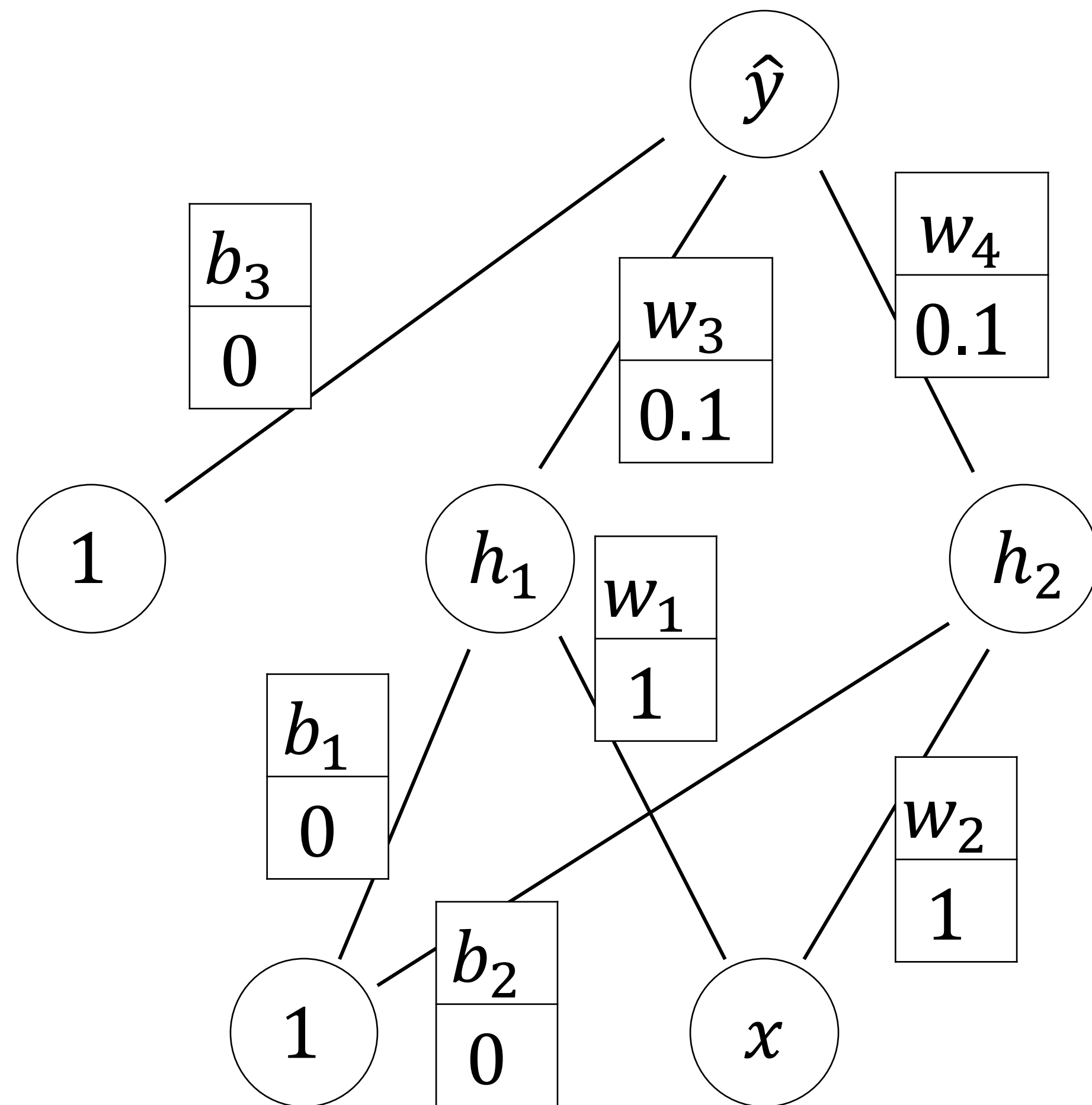
The Same Network, With Numbers



Compute: \hat{y} on input $x = 2$

This is a “forward pass”
through the network

Making Predictions Match the Labels

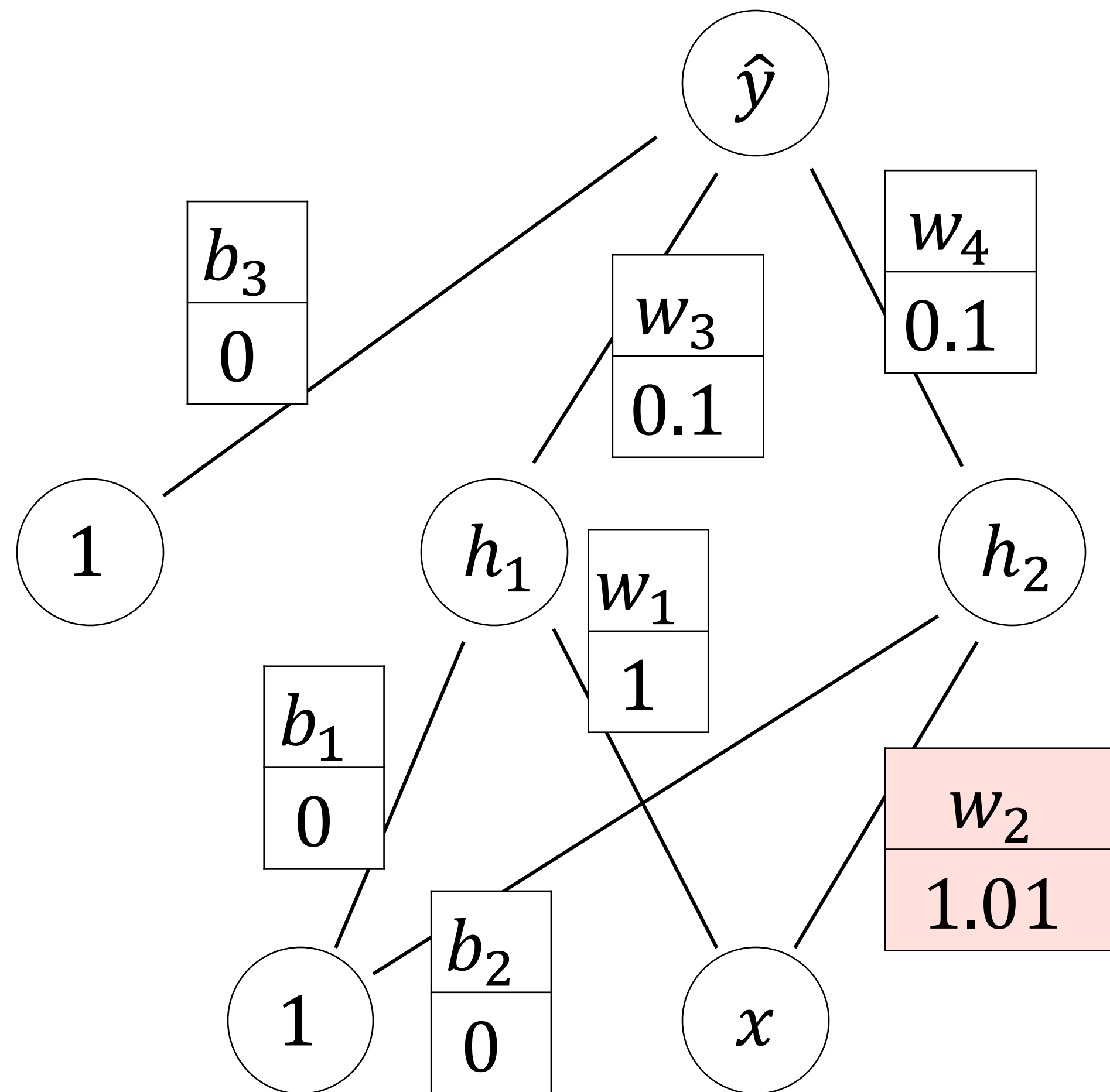


Suppose we have one labeled example

$$(x, y) = (2, 1)$$

If we increase w_2 a little, will our prediction get better or worse?

Change Weights, Change Prediction



Compute: \hat{y} on input $x = 2$

Formalize the Intuition with Calculus

Use the squared loss:

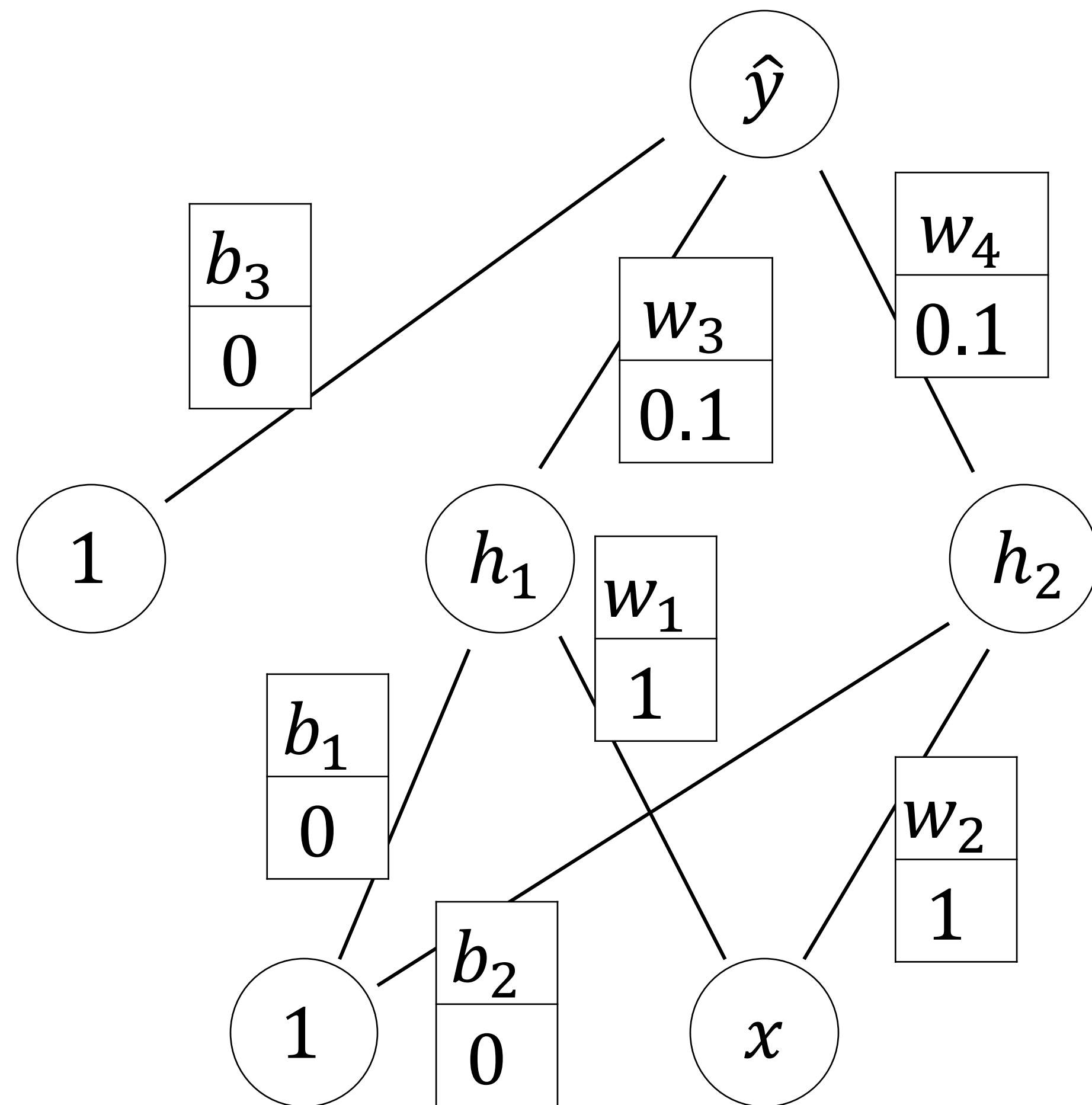
$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

Partial derivative: $\frac{\partial \ell}{\partial w_2}$

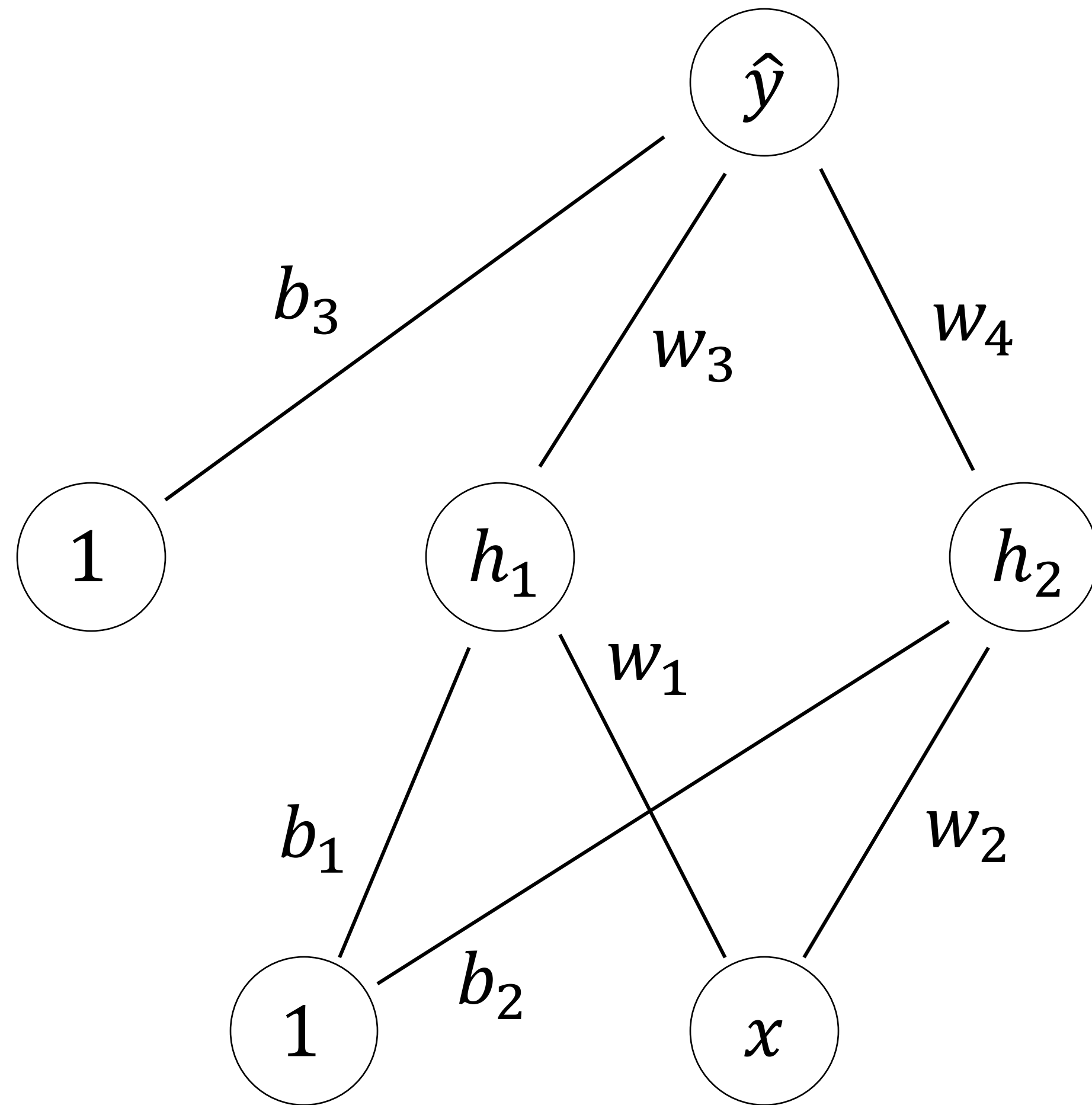
How does loss change when w_2 changes?

For these weights, with $(x, y) = (2, 1)$

$$\frac{\partial \ell}{\partial w_2} = -0.24$$



Notation for all the parameters at once



Collect all
parameters into a
single vector θ

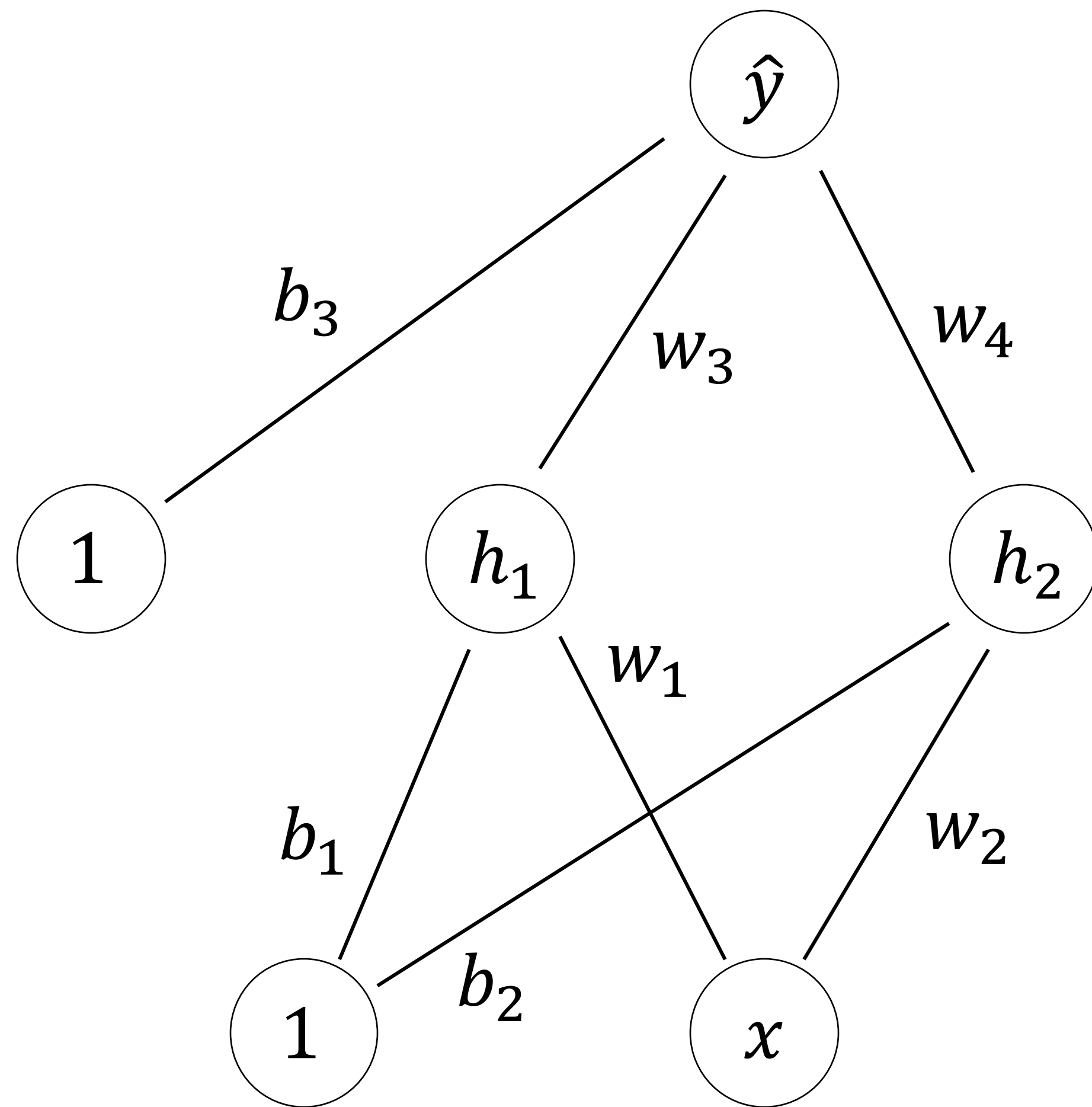
$$\theta = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Write prediction function: $\hat{y} = f(x; \theta)$

(Squared) loss on a dataset:

$$L(\theta; X, y) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i, y_i) = \frac{1}{n} \sum_{i=1}^n (f(x; \theta) - y)^2$$

Notation for all the parameters at once



Collect all
parameters into a
single vector θ

$$\theta = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Write all partial
derivatives as one
vector: the **gradient**

$$\nabla_{\theta} L(\theta; X, y) = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial b_3} \end{bmatrix}$$

Dissecting our Gradient Notation

- When f maps vector $x \in \mathbb{R}^d$ to scalar $f(x) \in \mathbb{R}$

$$\nabla f(x) \in \mathbb{R}^d$$

- We will work with functions of many vectors (and matrices!)

$$L(\theta; X, y)$$

- Write gradient with respect to θ

$$\nabla_{\theta} L(\theta_{t-1}; X, y)$$

Gradient Descent

Input: dataset (X, y) , loss function L , number of steps T , step size η

1. Initialize θ_0
2. For $t = 1, 2, \dots, T$
3. Calculate $g_t = \nabla_{\theta} L(\theta_{t-1}; X, y)$
4. Update $\theta_t \leftarrow \theta_{t-1} - \eta g_t$
5. Return θ_T

(Minibatch) Stochastic Gradient Descent

- Gradient descent uses loss on entire dataset every step:

$$\nabla L(\theta; X, y) = \sum_{i=1}^n \nabla \ell(\theta; x_i, y_i)$$

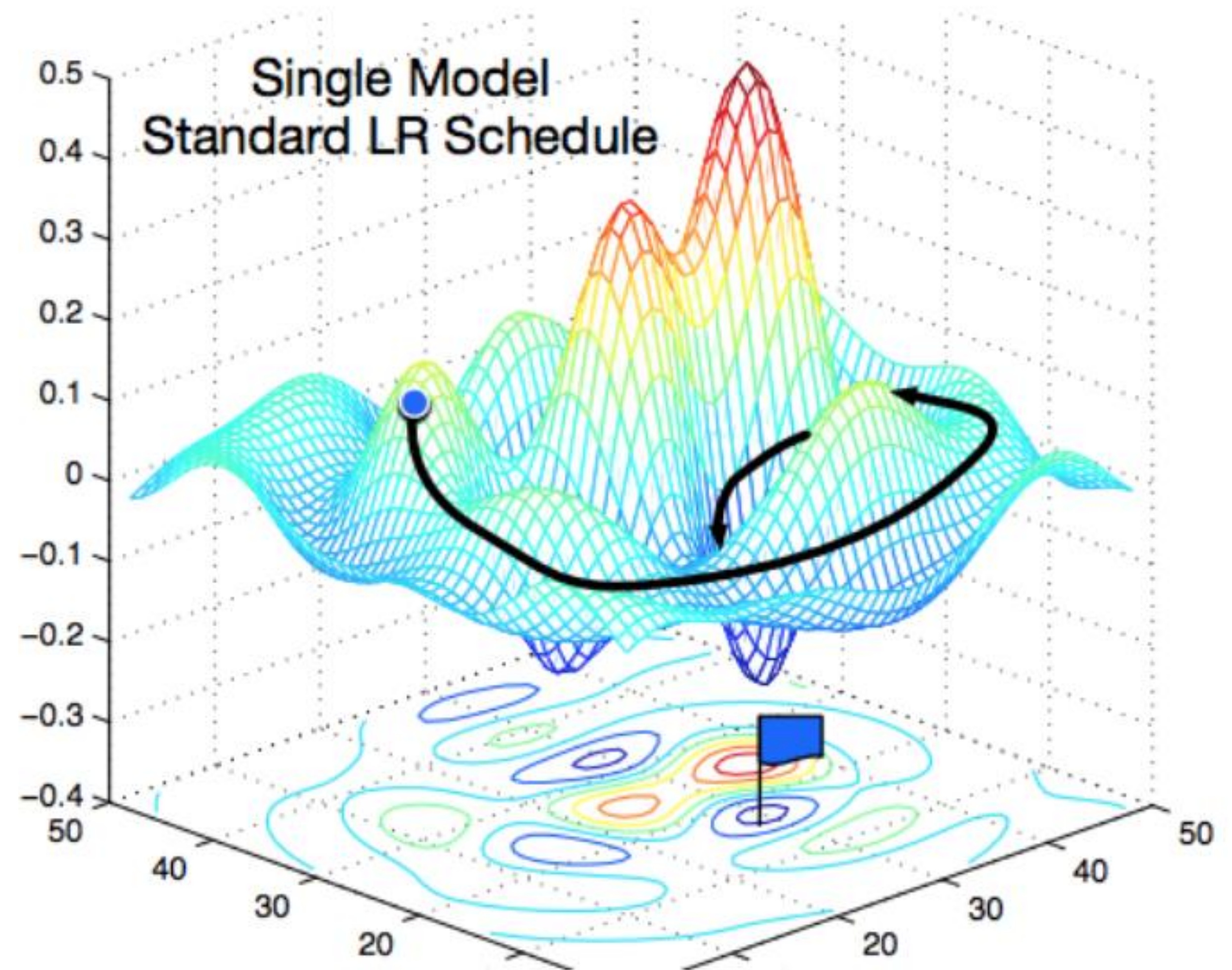
- This is extremely inefficient!
- On big datasets, better to update based on a few examples

Stochastic Gradient Descent

Input: dataset (X, y) , loss function L , number of steps T , step size η , batch size m

1. Initialize θ_0
2. For $t = 1, 2, \dots, T$
3. Select random $(x_1, y_1), \dots, (x_m, y_m)$
4. Calculate $g_t = \sum_{i=1}^m \nabla_{\theta} \ell(\theta_{t-1}; x_i, y_i)$
5. Update $\theta_t \leftarrow \theta_{t-1} - \eta g_t$
6. Return θ_T

How do we get the gradient?

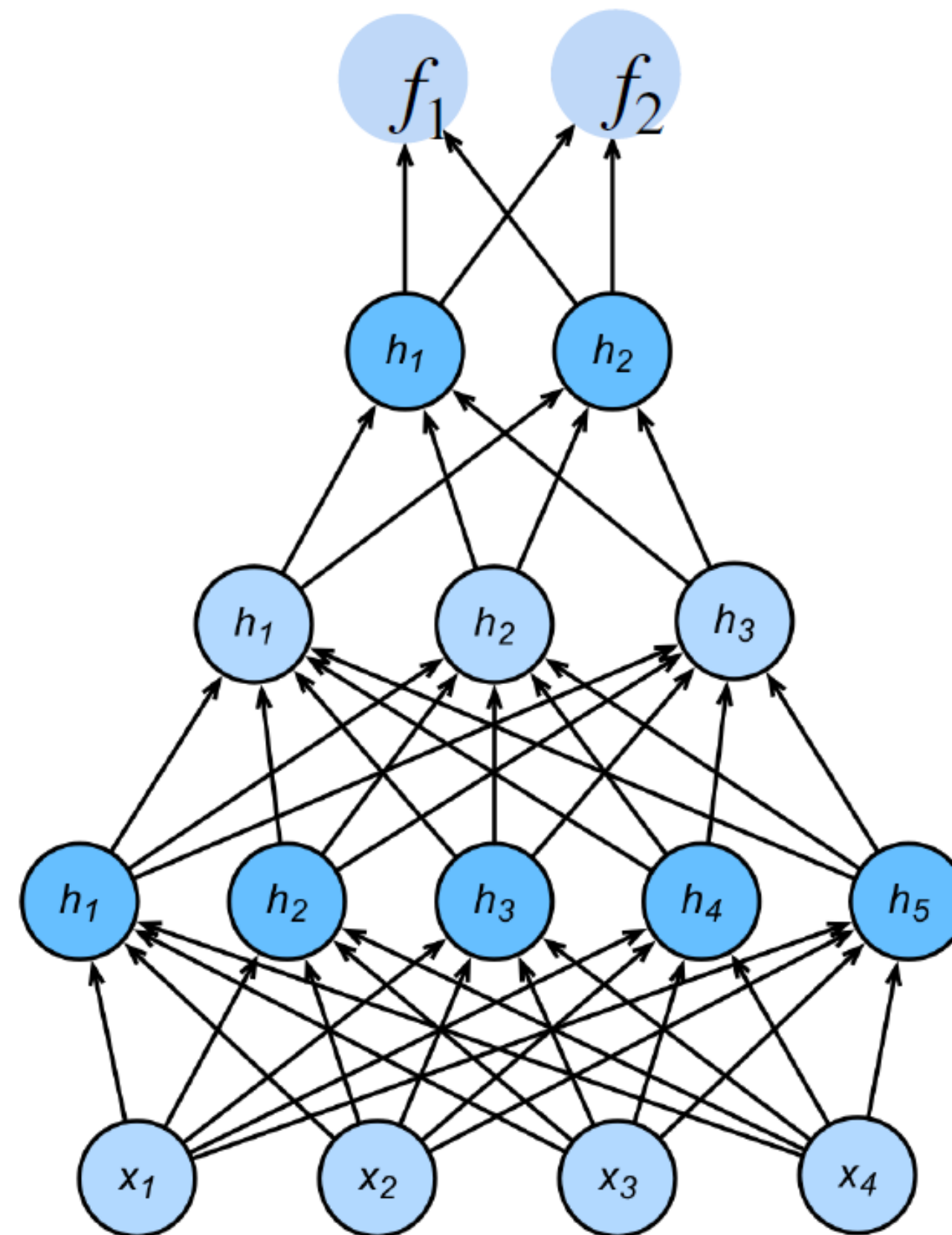


[Gao and Li et al., 2018]

Backpropagation: An Efficient Algorithm for Gradients in Neural Networks

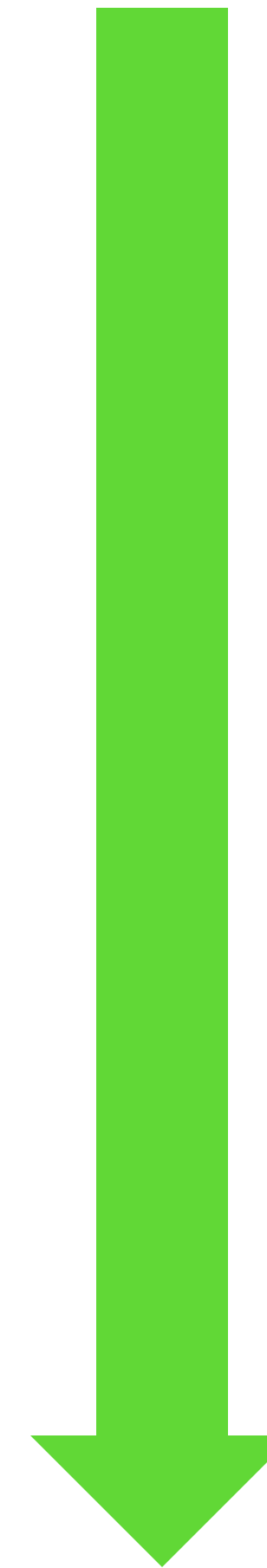
Forward pass:

Start with input layer,
compute all hidden
nodes and outputs
layer-by-layer



Backward pass:

Start with output layer,
compute all partial
derivatives
layer-by-layer



Derivatives of functions of single variables

- For many functions $f(x)$ we have simple rules to find the derivative.
- If $f(x) = cx^2$ then $\frac{df}{dx} = 2cx$
 - Or more generally, if $f(x) = cx^p$ then $\frac{df}{dx} = cpx^{p-1}$
- If $f(x) = \log x$ then $\frac{df}{dx} = \frac{1}{x}$
- If $f(x) = c$ then $\frac{df}{dx} = 0$

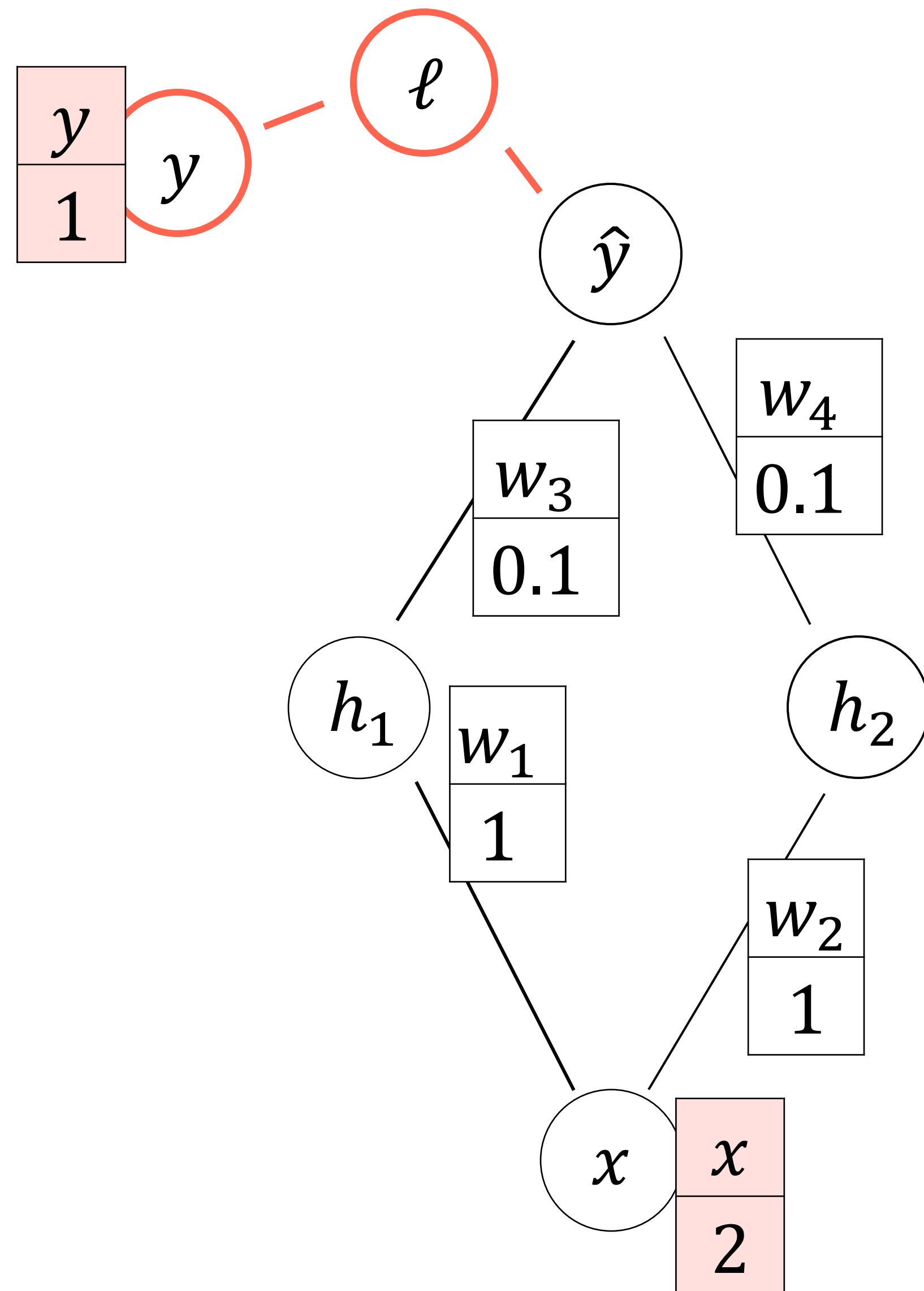
More Complex Functions

- Derivation rules can be applied hierarchically to find derivatives of more complex functions.
- **Sum rule:** If $f(x) = h(x) + g(x)$ then $\frac{df}{dx} = \frac{dh}{dx} + \frac{dg}{dx}$
- **Product rule:** If $f(x) = h(x) \cdot g(x)$ then $\frac{df}{dx} = h(x) \frac{dg}{dx} + g(x) \frac{dh}{dx}$
- **Chain rule:** If $f(x) = h(g(x))$ then $\frac{df}{dx} = \frac{dh}{dg} \frac{dg}{dx}$
- **More complex chain rule:** if $f(x) = f_1(f_2(\cdots f_n(x) \cdots))$ then $\frac{df}{dx} = \frac{df_1}{df_2} \frac{df_2}{df_3} \cdots \frac{df_n}{dx}$

Computing Partial Derivatives

- Rules from single-variable calculus directly extend to multivariate calculus with one small change.
- When computing $\frac{\partial f}{\partial x_i}$, treat all other variables as constants.
- Examples:
 - If $f(x_1, x_2) = \log x_1 + \log x_2$ then $\frac{\partial f}{\partial x_1} = \frac{1}{x_1}$
 - If $f(x_1, x_2) = x_1(x_1 + x_2)$ then $\frac{\partial f}{\partial x_2} = x_1$

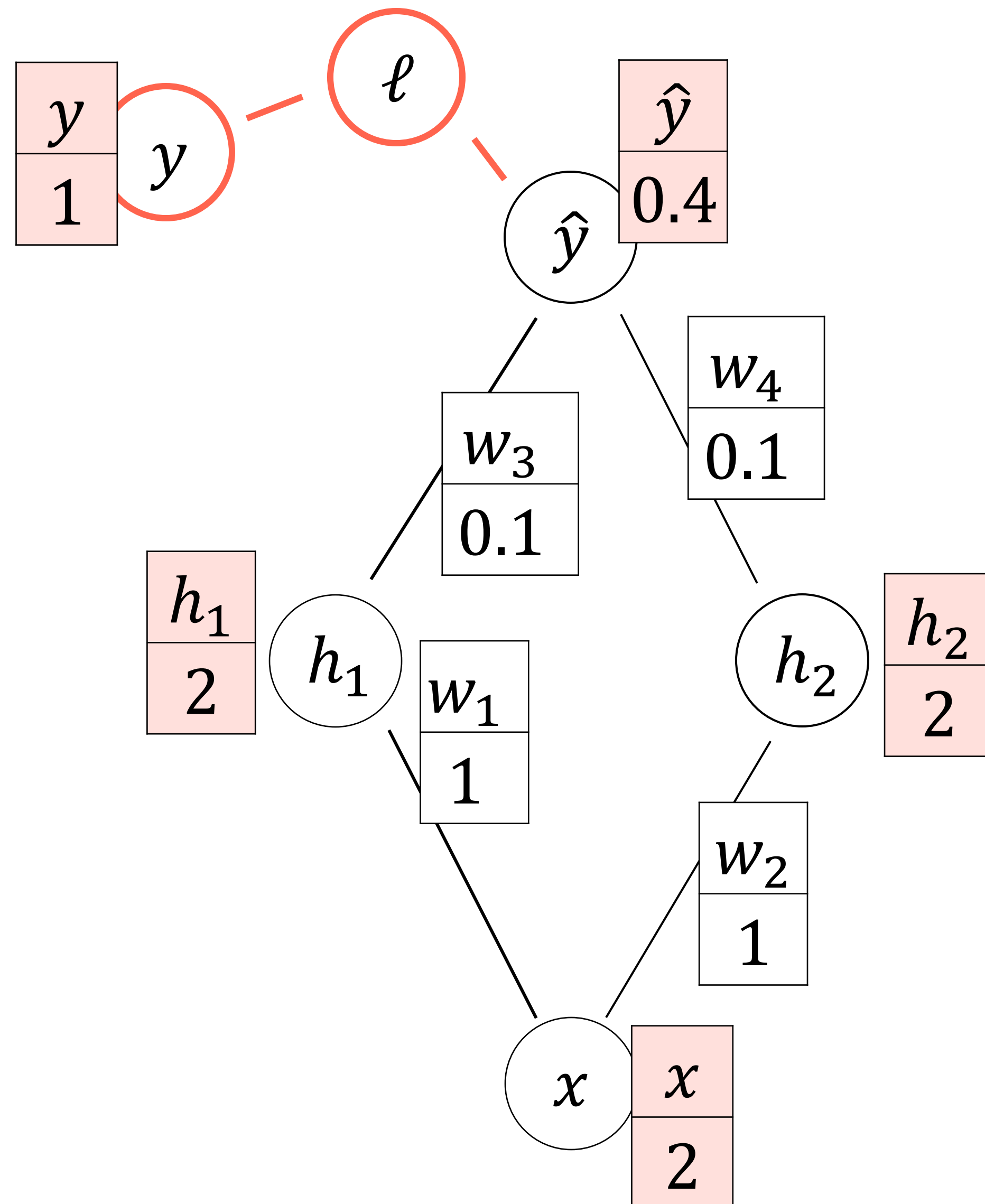
Back to Our Simple Example



Drop bias terms to simplify the picture

Add in notation for loss, and data point
 $(x, y) = (2, 1)$

Back to Our Simple Example

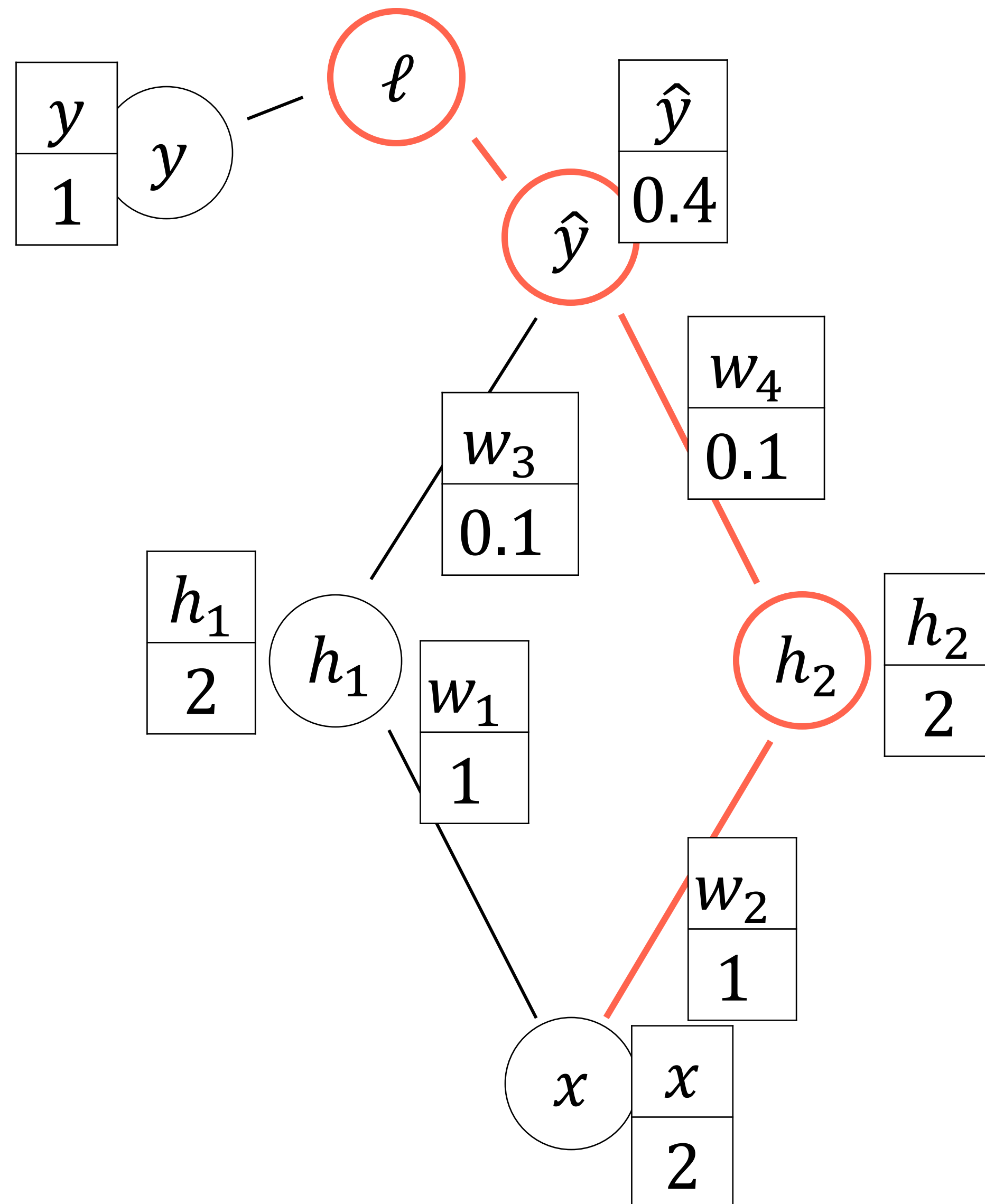


Drop bias terms to simplify the picture

Add in notation for loss, and data point
 $(x, y) = (2, 1)$

We already performed forward pass,
our algorithm stored output,
intermediate hidden node values

Back to Our Simple Example



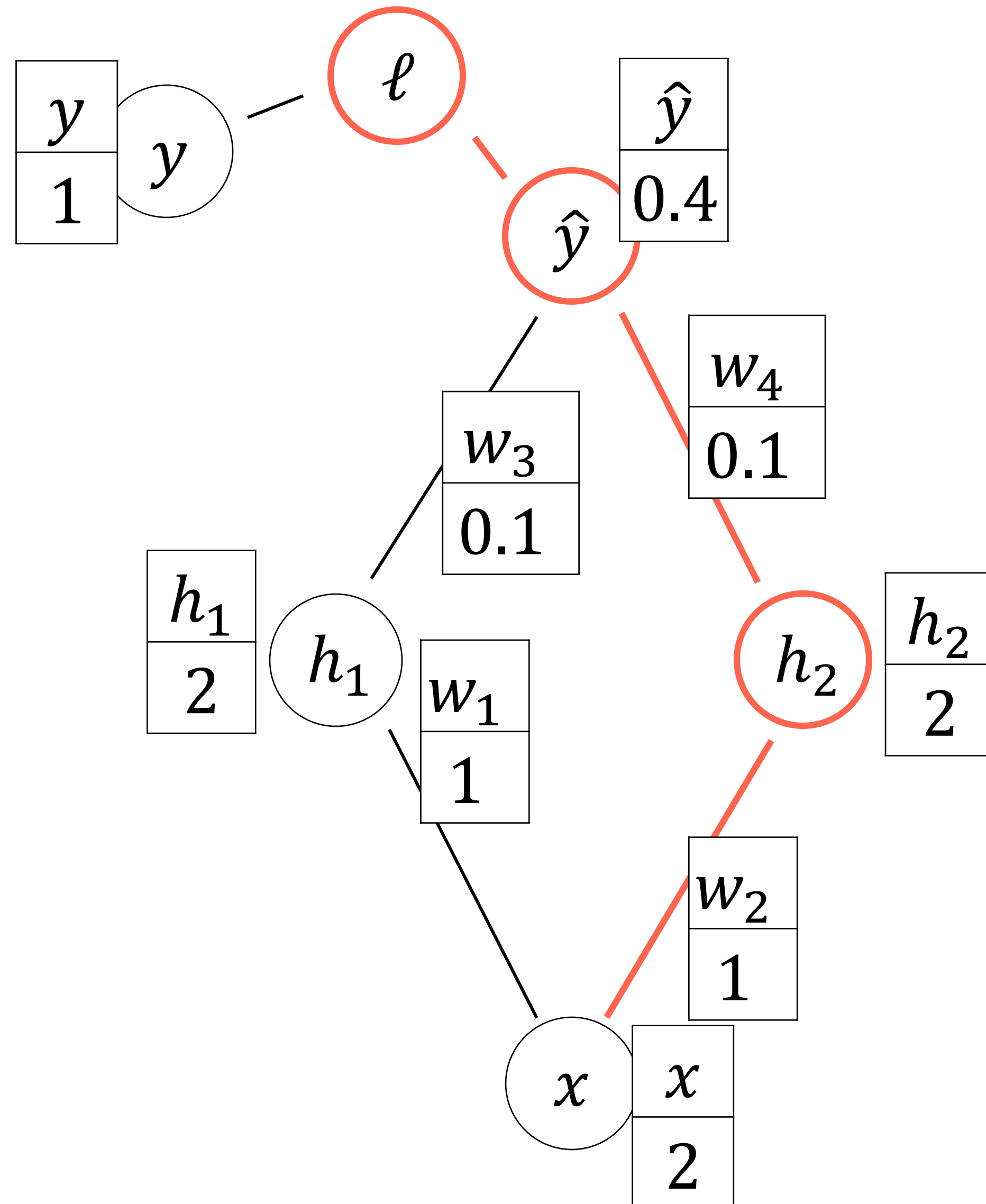
Use the squared loss: $\ell(y, \hat{y}) = (\hat{y} - y)^2$

Let's compute $\frac{\partial \ell}{\partial w_2}$

Orange lines: single path from w_2 to loss

Chain rule: $\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial h_2} \times \frac{\partial h_2}{\partial w_2}$

Back to Our Simple Example



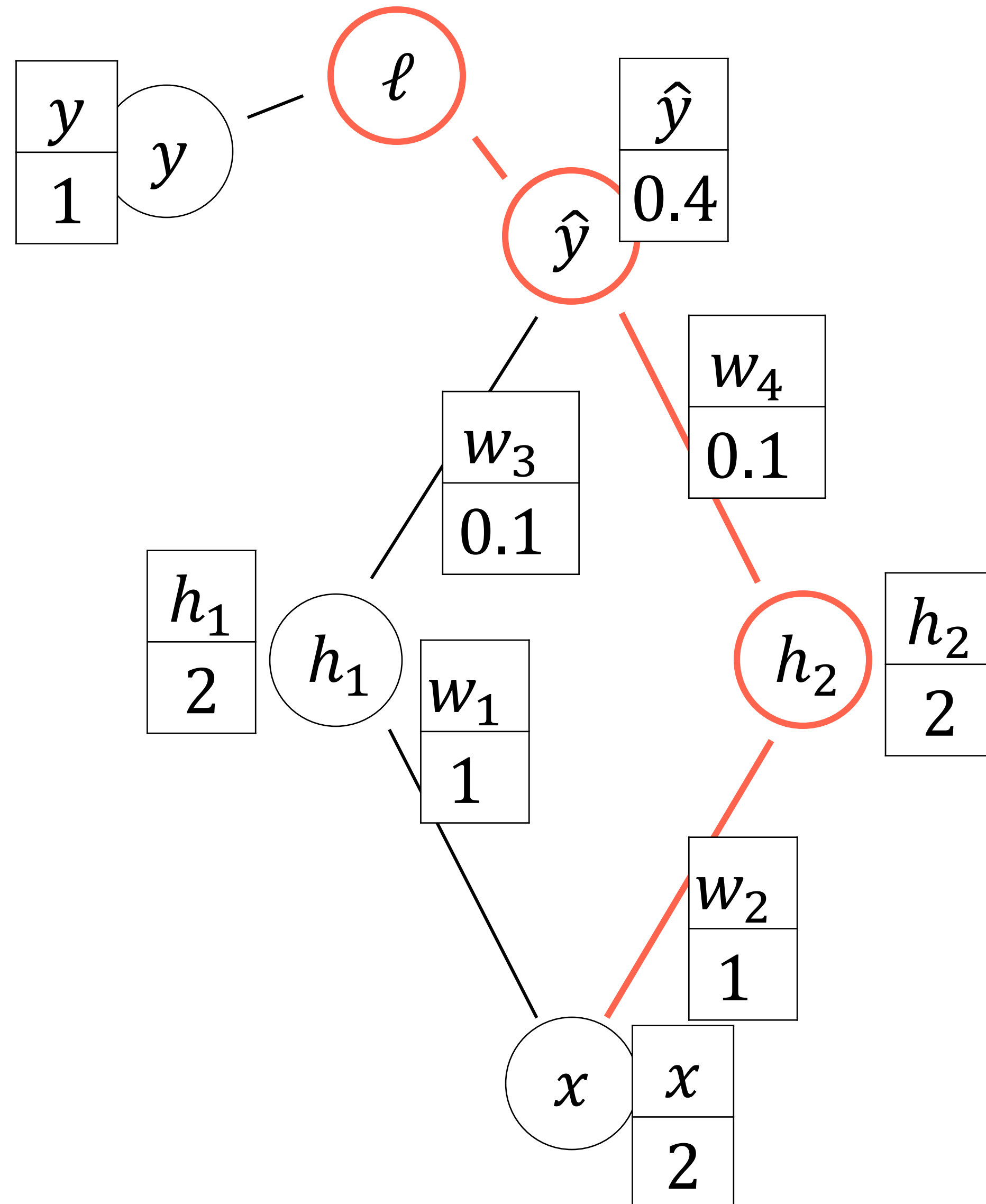
Chain rule: $\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial h_2} \times \frac{\partial h_2}{\partial w_2}$

First term:

$$\ell = (\hat{y} - y)^2$$

$$\frac{\partial \ell}{\partial \hat{y}} = 2(\hat{y} - y) = -1.2$$

Back to Our Simple Example



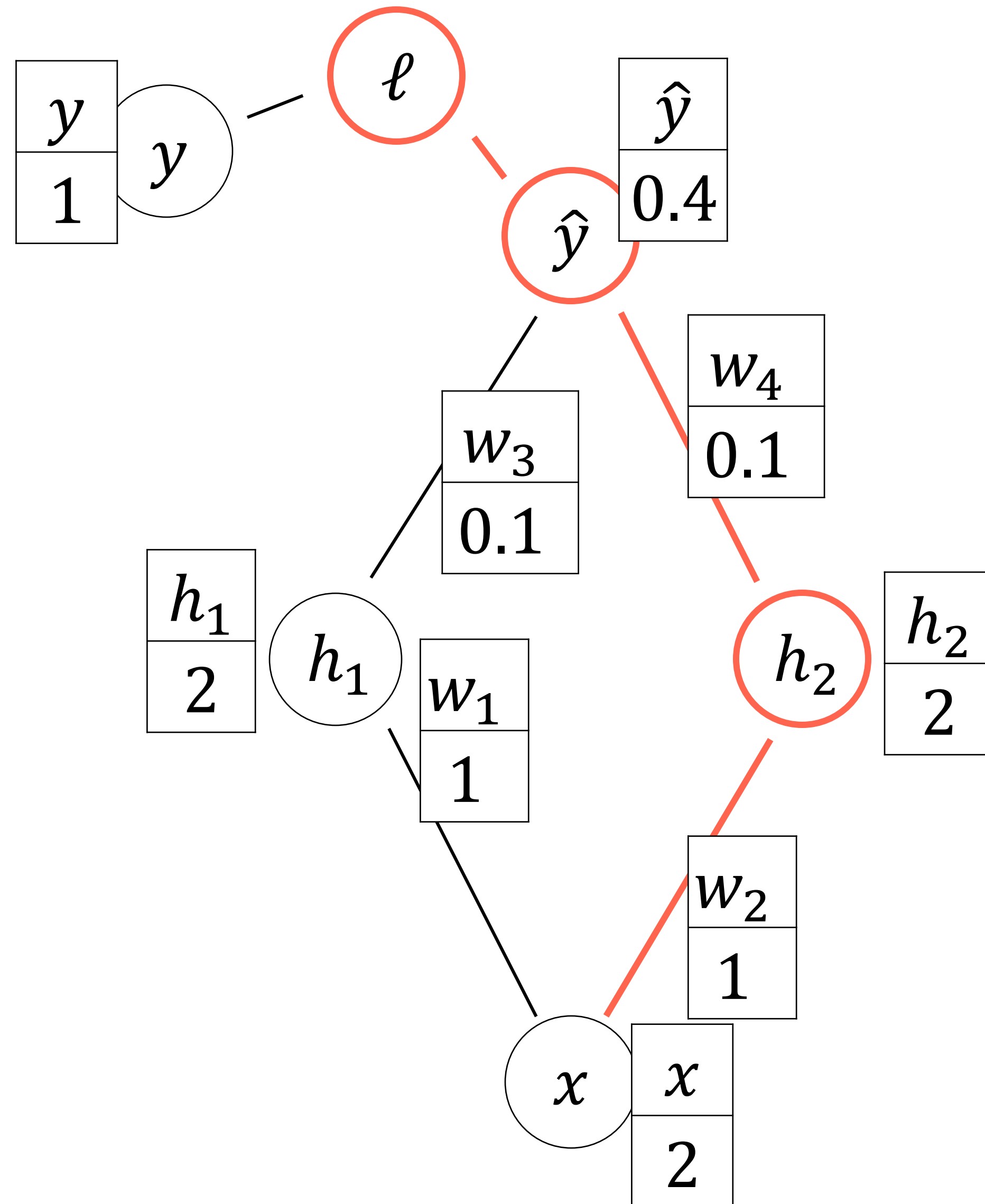
Chain rule: $\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial h_2} \times \frac{\partial h_2}{\partial w_2}$

Second term:

$$\hat{y} = w_3 h_1 + w_4 h_2$$

$$\frac{\partial \hat{y}}{\partial h_2} = w_4 = 0.1$$

Back to Our Simple Example



Chain rule: $\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial h_2} \times \frac{\partial h_2}{\partial w_2}$

Putting it all together:

$$\frac{\partial \ell}{\partial w_2} = (-1.2) \times (0.1) \times (2) = -.24$$

The Backpropagation Algorithm: Layer-by-Layer, Backwards

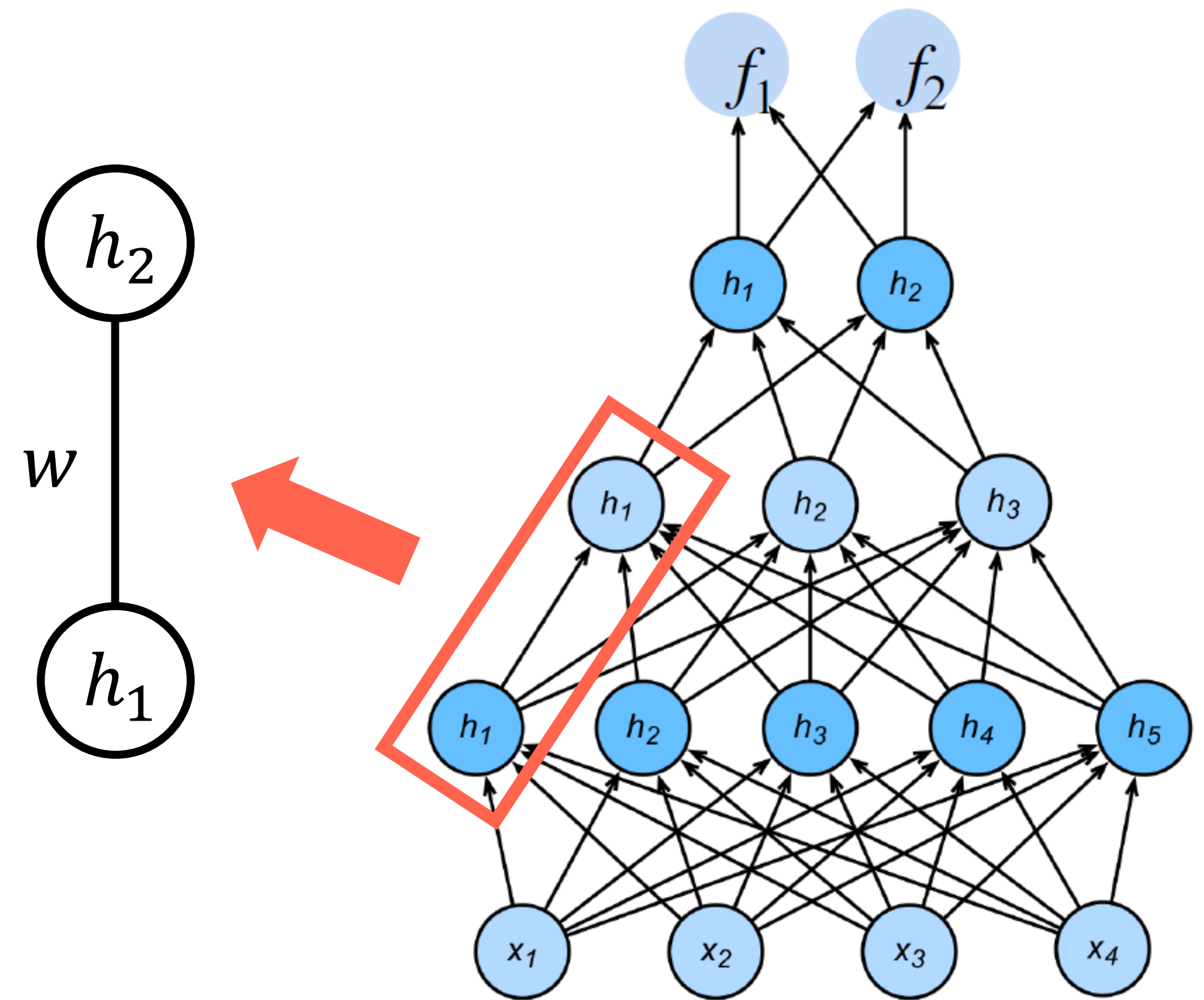
Want to find partial derivative for
weight w in middle of network

Connects from h_1 to h_2

Chain rule: $\frac{\partial \ell}{\partial w} = \frac{\partial \ell}{\partial h_2} \times \frac{\partial h_2}{\partial w}$

We already
computed this

This is simple
to compute



The spelled-out intro to neural n X

www.youtube.com/watch?v=VMj-3S1tku0

YouTube

Search

Create

micrograd_from_scratch_yay Last Checkpoint: 8 minutes ago (autosaved)

```
for child in v._prev:
    edges.add((child, v))
    build(child)
build(root)
return nodes, edges

def draw_dot(root):
    dot = Digraph(format='svg', graph_attr={'rankdir': 'LR'}) # LR = left to right
    nodes, edges = trace(root)
    for n in nodes:
        uid = str(id(n))
        # for any value in the graph, create a rectangular ('record') node for it
        dot.node(name = uid, label = "{ %s | data %.4f }" % (n.label, n.data), shape='record')
        if n._op:
            # if this value is a result of some operation, create an op node for it
            dot.node(name = uid + n._op, label = n._op)
            # and connect this node to it
            dot.edge(uid + n._op, uid)

    for n1, n2 in edges:
        # connect n1 to the op node of n2
        dot.edge(str(id(n1)), str(id(n2)) + n2._op)

    return dot

In [112]: draw_dot(L)
Out[112]:
```

In []:

30:07 / 2:25:51 • starting the core Value object of micrograd and its visualization >

The spelled-out intro to neural networks and backpropagation: building micrograd



Andrej Karpathy
1.07M subscribers

Subscribe

59K



Share

Ask

Save

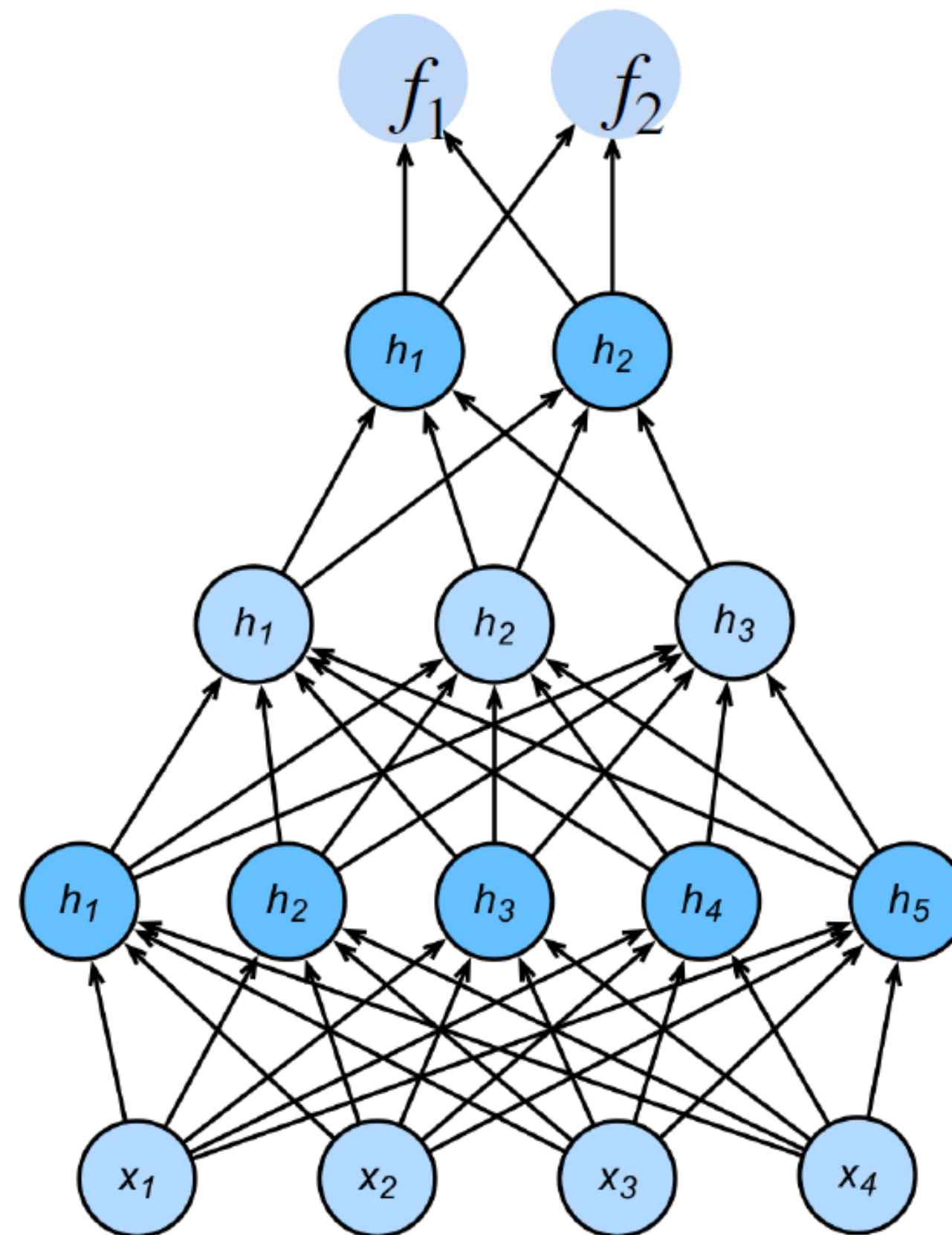


Highly recommended: <https://karpathy.ai/zero-to-hero.html>

Backpropagation: An Efficient Algorithm for Gradients in Neural Networks

Forward pass:

Start with input layer,
compute all hidden
nodes and outputs
layer-by-layer

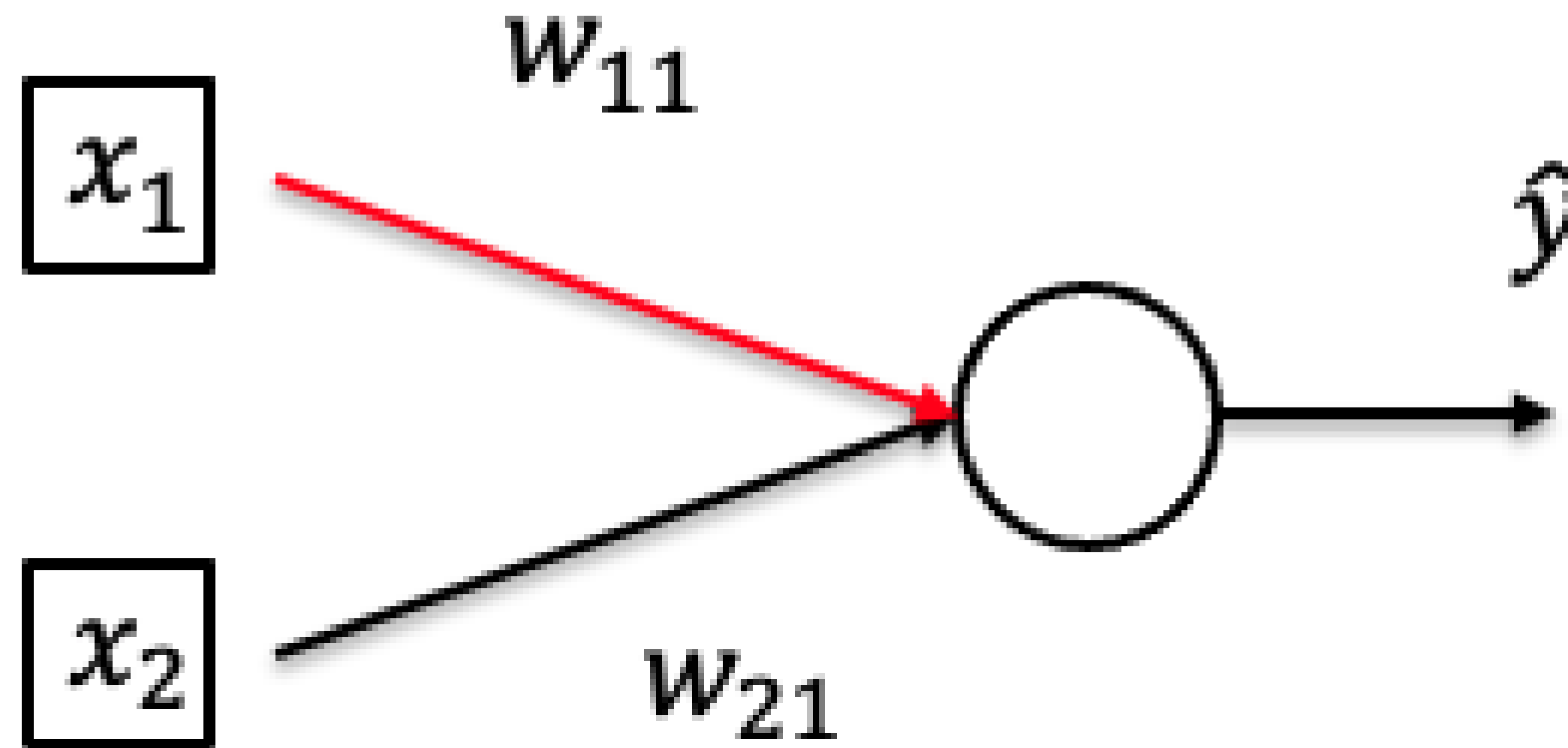


Backward pass:

Start with output layer,
compute all partial
derivatives
layer-by-layer

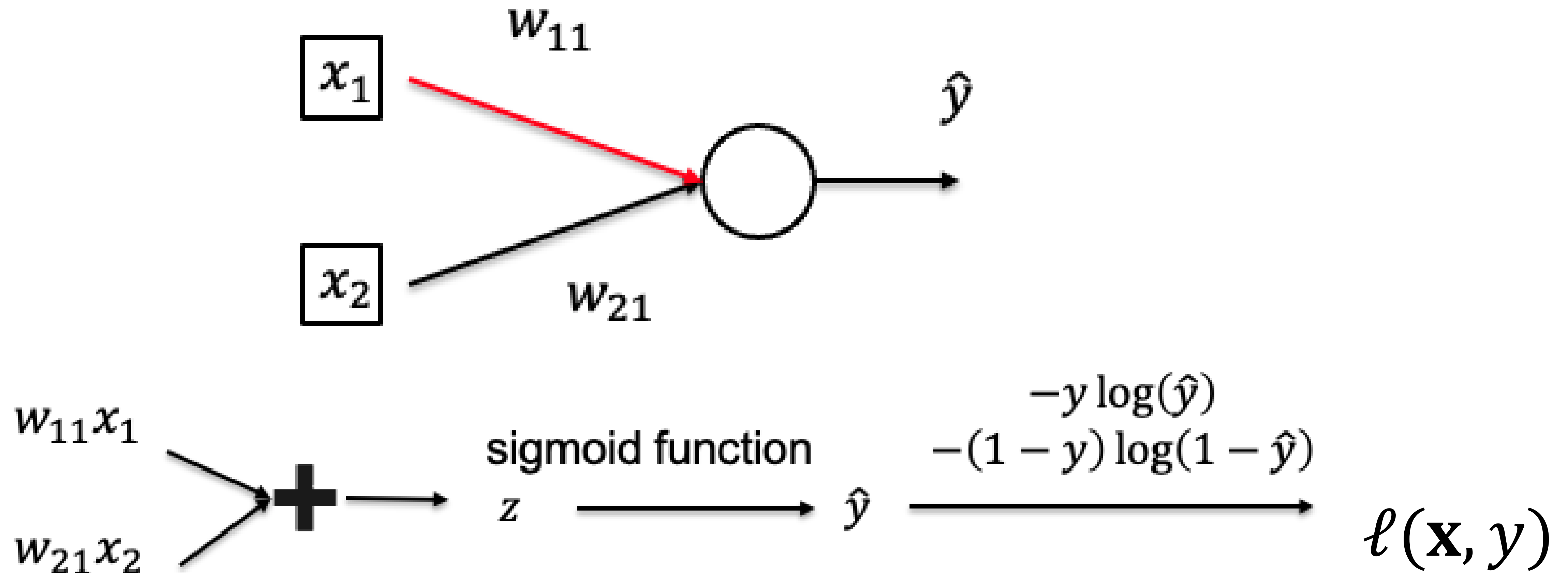
Bonus: More Backprop Math

Calculate Gradient (on one data point)



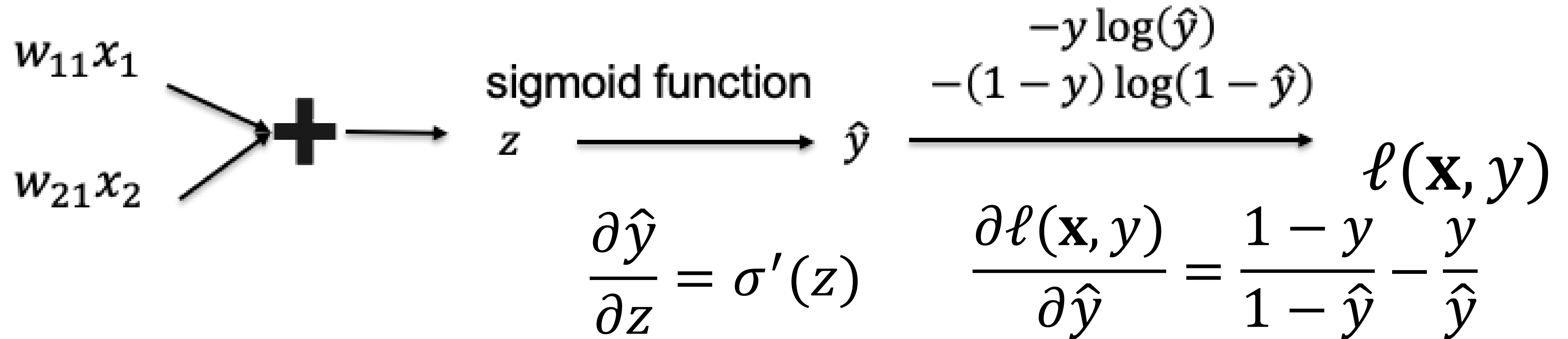
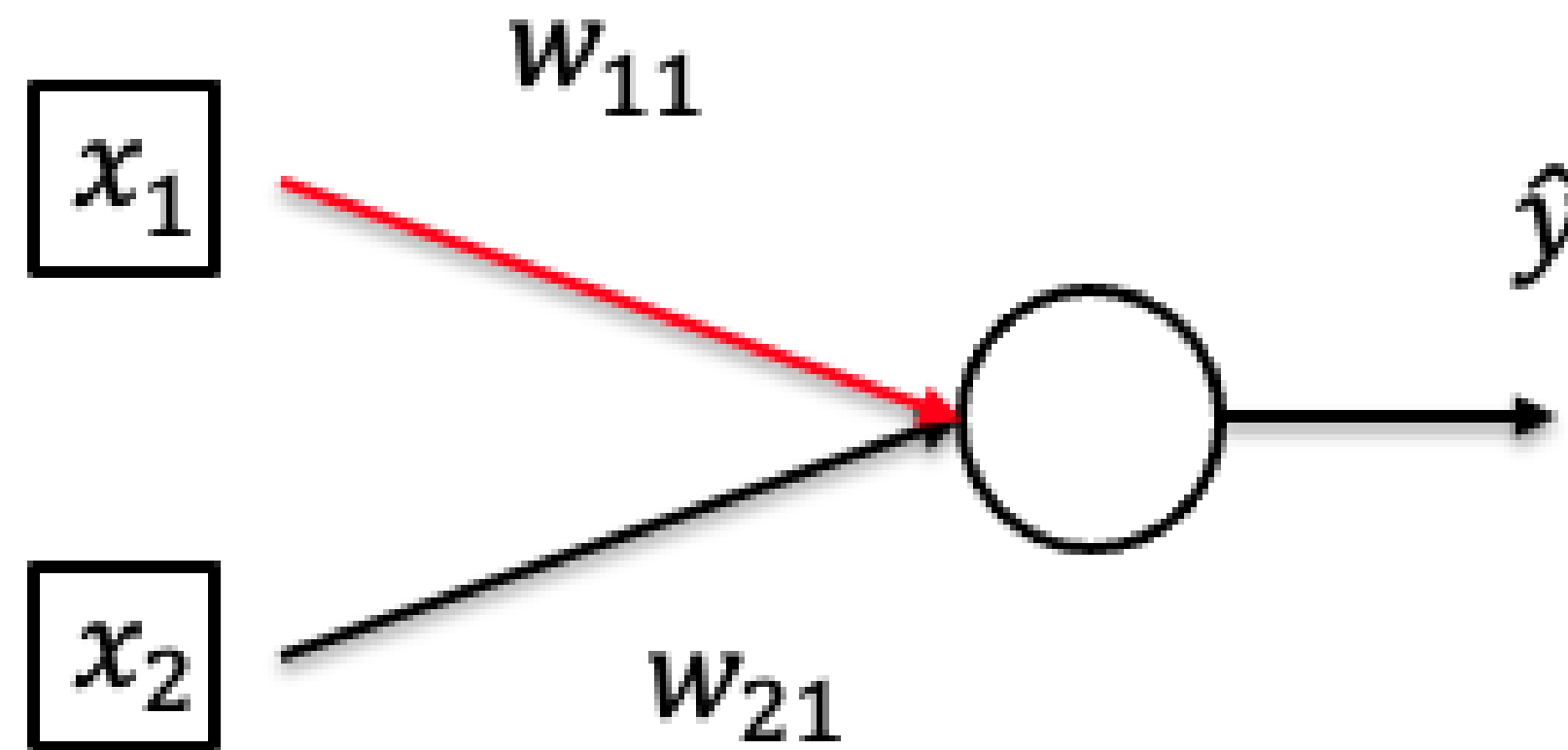
- Want to compute $\frac{\partial \ell(\mathbf{x}, y)}{\partial w_{11}}$
- Data point: $((x_1, x_2), y)$

Calculate Gradient (on one data point)



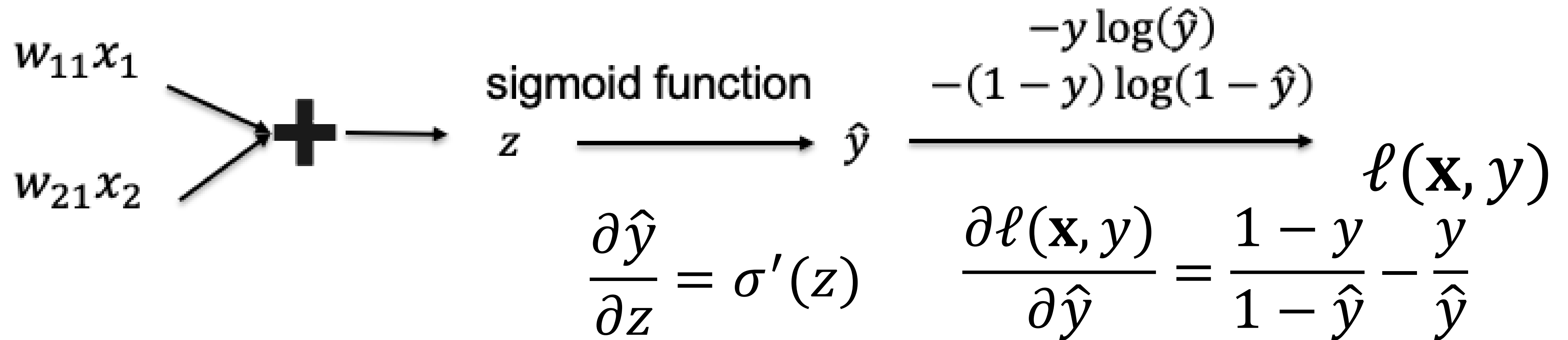
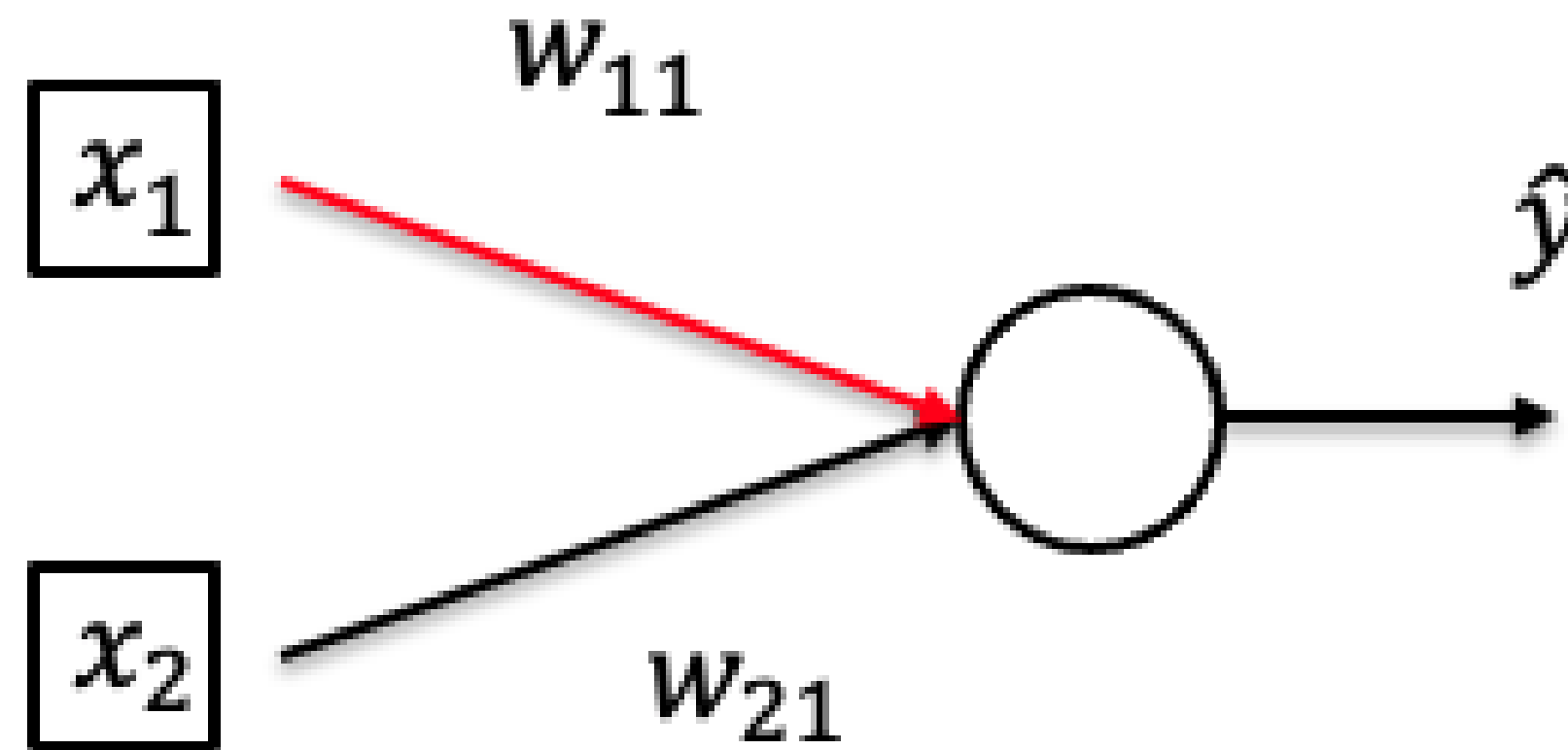
Use chain rule!

Calculate Gradient (on one data point)



- By chain rule: $\frac{\partial \ell}{\partial w_{11}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$

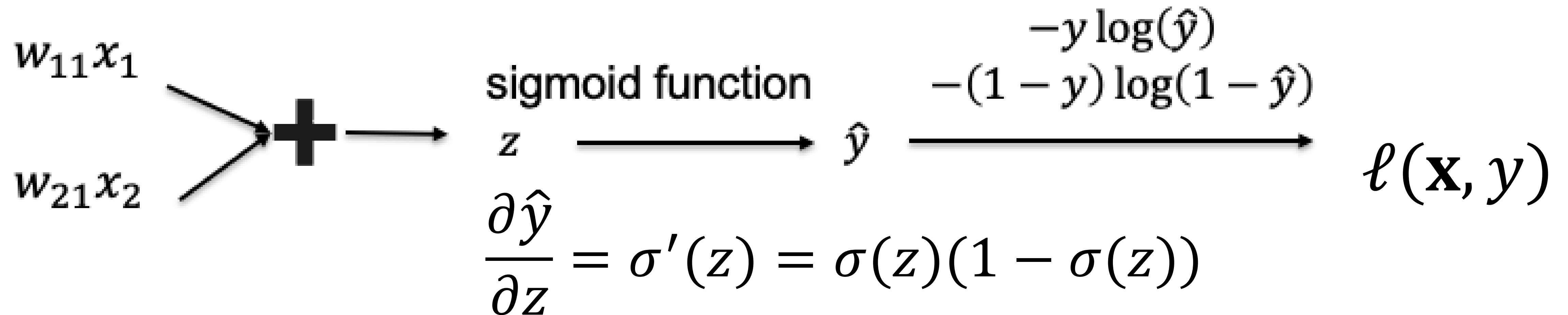
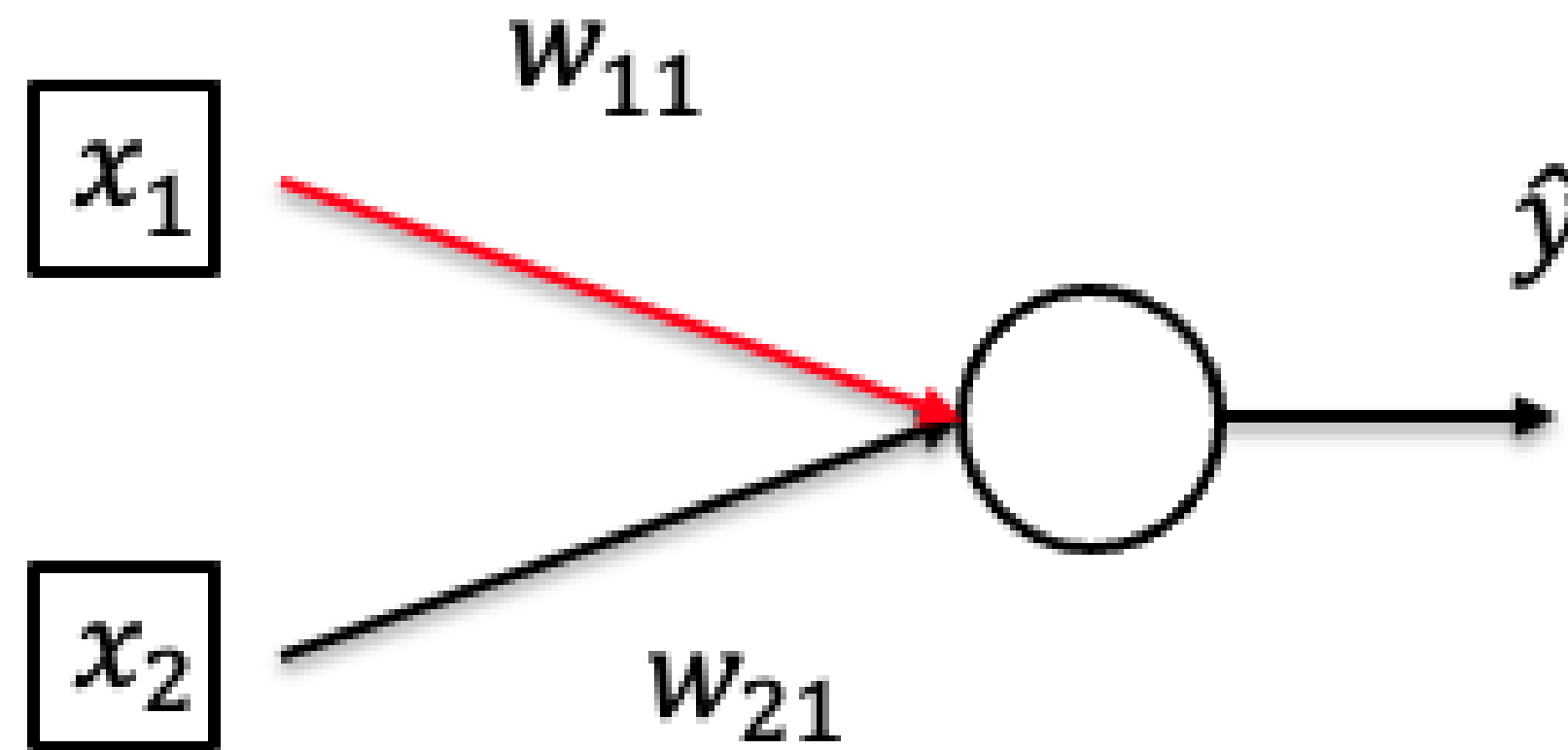
Calculate Gradient (on one data point)



- By chain rule:

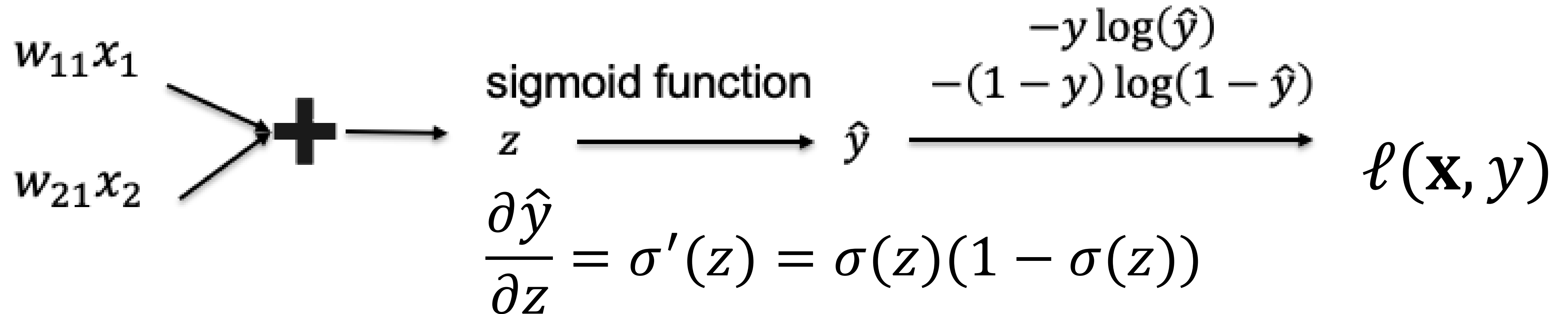
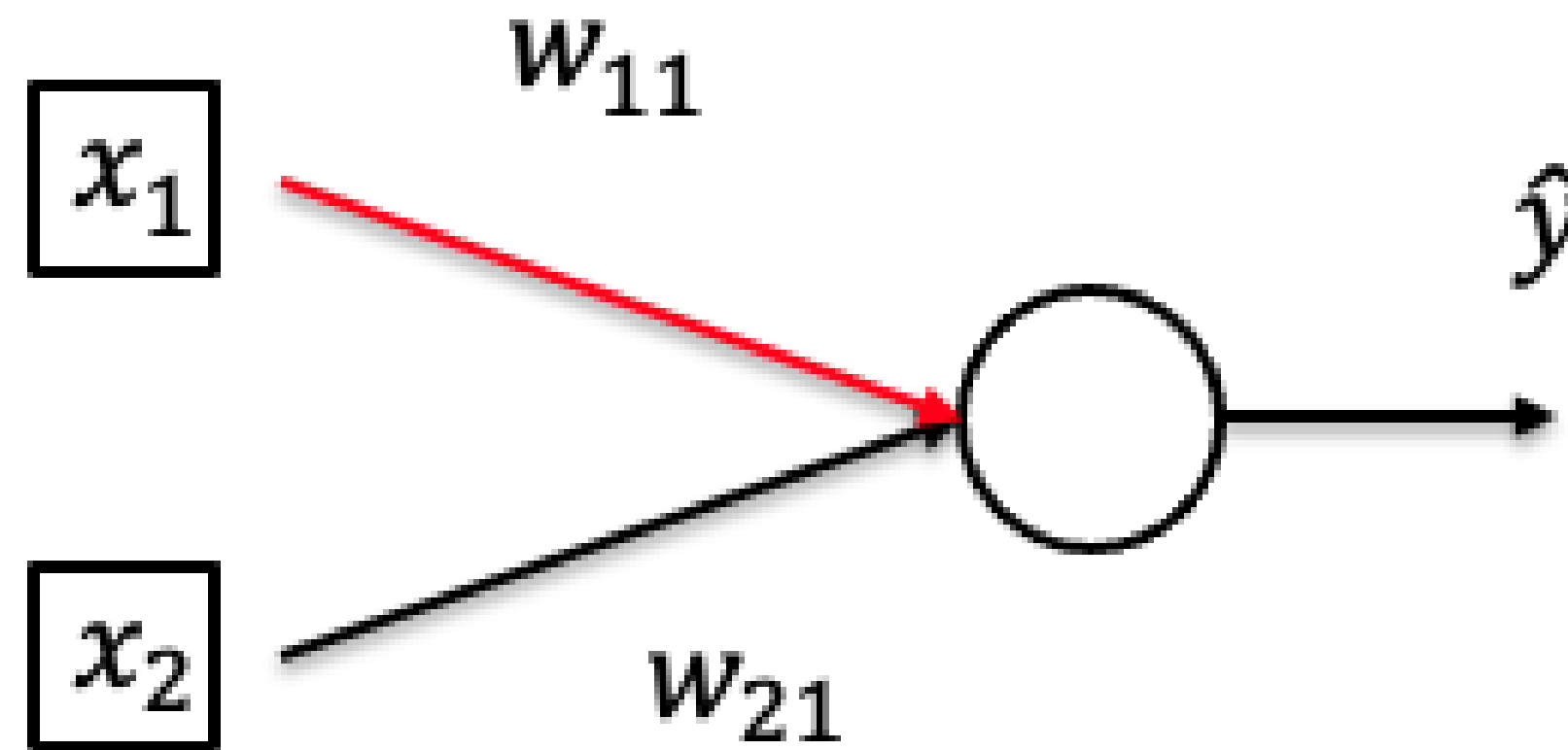
$$\frac{\partial \ell}{\partial w_{11}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} x_1$$

Calculate Gradient (on one data point)



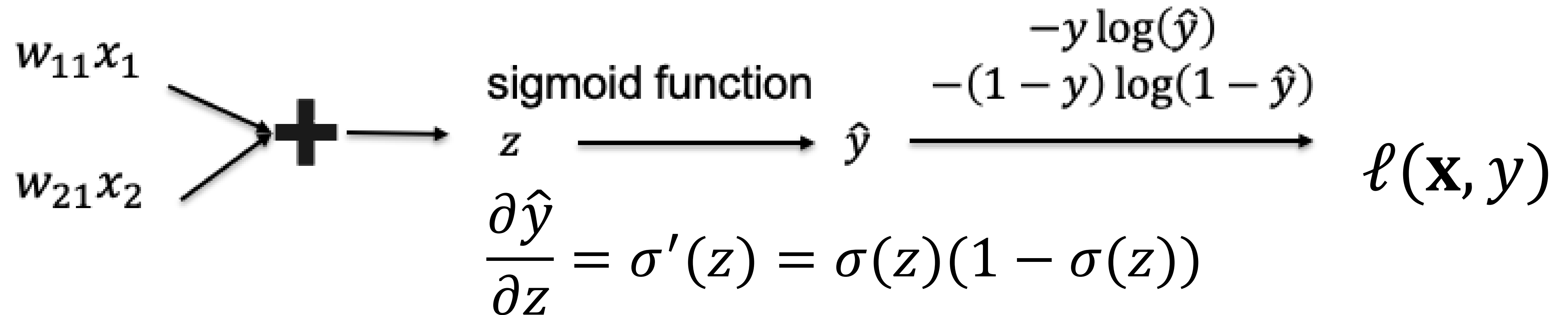
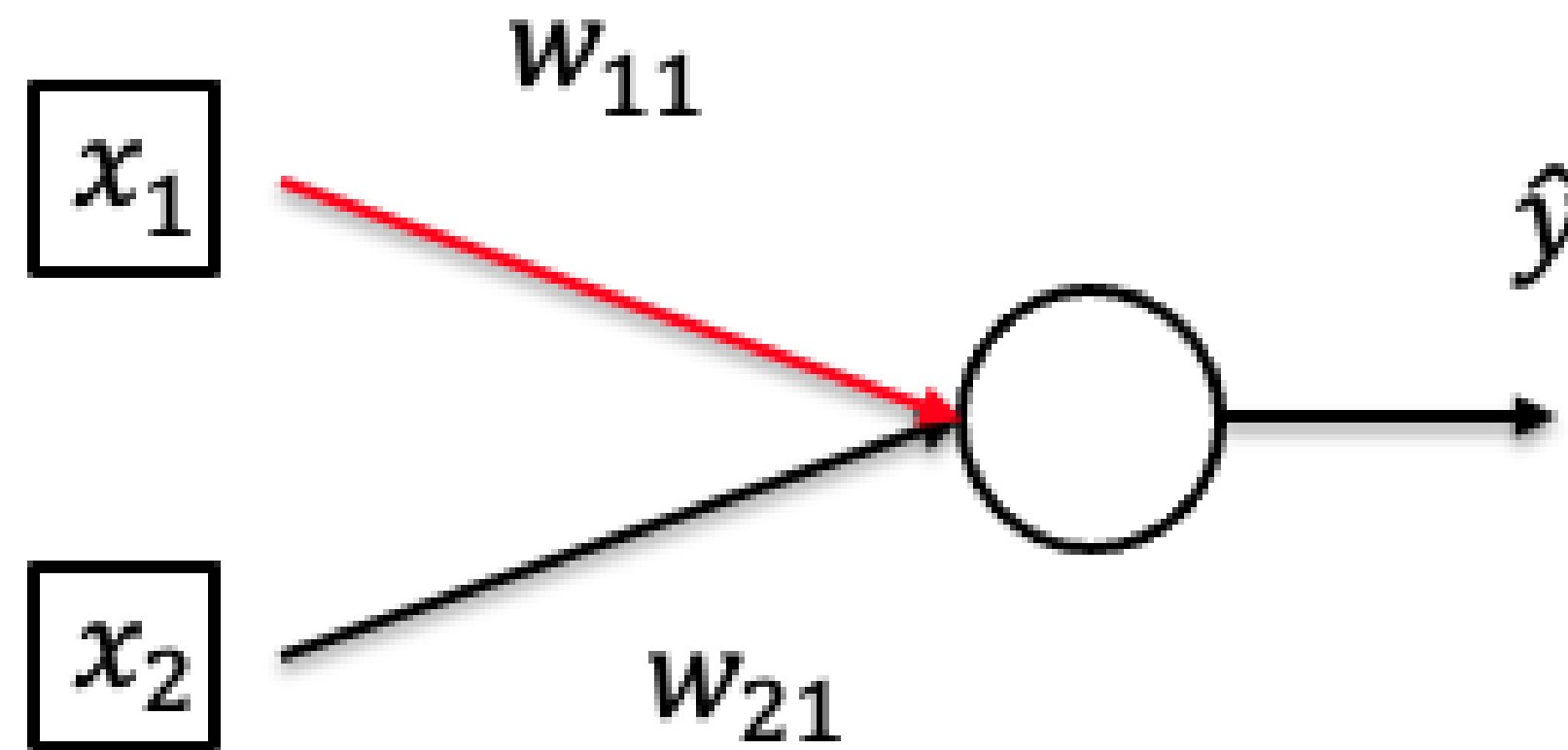
- By chain rule: $\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \hat{y}(1 - \hat{y})x_1$

Calculate Gradient (on one data point)



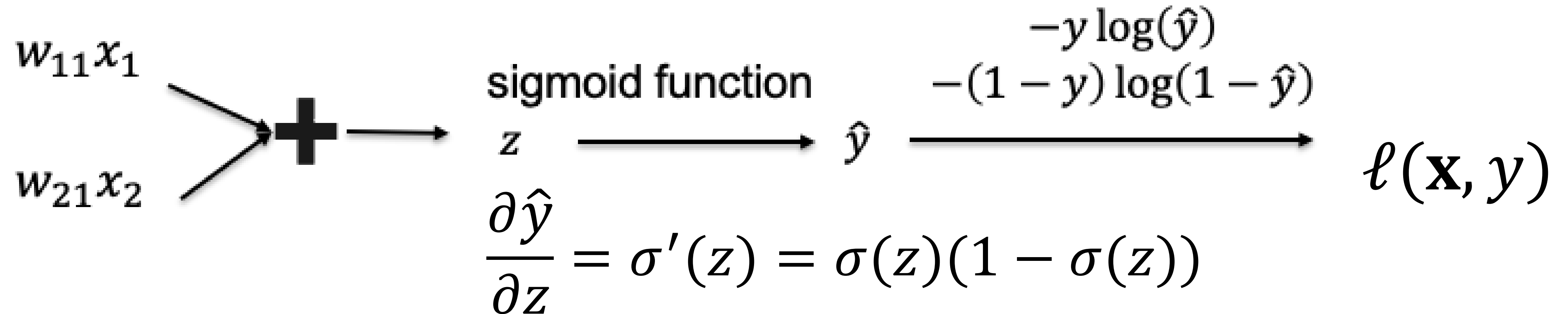
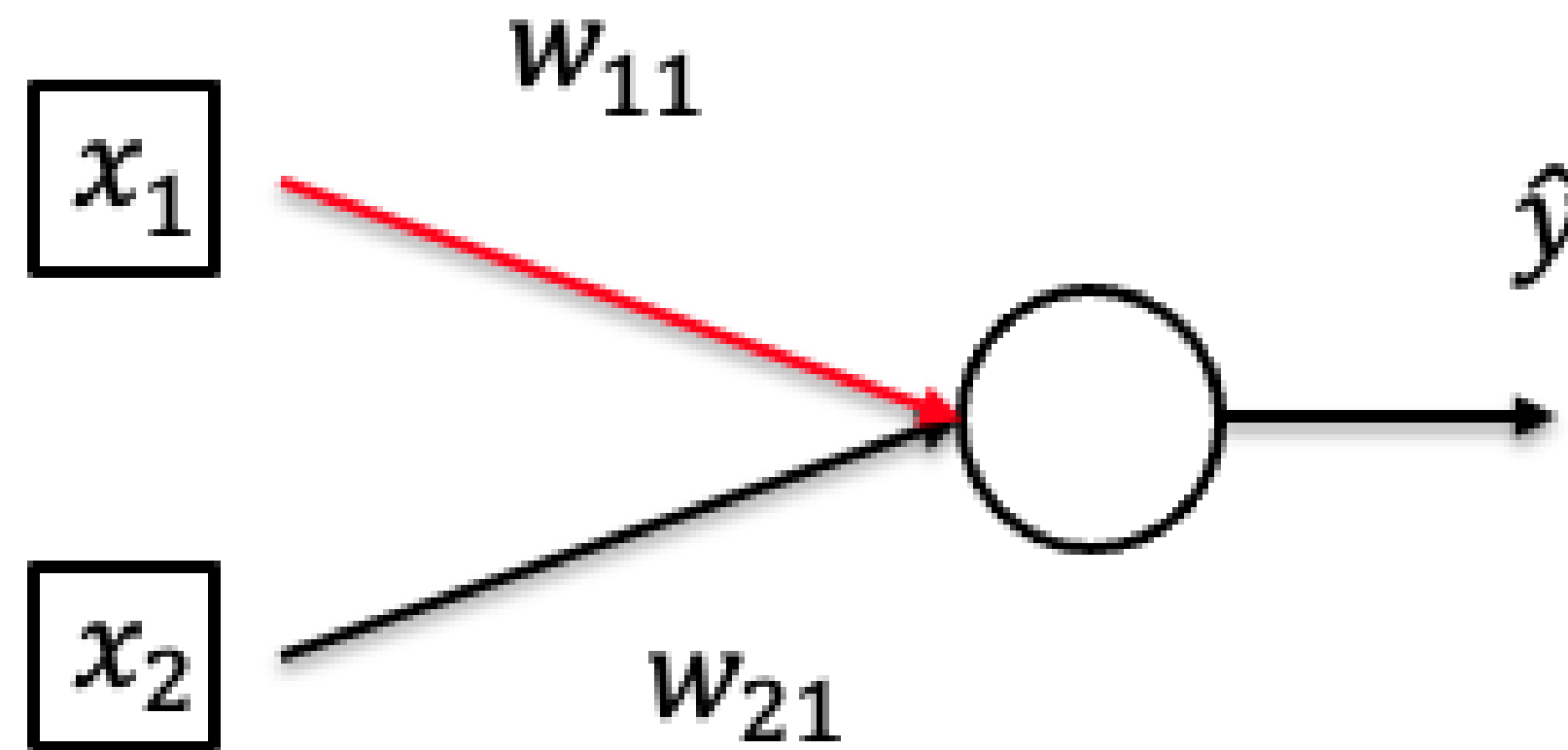
- By chain rule:
$$\frac{\partial l}{\partial w_{11}} = \left(\frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right) \hat{y} (1 - \hat{y}) x_1$$

Calculate Gradient (on one data point)



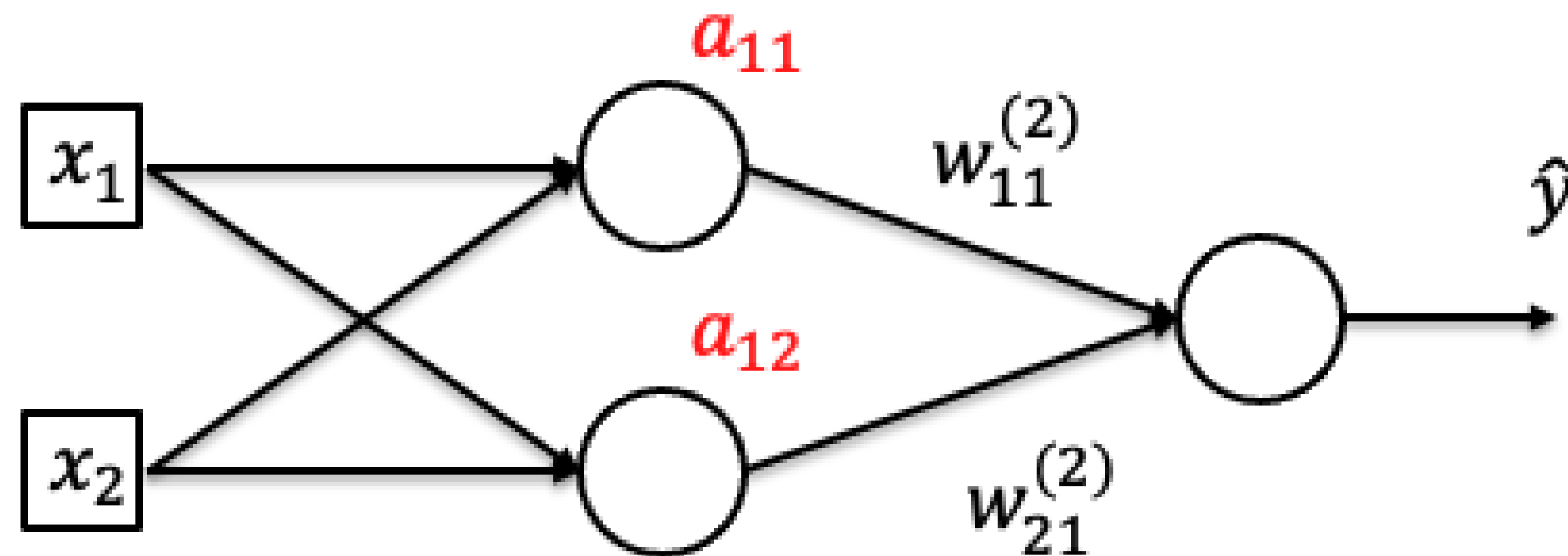
- By chain rule:
$$\frac{\partial l}{\partial w_{11}} = (\hat{y} - y)x_1$$

Calculate Gradient (on one data point)

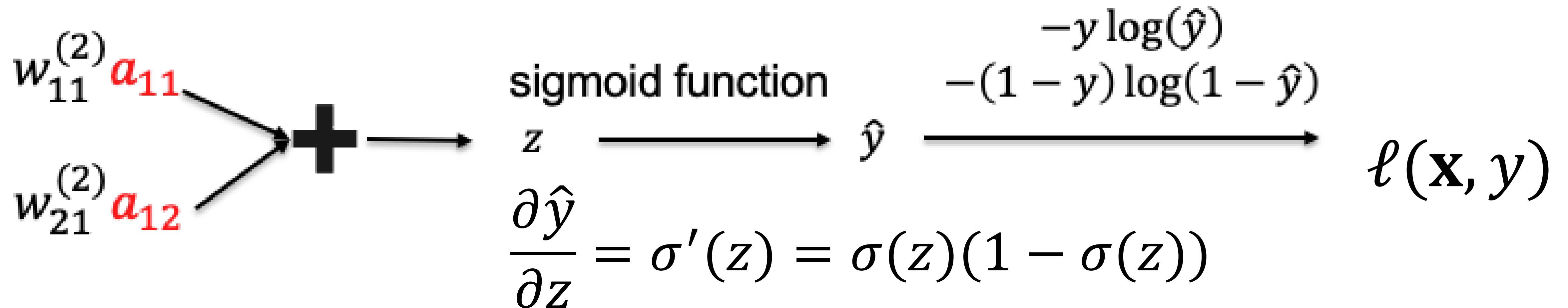


- By chain rule:
$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y)w_{11}$$

Calculate Gradient (on one data point)

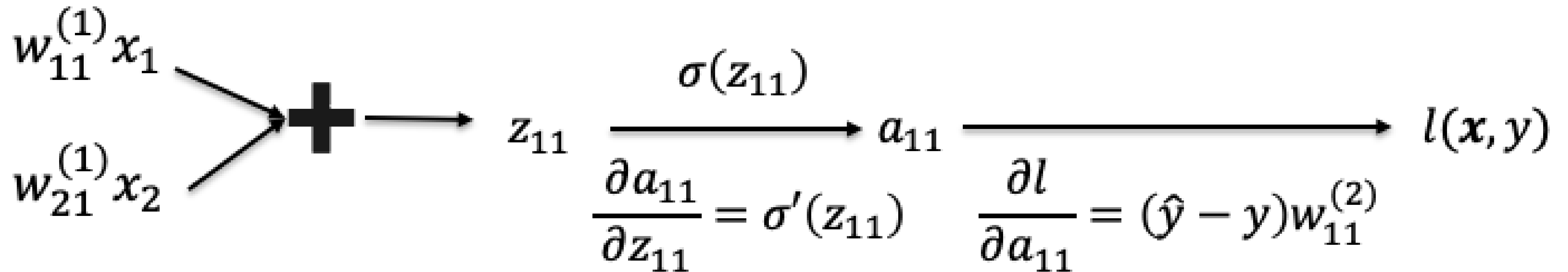
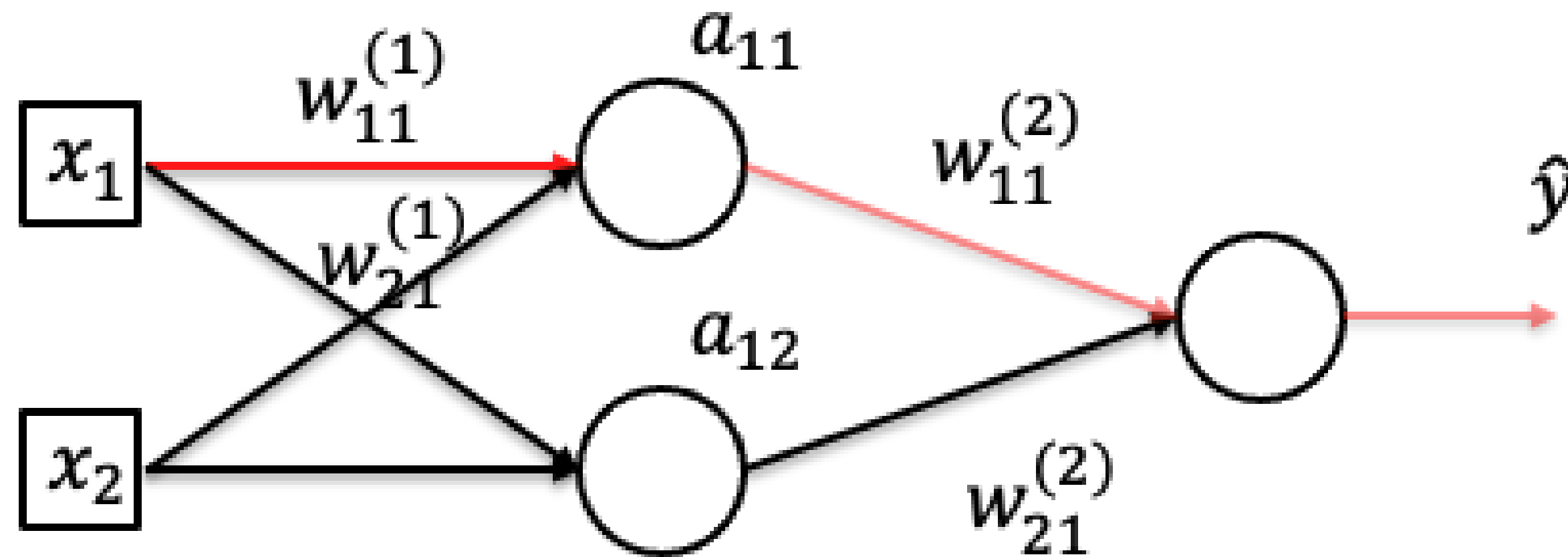


Make it deeper



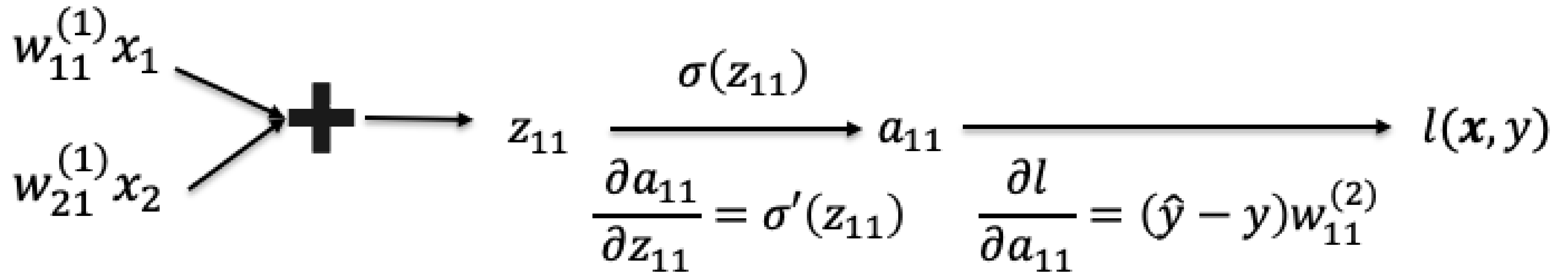
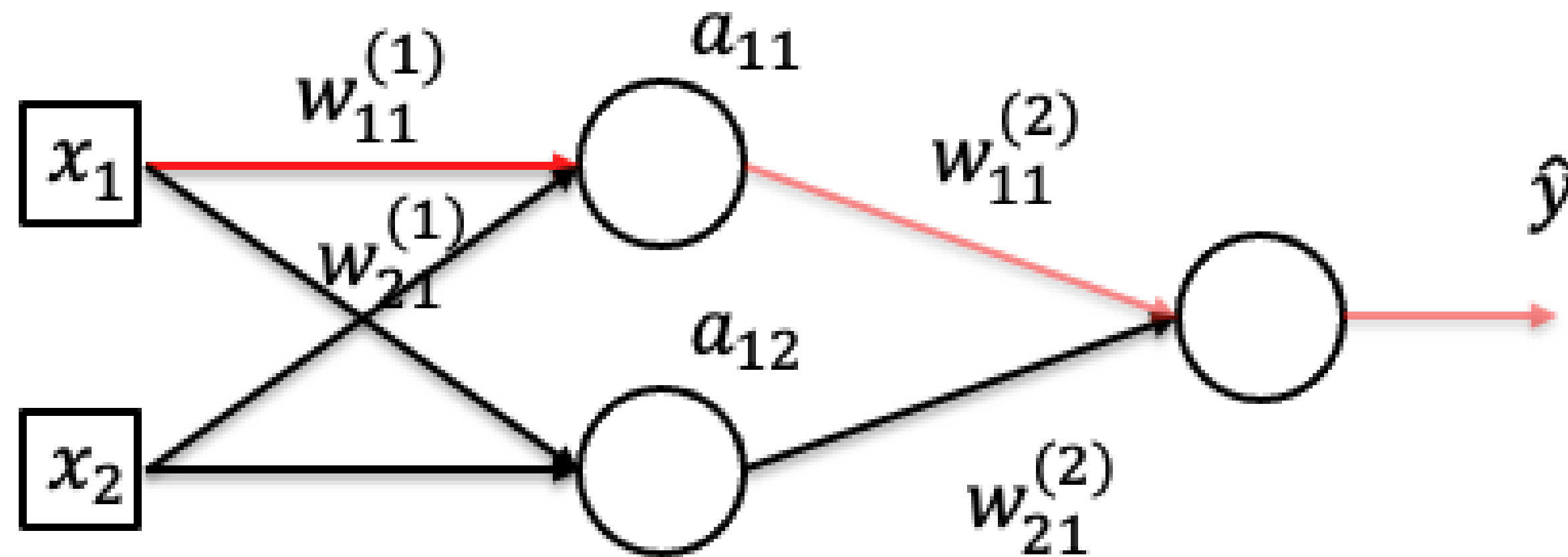
- By chain rule: $\frac{\partial \ell}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}, \frac{\partial \ell}{\partial a_{12}} = (\hat{y} - y)w_{21}^{(2)}$

Calculate Gradient (on one data point)



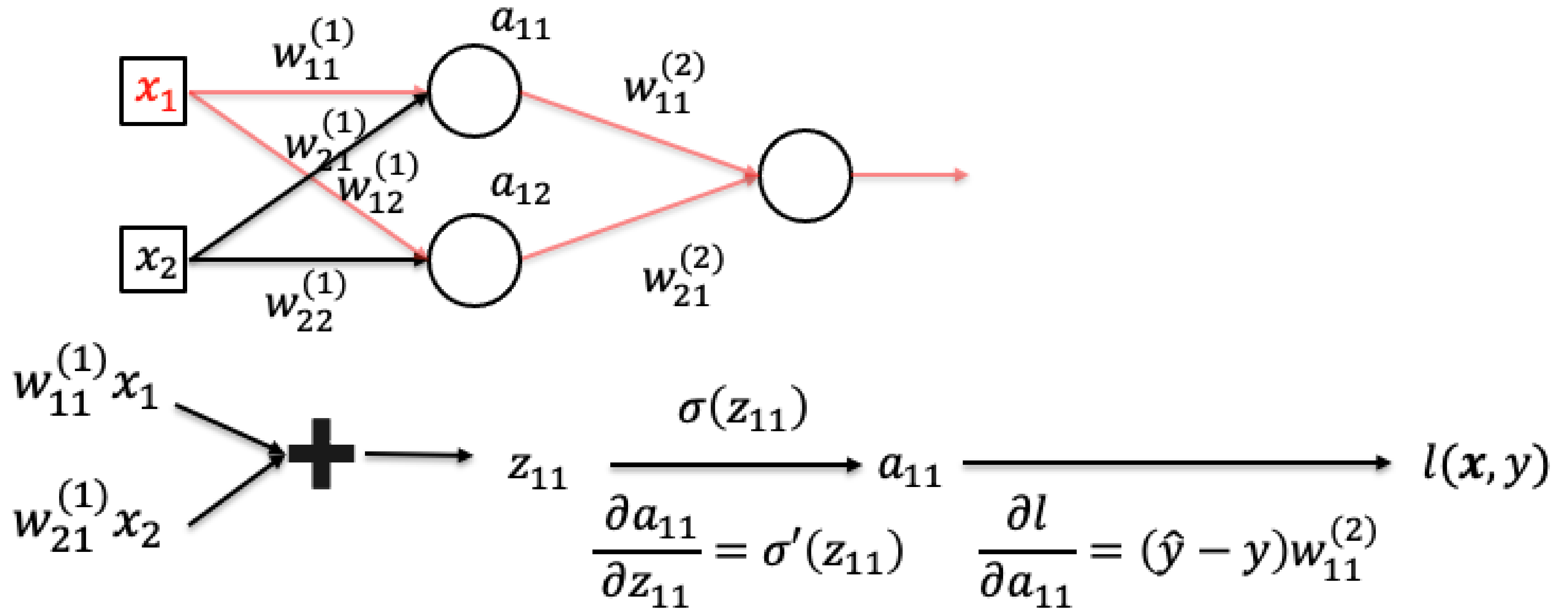
- By chain rule:
$$\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y)w_{11}^{(2)} \frac{\partial a_{11}}{\partial w_{11}^{(1)}}$$

Calculate Gradient (on one data point)



- By chain rule:
$$\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y)w_{11}^{(2)} a_{11} (1 - a_{11}) x_1$$

Calculate Gradient (on one data point)



- By chain rule:

$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial x_1} + \frac{\partial l}{\partial a_{12}} \frac{\partial a_{12}}{\partial x_1}$$