

Deep Learning

How to use Deep Learning for Marketing?



Making machines think like humans.

How (Shallow) Neural Networks work

This session has 3 learning goals

3

After this lecture you should be able to:

1. Understand the concept of neural networks.
2. Understand the basic design building blocks of neural networks.
3. Understand how to estimate neural networks.



Design and initialize a neural network

What is an artificial neural network?

5



A portrait of Ian Goodfellow, a young man with dark hair and glasses, wearing a black t-shirt. He is smiling and looking towards the camera. A speech bubble originates from his head, containing his quote.

Start by learning the basics
really well... [Andrew Ng]
always told me to work on
thorough mastery of these
basics.

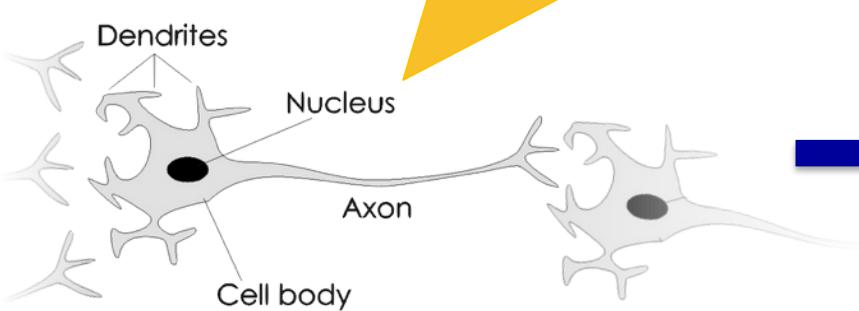
<https://hackernoon.com/interview-with-deep-learning-researcher-and-the-ganfather-dr-ian-goodfellow-cd300863ecff>

Hello
my name is
Ian Goodfellow
The GANfather

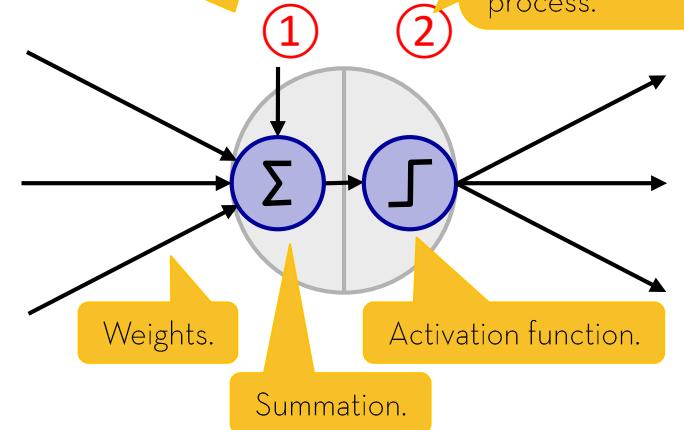
Why neural...

The basic building block of an artificial neural network is a neuron.

For the **biological neuron**, electrical signals from other neurons are conveyed to the cell body by dendrites; resultant electrical signals are sent along the axon to be distributed to other neurons.



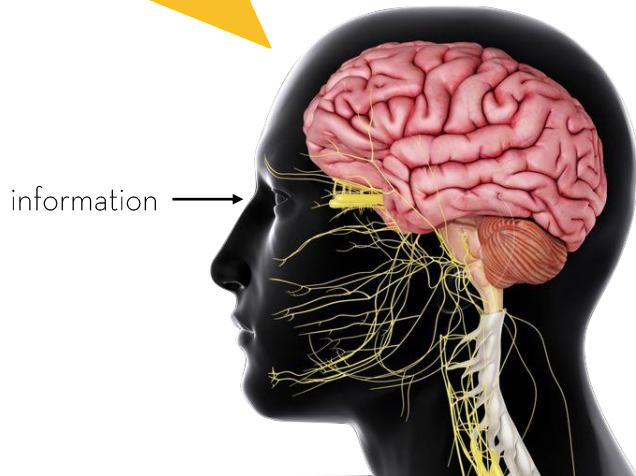
The operation of the **artificial neuron** is analogous to (though much simpler than) the operation of the biological neuron: Activations from other neurons are summed at the neuron and passed through an activation function, after which the value is sent to other neurons.



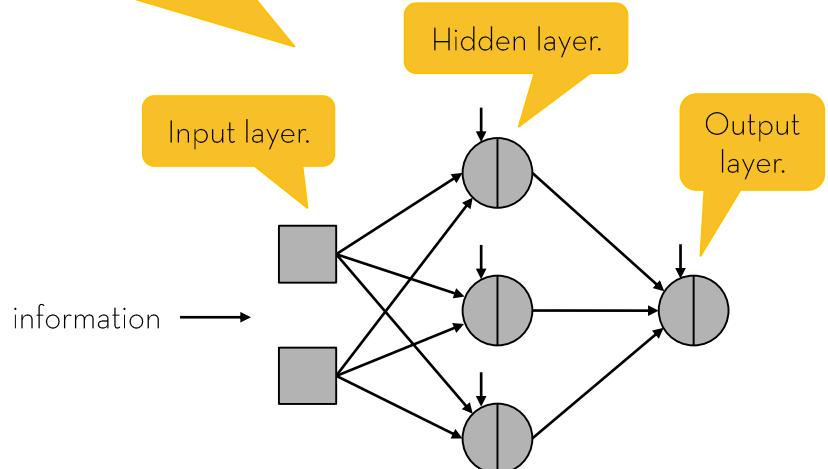
Why network...

We can stack layers of neurons together to form networks.

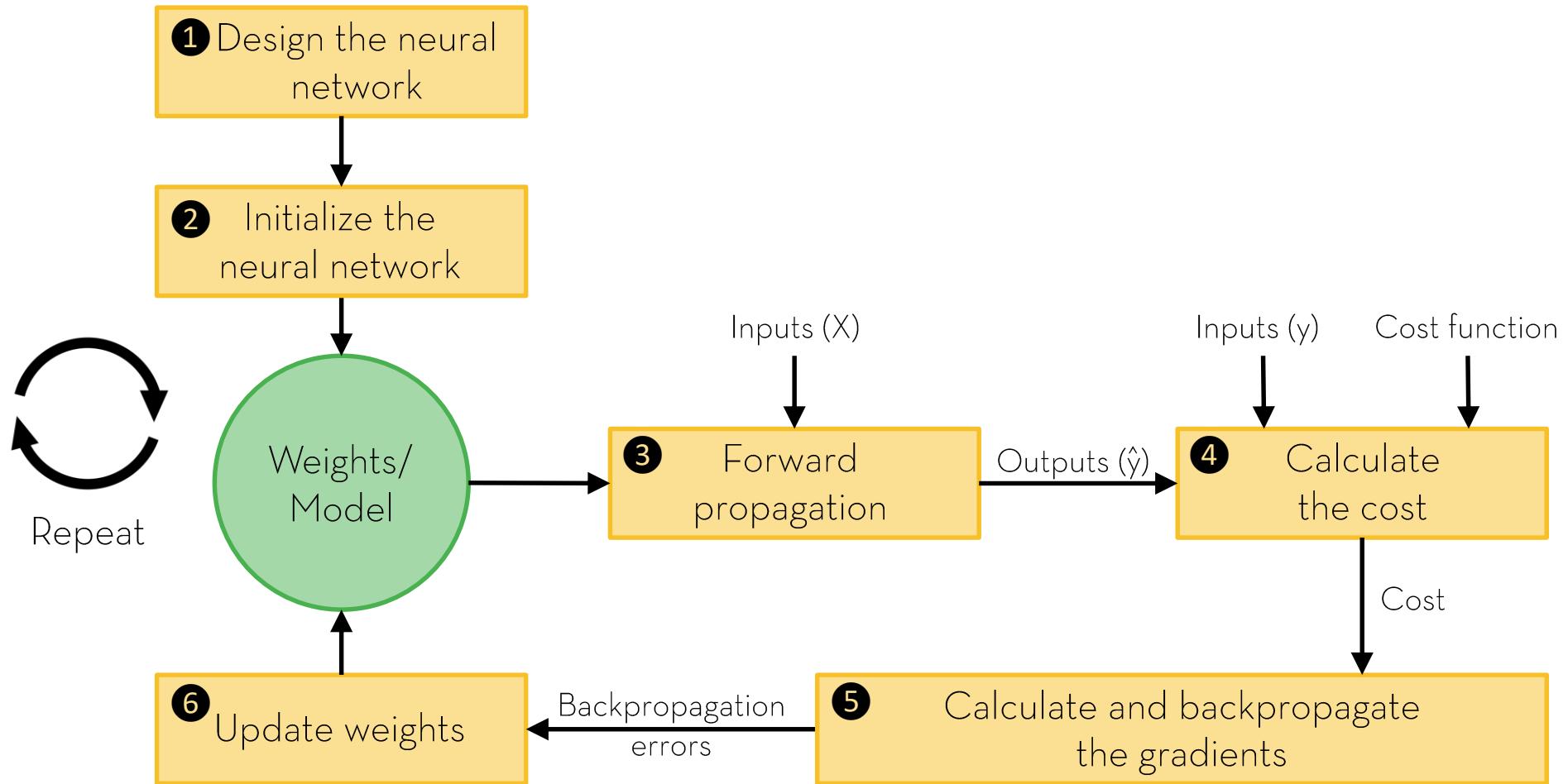
The **brain** is a collection of biological neurons that form pathways. Along these pathways flow chemical and electric signals. The brain processes information which is taken in through sensory organs.



The operation of an **artificial neural networks** is analogous to (though much simpler than) the operation of the brain: input information flows through the pathways of the network and is processed by weighted summation and activation operations of the neurons.



Steps of estimating neural networks

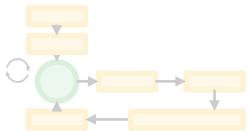


Steps of estimating neural networks

9

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



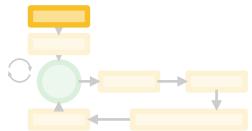
<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

Steps of estimating neural networks

10

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

Step 1: Design the neural network

- 1) Normalize the input data.
- 2) Based on the data, design the neural network:
 - Output Layer: Define number of output nodes.
 - Input Layer: Define number of input nodes.
 - Hidden Layer: Define number of hidden nodes.
 - Define the activation function on the hidden and output layer.
 - Define connections between nodes.

Step 1: Design the neural network

a) Normalize the input data

Y	X1	X2
0	2	1.5
0	2.8	1.2
0	1.5	1
...
1	2.5	2
0	7.7	3.5



Y	X1	X2
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

We **normalize our data** using the Z-Score algorithm:

$$X' = (X - X_{\text{mean}}) / (X_{\text{standard_deviation}})$$

Step 1: Design the neural network

b) Based on the data, design the neural network (1/8)

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34



Output nodes define what kind of data our neural network will result in.

As we have one positive output class, i.e. a **binary classification problem**, we choose one output node. The data is transformed in the correct output specification in this step.

Step 1: Design the neural network

b) Based on the data, design the neural network (2/8)

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34



IN1



ON



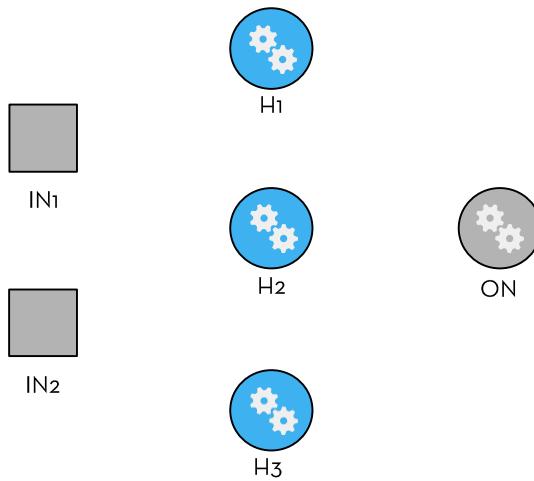
IN2

Input nodes make sure our neural network receives the data it needs. Because we have two inputs, we choose two input nodes.

Step 1: Design the neural network

b) Based on the data, design the neural network (3/8)

Y	X1	X2
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

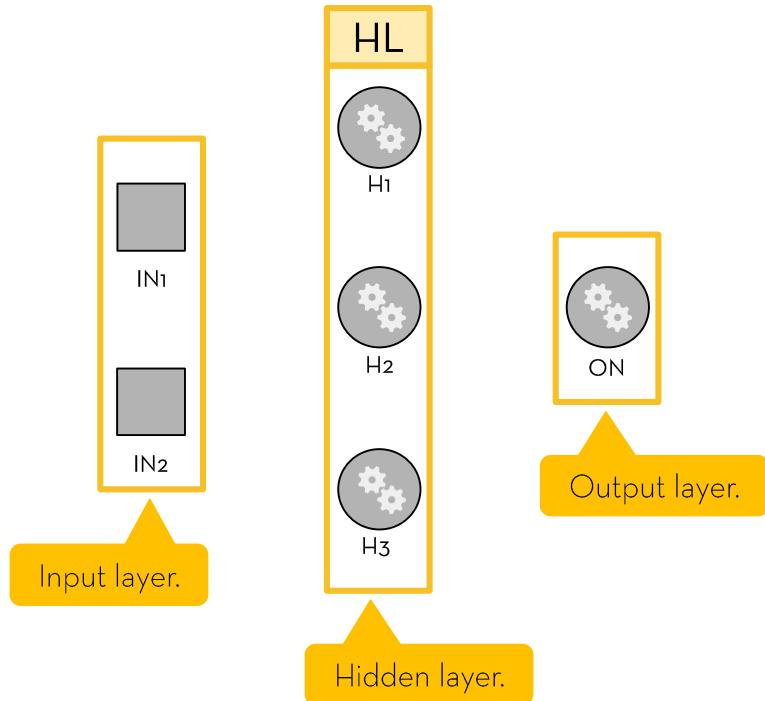


Hidden nodes further transform the data in order to make it easier to «solve the problem». The amount of hidden nodes is up to you.

Step 1: Design the neural network

b) Based on the data, design the neural network (4/8)

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

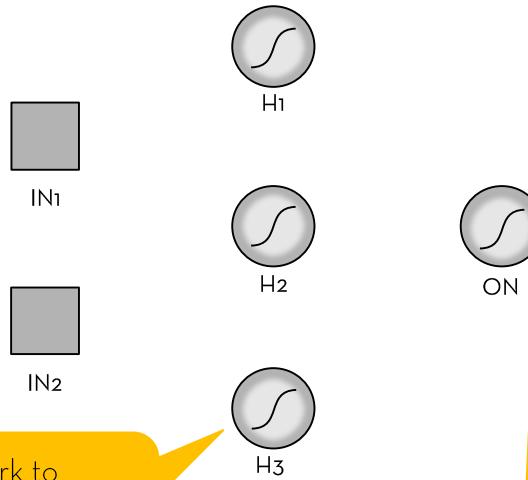


Step 1: Design the neural network

b) Based on the data, design the neural network (5/8)

17

Y	X1	X2
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

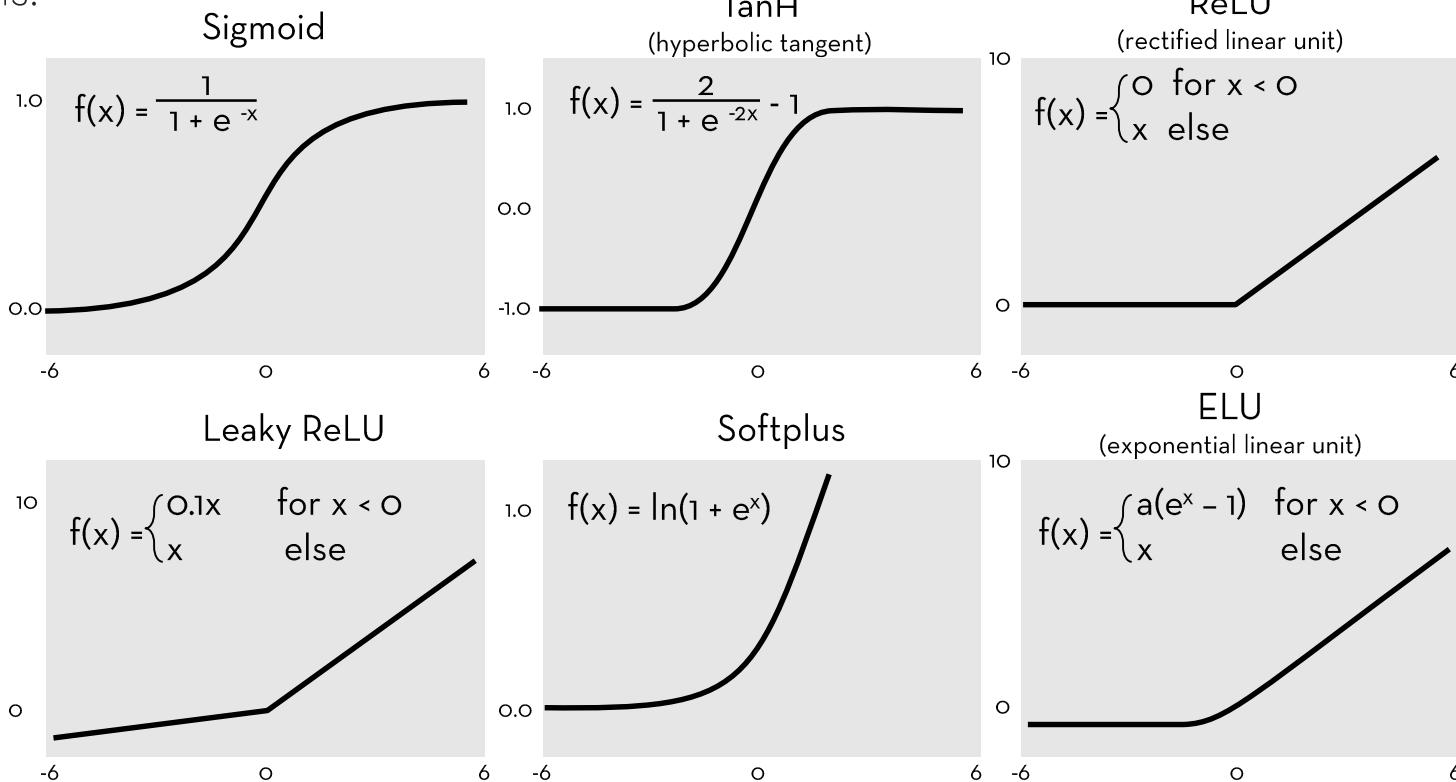


Activation functions allow the neural network to make non-linear approximations. We use the **sigmoid function** for the hidden and output layer.

Since we have a binary classification problem, the output node must have an activation function that restricts the output values to a range between 0 and 1

There are multiple activation functions, i.e. ways to introduce non-linearities in the neural network

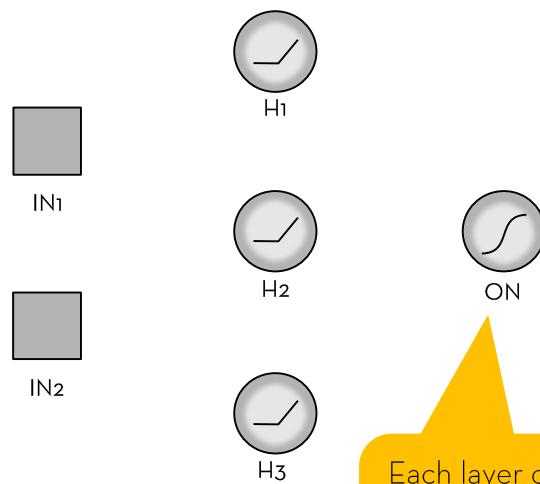
Instead of the sigmoid function, we could also use any other of the many activation functions:



Step 1: Design the neural network

b) Based on the data, design the neural network (6/8)

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34



Each layer can have a different activation function; however, we use sigmoid for each.

Step 1: Design the neural network

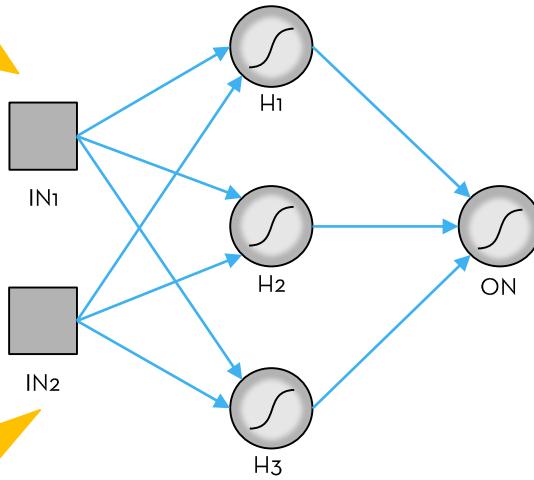
b) Based on the data, design the neural network (7/8)

20

Y	X1	X2
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

In our example of a traditional neural network, we **connect all nodes of a layer with all nodes from the next layer**. This is called a **dense** network. In advanced neural networks, some nodes are connected with only a few specific other nodes.

The nodes are **connected from one layer to each other**, forming a flow from the start to the end. This is called a **feedforward** network.



Step 1: Design the neural network

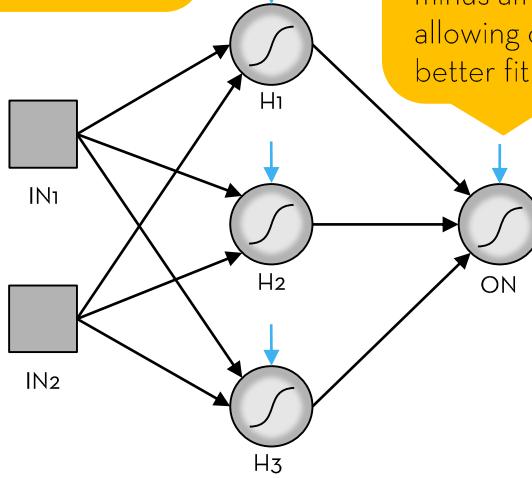
b) Based on the data, design the neural network (8/8)

21

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

In order to assist the hidden nodes and output node in transforming the data, we give each node an **individual bias**.

The **bias** allows us to further transform the data by shifting it by a plus or minus amount, therefore allowing our network to better fit the data.

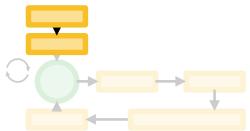


Steps of estimating neural networks

22

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

Step 2: Initialize the neural network

23

- a) Set the weights and biases: Set to random values between 1 and -1.
- b) Set the inputs: Set to available values from the data.

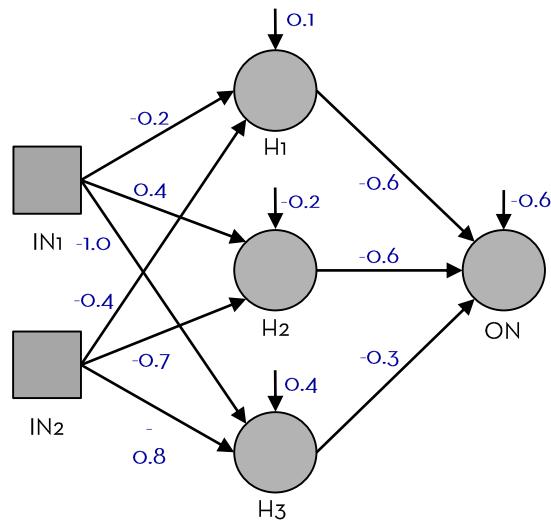
Step 2: Initialize the neural network

a) Set the weights and biases

Y	X₁	X₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

	W₁	W₂	W₃
IN₁ → HL	-0.2	0.4	-1.0
IN₂ → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

Weights and biases
are set randomly.



	H₁	H₂	H₃	ON
Bias	0.1	-0.2	0.4	-0.6

Step 2: Initialize the neural network

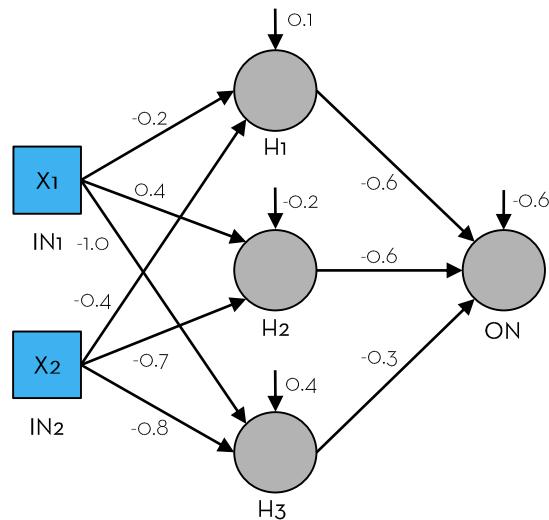
b) Set the inputs

Y	X1	X2
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

As inputs, we take X1 and X2.

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6



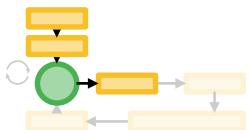
Estimating a neural network: Forward Propagation & Error calculation

Steps of estimating neural networks

27

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

We will first illustrate the general principle of the forward pass before looking at a numerical example.



Repeat 3. - 6. until convergence criteria is met.

Step 3: Forward propagation

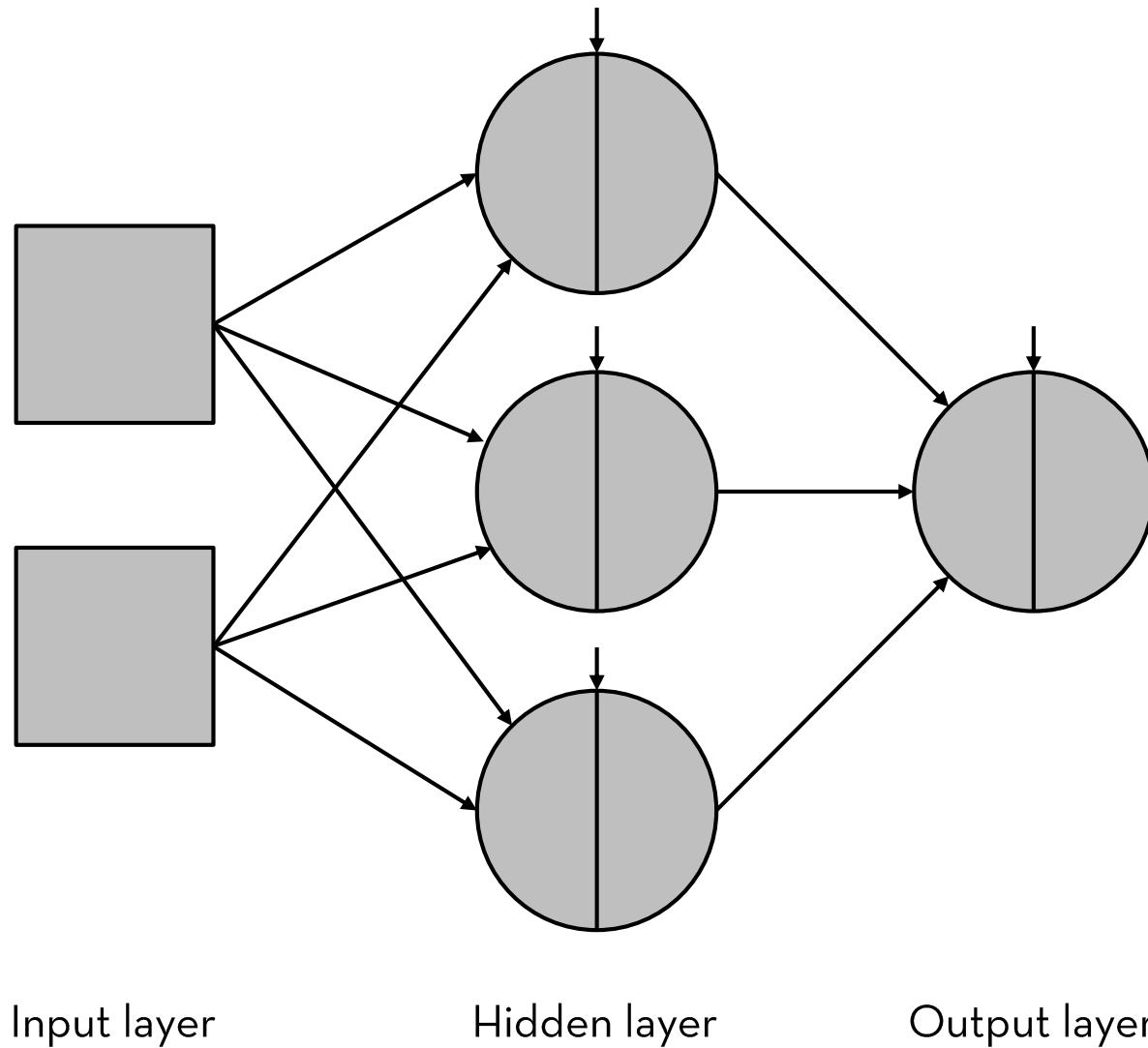
General principle

- a) Calculate the values of the hidden layer.
- b) Calculate the values of the output layer.

We work our way forward through the network.

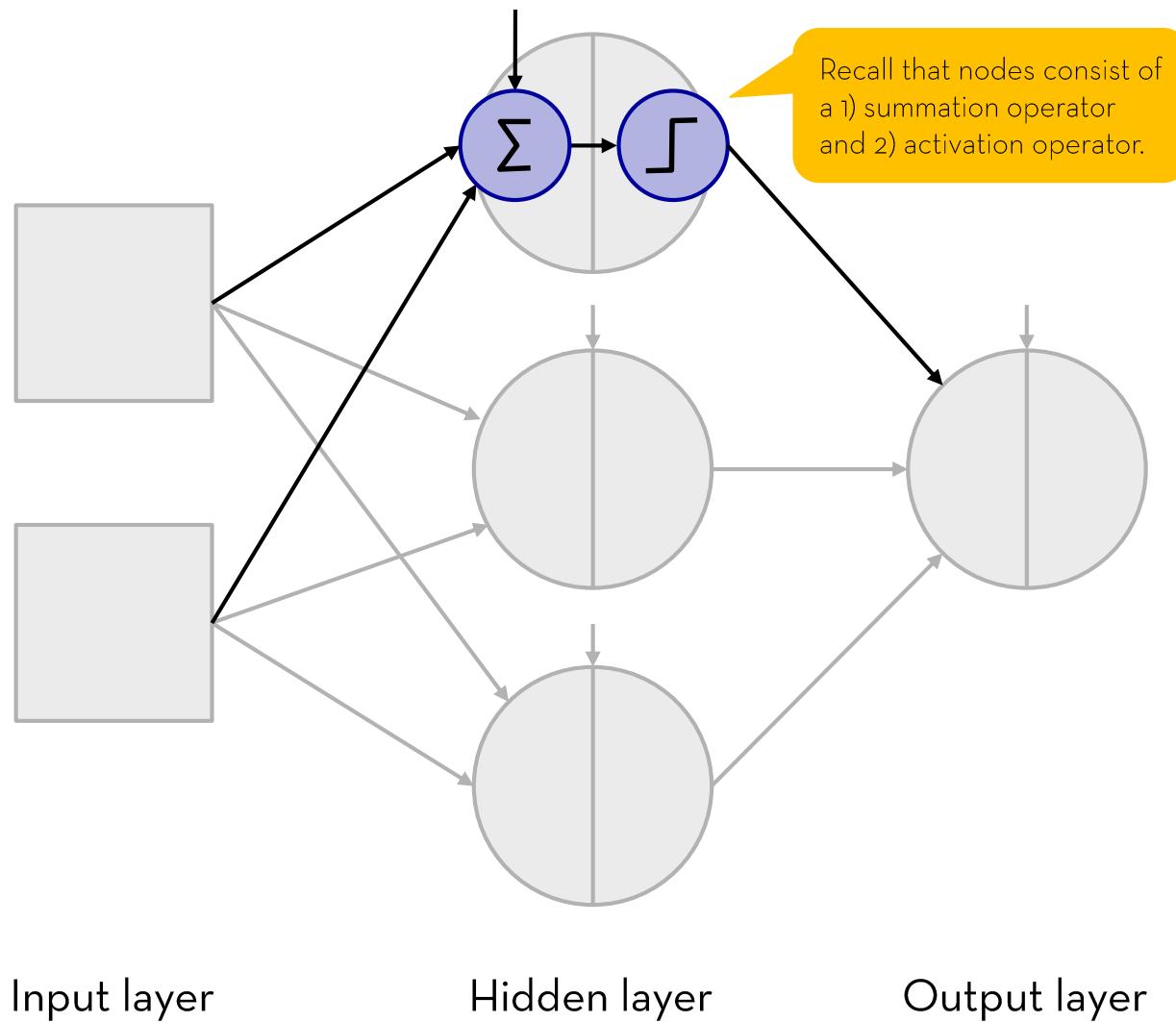
The general principle of forward propagation

29



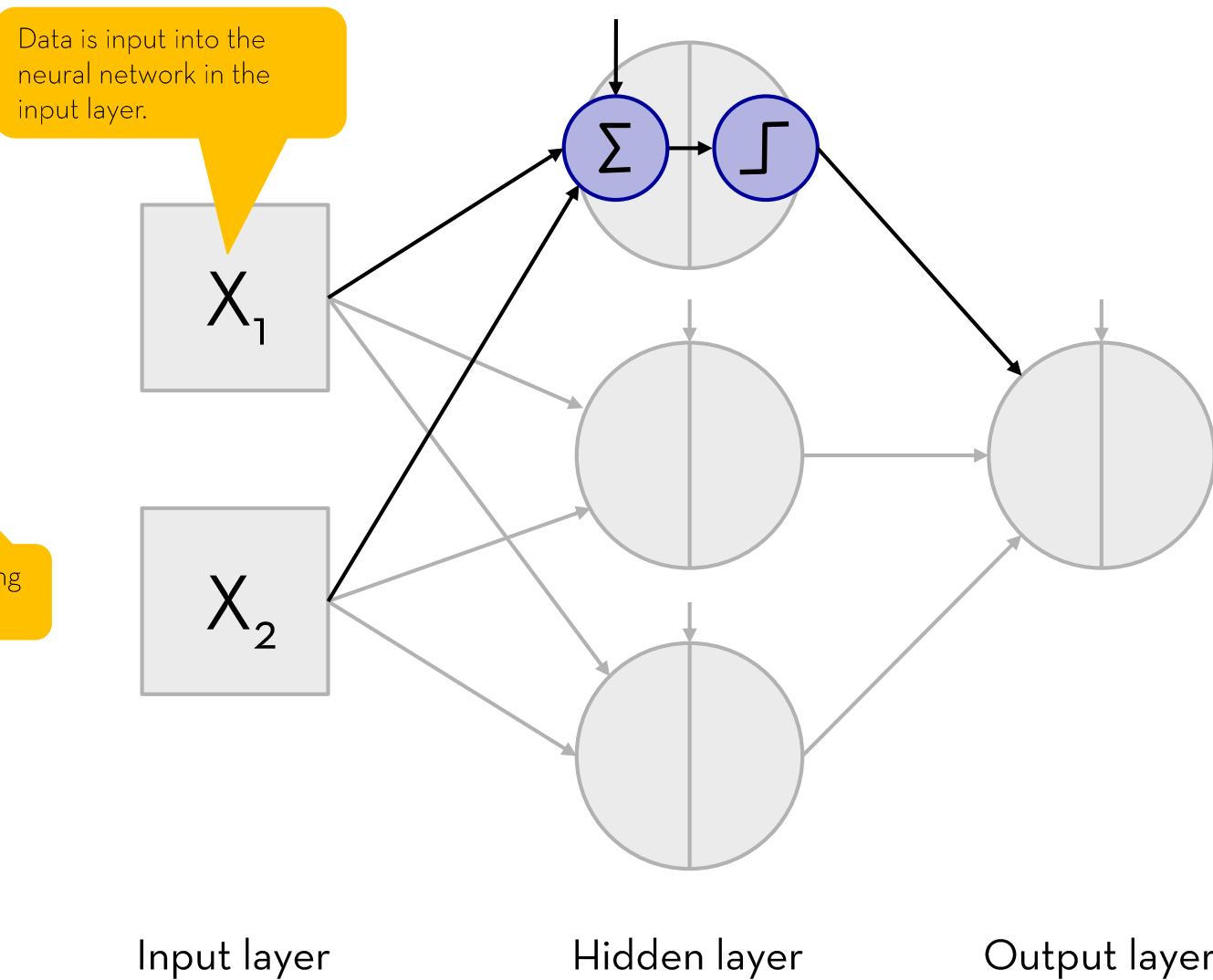
The general principle of forward propagation

Nodes process the data in 2 steps



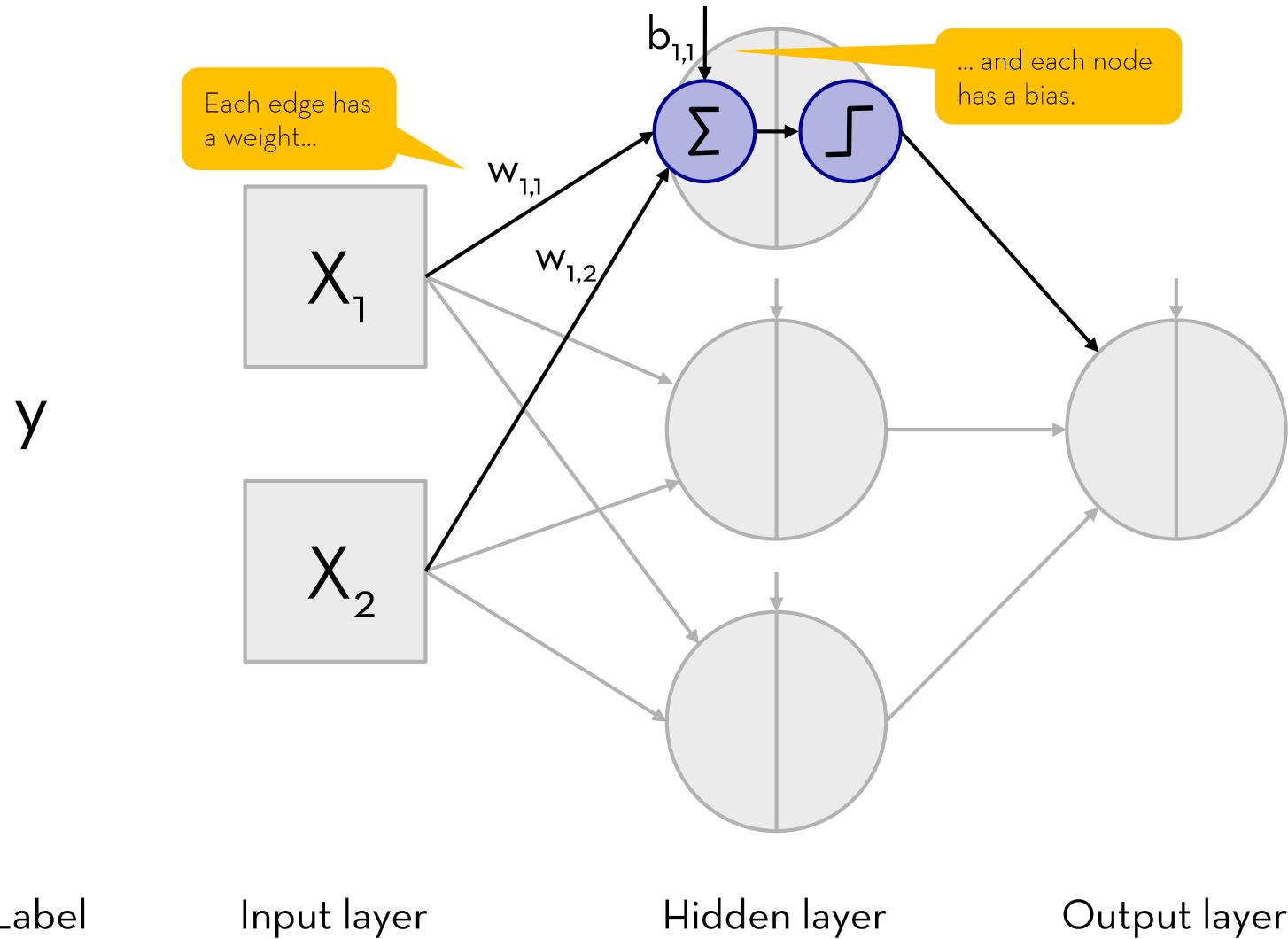
The general principle of forward propagation

31



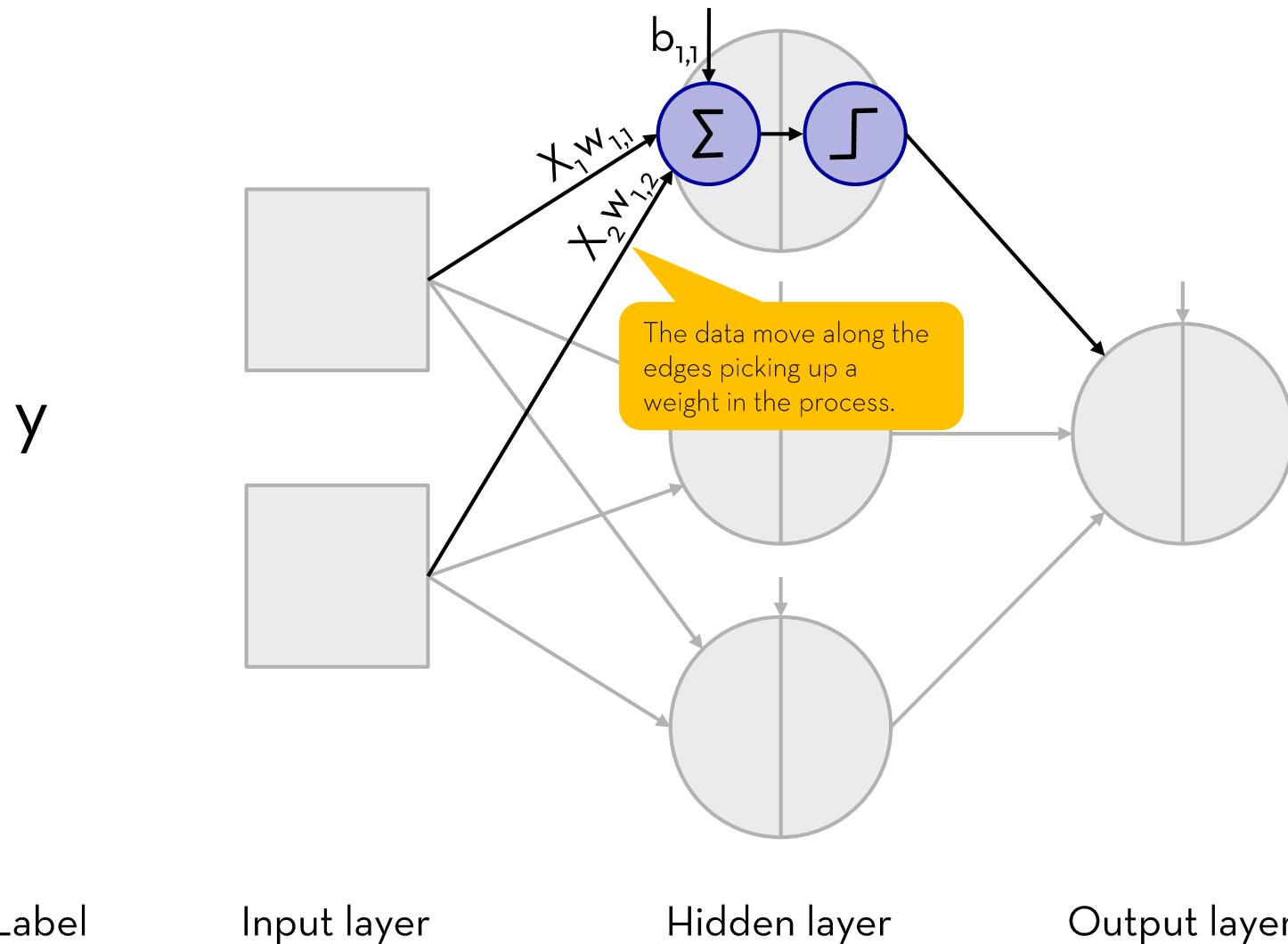
The general principle of forward propagation

a) Calculate the values of the hidden layer (1/5)



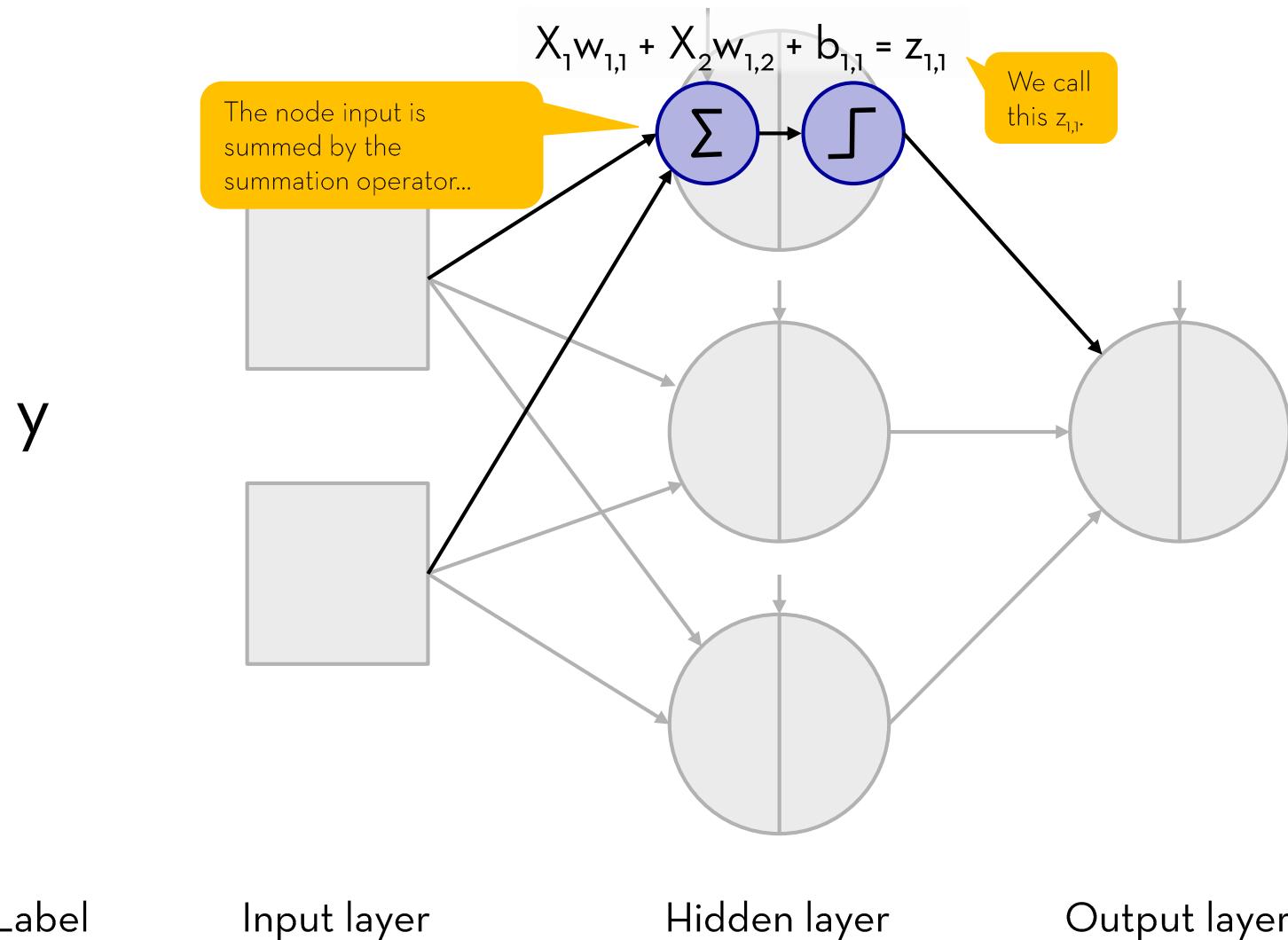
The general principle of forward propagation

a) Calculate the values of the hidden layer (2/5)



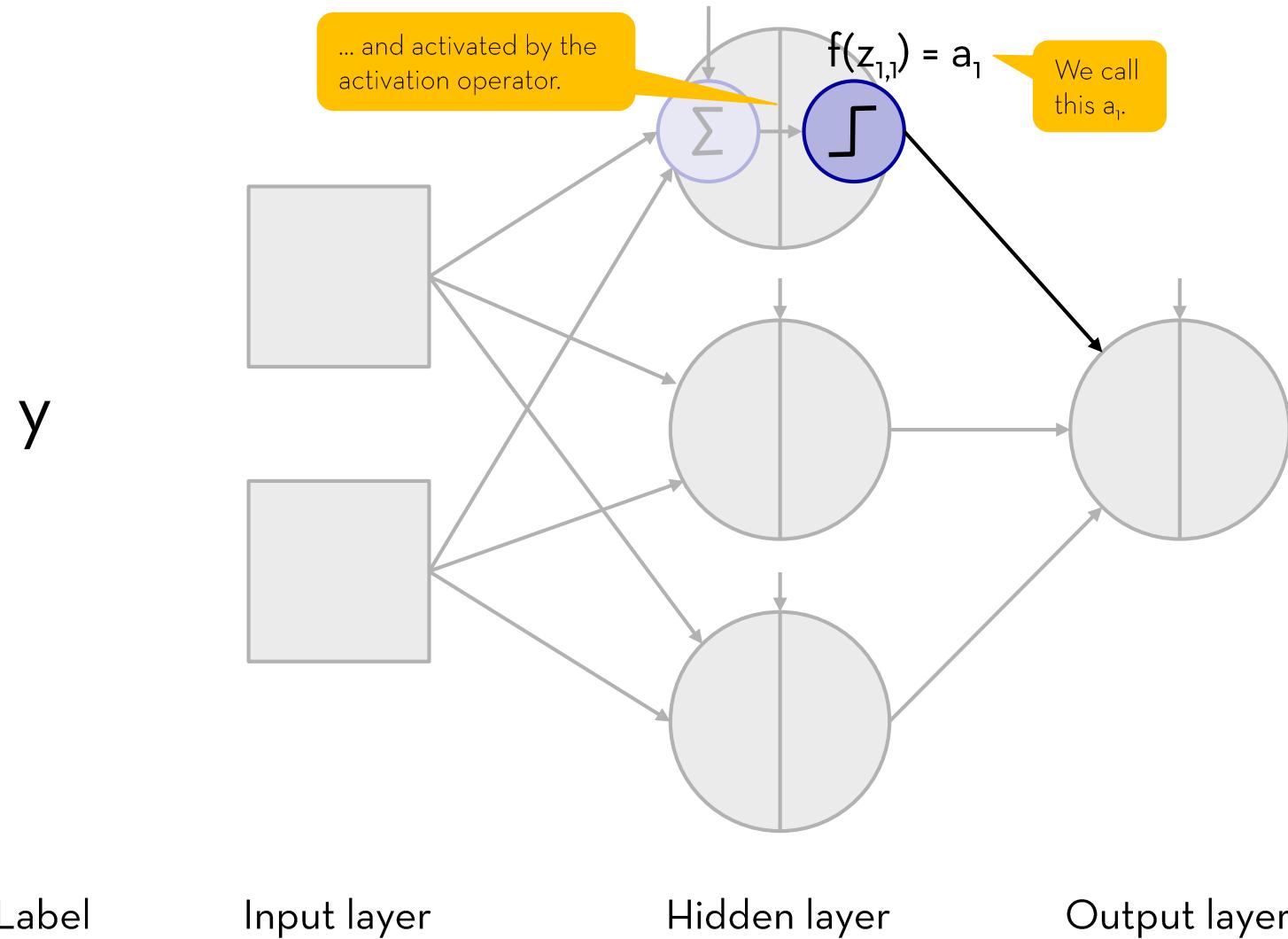
The general principle of forward propagation

a) Calculate the values of the hidden layer (3/5)



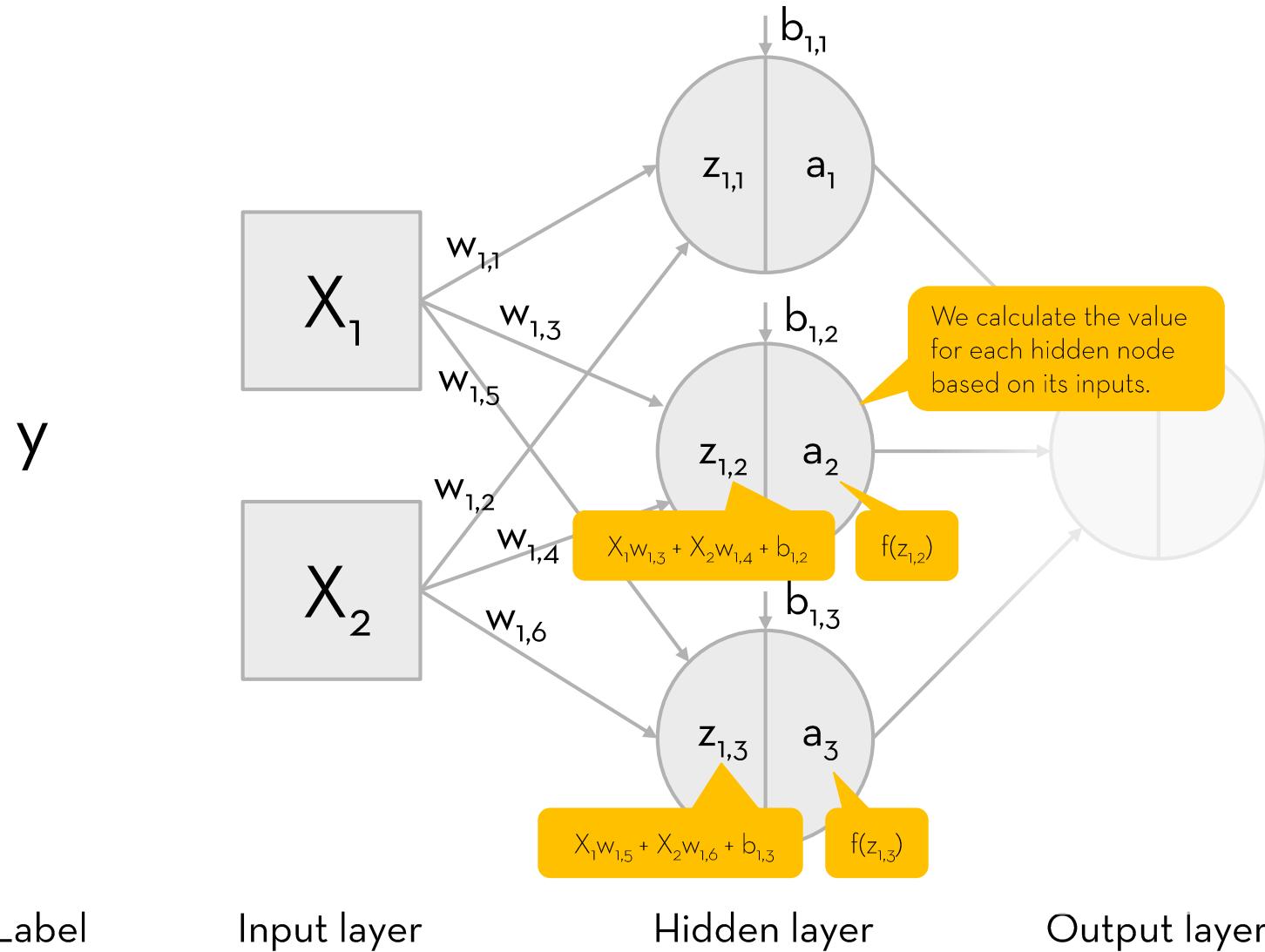
The general principle of forward propagation

a) Calculate the values of the hidden layer (4/5)



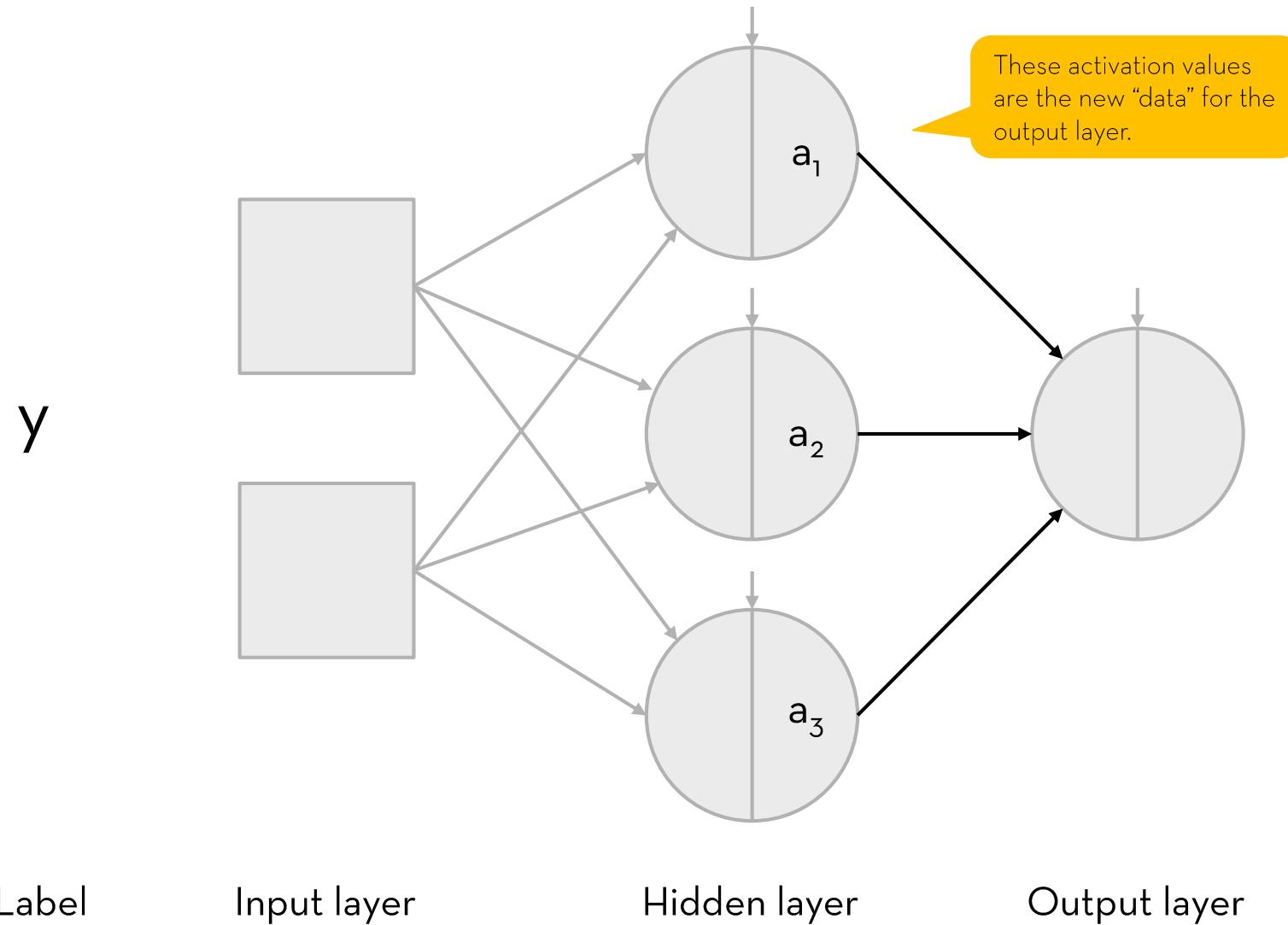
The general principle of forward propagation

a) Calculate the values of the hidden layer (5/5)



The general principle of forward propagation

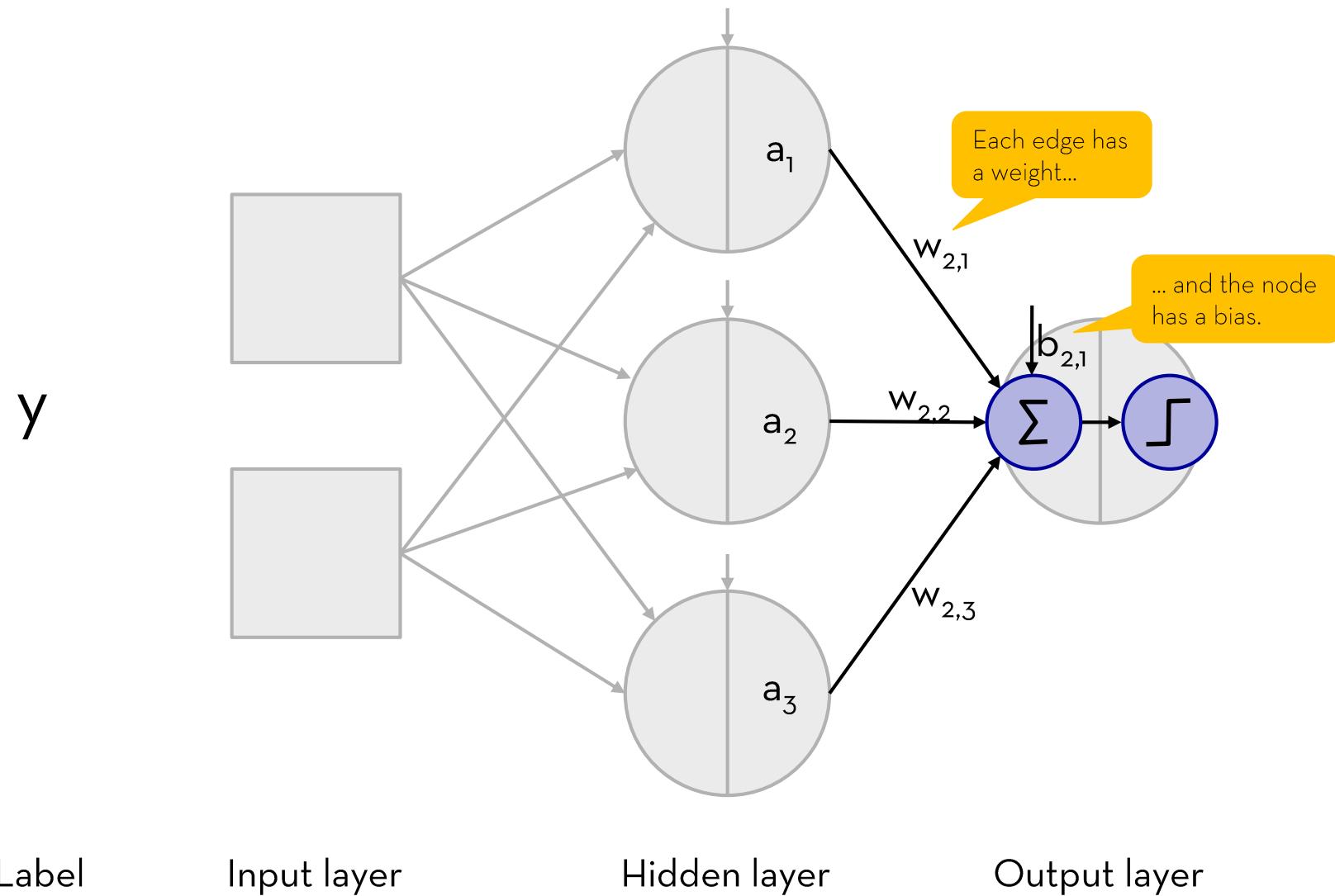
b) Calculate the values of the output layer (1/5)



The general principle of forward propagation

b) Calculate the values of the output layer (2/5)

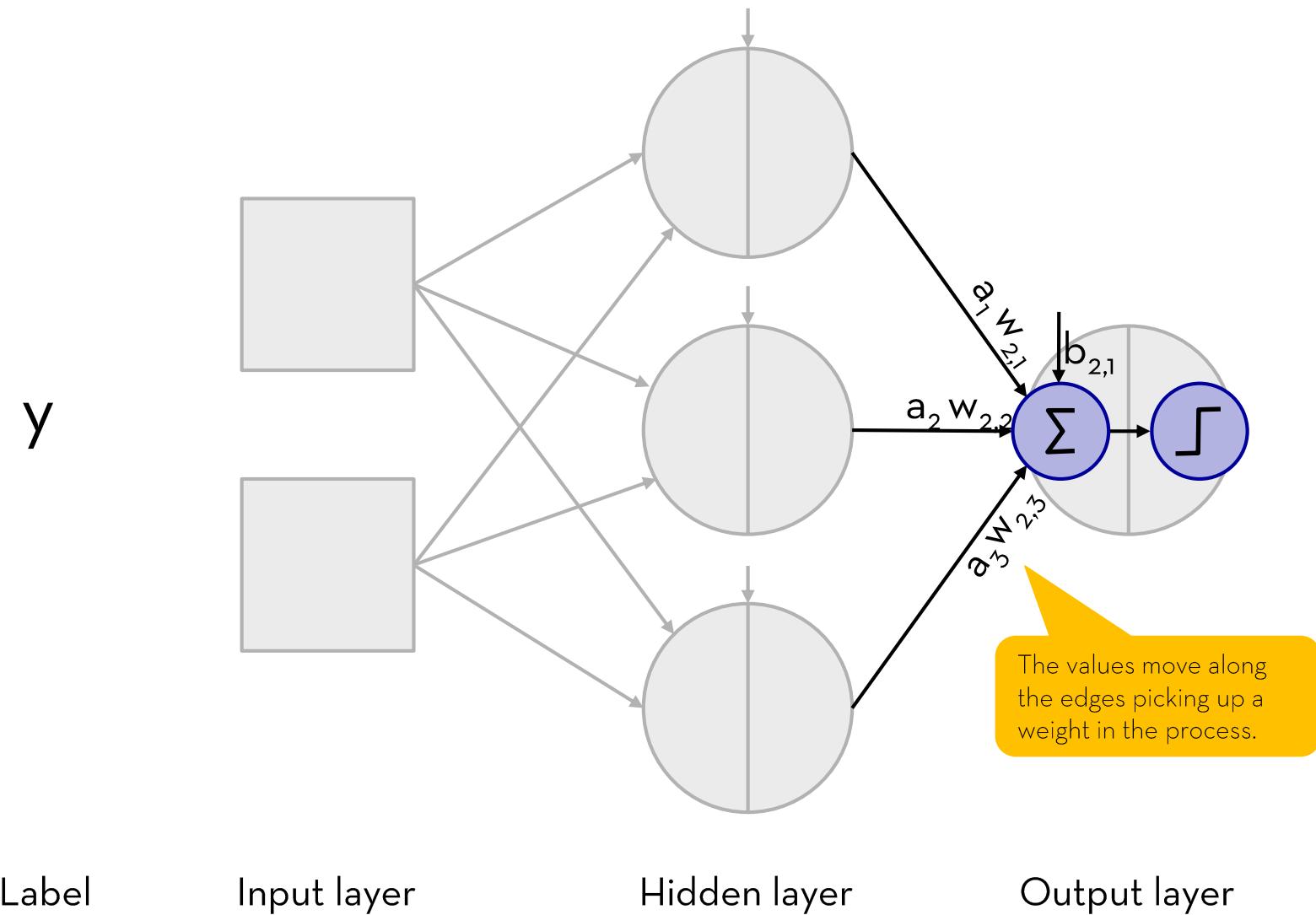
38



The general principle of forward propagation

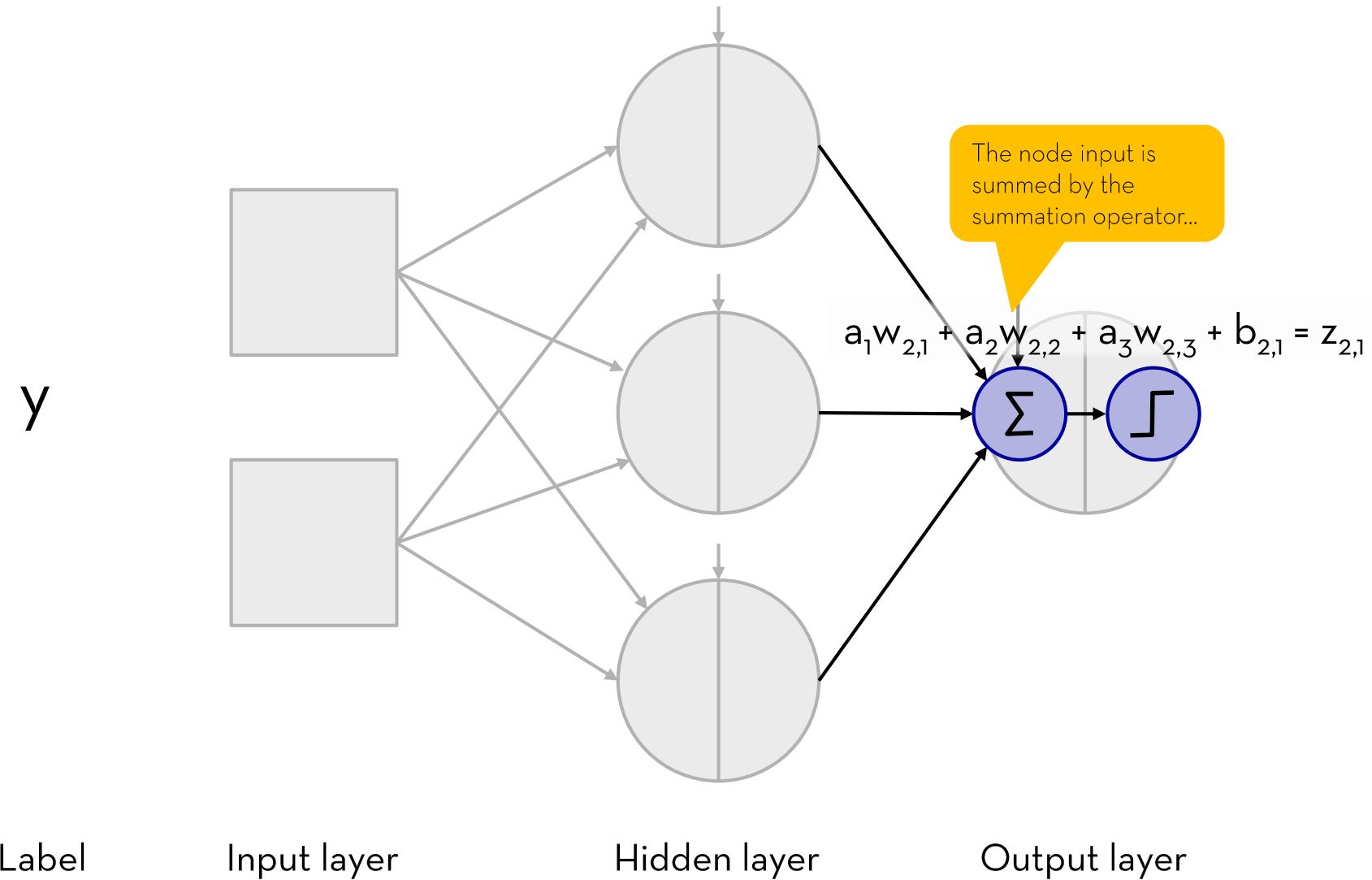
b) Calculate the values of the output layer (3/5)

39



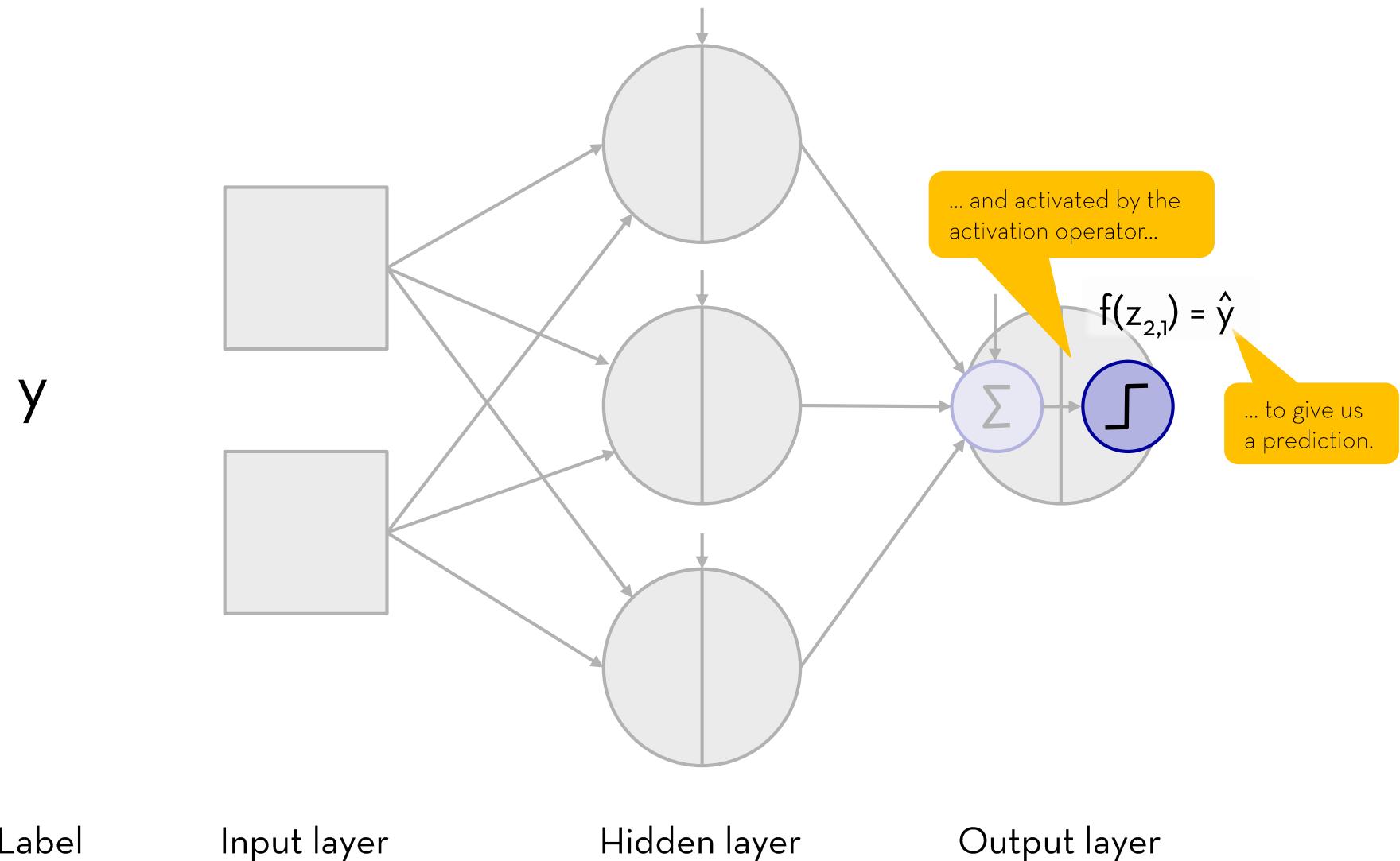
The general principle of forward propagation

b) Calculate the values of the output layer (4/5)



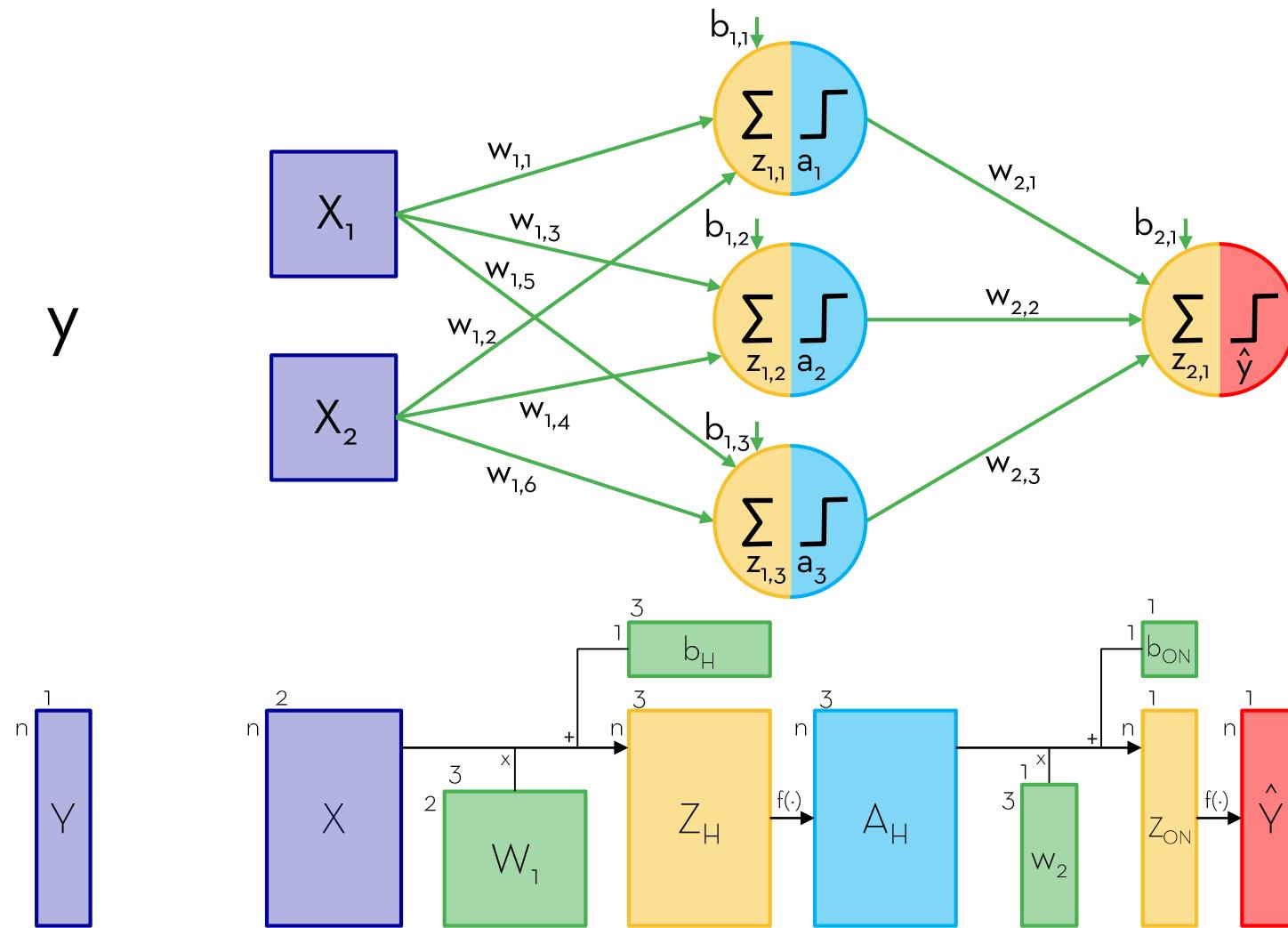
The general principle of forward propagation

b) Calculate the values of the output layer (5/5)



A summary of the general principle of forward propagation

42



Step 3: Forward propagation

Numerical example

- a) Calculate the values of the hidden layer.
- b) Calculate the values of the output layer.

We work our way forward through the network.

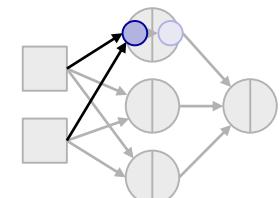
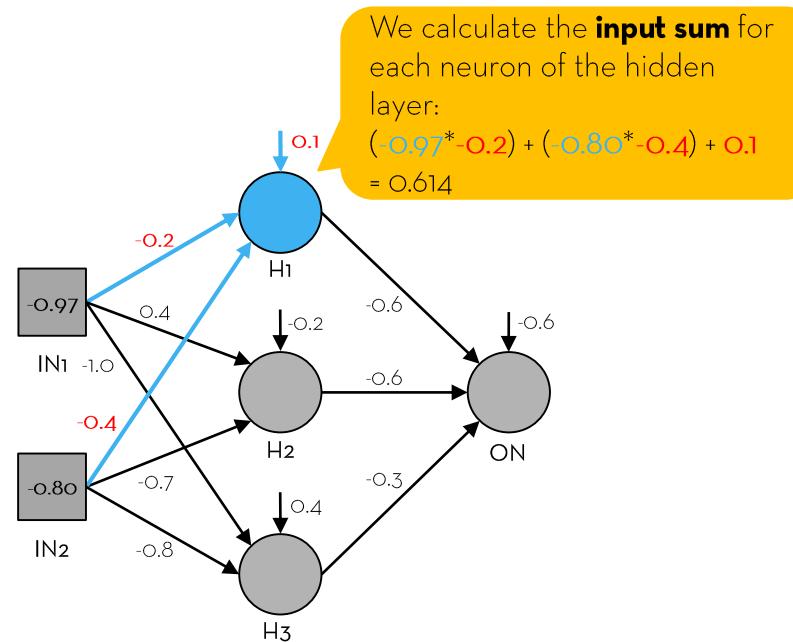
Step 3: Forward propagation

a) Calculate the values of the hidden layer (1/5)

Y	X ₁	X ₂
0	-0.97	-0.80
0	-0.69	-0.97
0	-1.14	-1.09
...
1	-0.79	-0.52
0	1.04	0.34

	W ₁	W ₂	W ₃
IN ₁ → HL	-0.2	0.4	-1.0
IN ₂ → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H ₁	H ₂	H ₃	ON
Bias	0.1	-0.2	0.4	-0.6



Step 3: Forward propagation

a) Calculate the values of the hidden layer (2/5)

Y	X ₁	X ₂	H ₁	H ₂	H ₃
0	-0.97	-0.80	0.65		
0	-0.69	-0.97			
0	-1.14	-1.09			
...			
1	-0.79	-0.52			
0	1.04	0.34			

	W ₁	W ₂	W ₃
IN ₁ → HL	-0.2	0.4	-1.0
IN ₂ → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

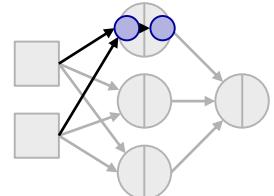
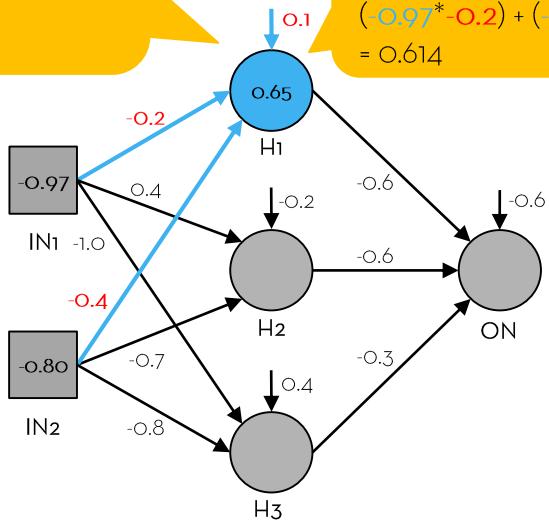
	H ₁	H ₂	H ₃	ON
Bias	0.1	-0.2	0.4	-0.6

We then apply the **sigmoid function**, to get the output of a node (which is its value):

$$f(x) = 1/(1+e^{-x})$$

Here: $f(0.614) = 0.65$

We calculate the **input sum** for each hidden layer:

$$(-0.97 * -0.2) + (-0.80 * -0.4) + 0.1 = 0.614$$


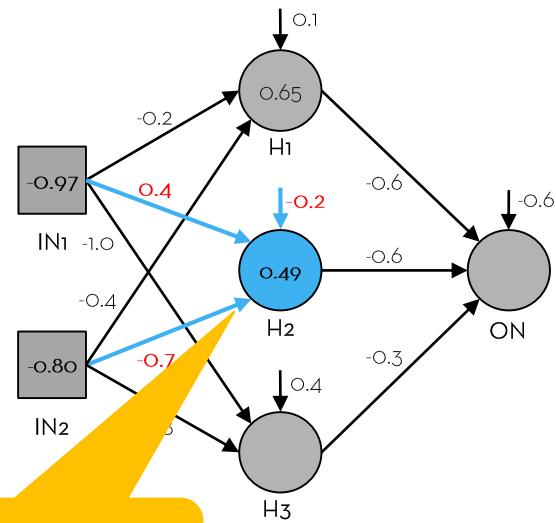
Step 3: Forward propagation

a) Calculate the values of the hidden layer (3/5)

Y	X1	X2	H1	H2	H3
O	-0.97	-0.80	0.65	0.49	
O	-0.69	-0.97			
O	-1.14	-1.09			
...			
1	-0.79	-0.52			
O	1.04	0.34			

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

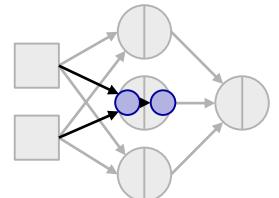


Input sum:

$$(-0.97 \cdot 0.4) + (-0.80 \cdot -0.7) + -0.2$$

$$= -0.026$$

$f(-0.026) = 0.49$



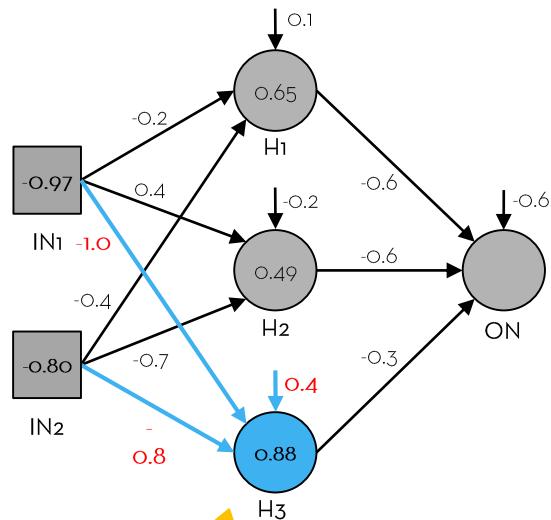
Step 3: Forward propagation

a) Calculate the values of the hidden layer (4/5)

Y	X ₁	X ₂	H ₁	H ₂	H ₃
O	-0.97	-0.80	0.65	0.49	0.88
O	-0.69	-0.97			
O	-1.14	-1.09			
...			
1	-0.79	-0.52			
O	1.04	0.34			

	W ₁	W ₂	W ₃
IN ₁ → HL	-0.2	0.4	-1.0
IN ₂ → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

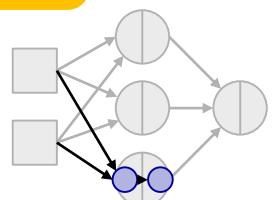
	H ₁	H ₂	H ₃	ON
Bias	0.1	-0.2	0.4	-0.6



Input sum:

$$(-0.97 * -1.0) + (-0.80 * -0.8) + 0.4 = 2.01$$

$$f(2.01) = 0.88$$

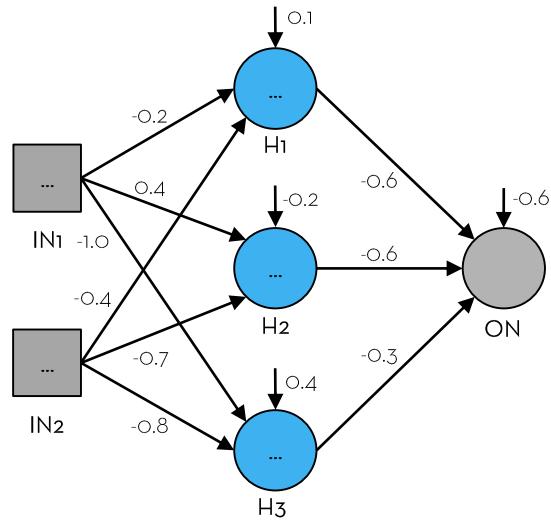


Step 3: Forward propagation

a) Calculate the values of the hidden layer (5/5)

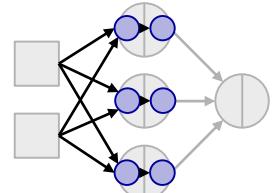
Y	X ₁	X ₂	H ₁	H ₂	H ₃
0	-0.97	-0.80	0.65	0.49	0.88
0	-0.69	-0.97	0.65	0.55	0.87
0	-1.14	-1.09	0.68	0.53	0.92
...
1	-0.79	-0.52	0.61	0.46	0.83
0	1.04	0.34	0.44	0.49	0.29

We do the same for all other observations.



	W ₁	W ₂	W ₃
IN ₁ → HL	-0.2	0.4	-1.0
IN ₂ → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H ₁	H ₂	H ₃	ON
Bias	0.1	-0.2	0.4	-0.6



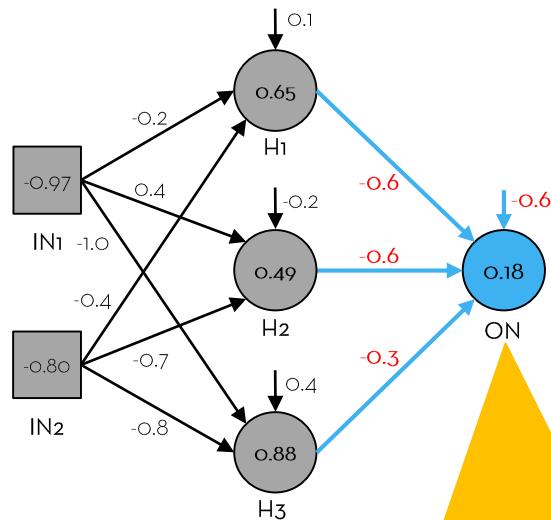
Step 3: Forward propagation

b) Calculate the values of the output layer (2/3)

Y	X1	X2	H1	H2	H3	Output
0	-0.97	-0.80	0.65	0.49	0.88	0.18
0	-0.69	-0.97	0.65	0.55	0.87	
0	-1.14	-1.09	0.68	0.53	0.92	
...	
1	-0.79	-0.52	0.61	0.46	0.83	
0	1.04	0.34	0.44	0.49	0.29	

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

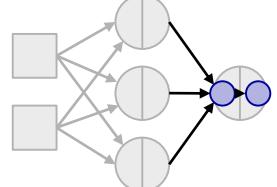
	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6



Input sum:

$$(0.65 * -0.6) + (0.49 * -0.6) + (0.88 * -0.3) + -0.6 = -1.55$$

$$f(-1.55) = 0.18$$

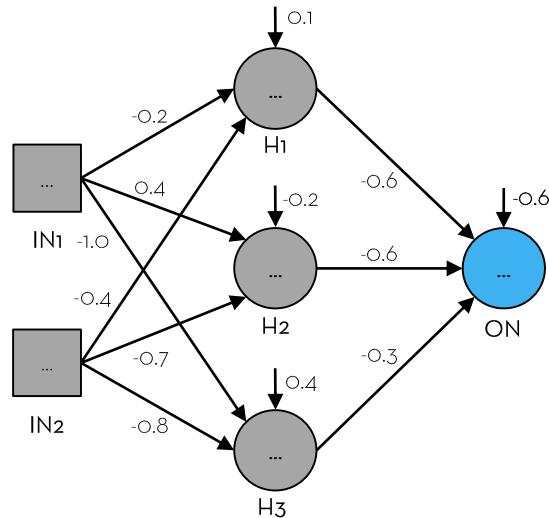


Step 3: Forward propagation

b) Calculate the values of the output layer (3/3)

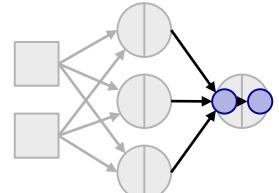
Y	X1	X2	H1	H2	H3	Output
0	-0.97	-0.80	0.65	0.49	0.88	0.18
0	-0.69	-0.97	0.65	0.55	0.87	0.17
0	-1.14	-1.09	0.68	0.53	0.92	0.17
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18
0	1.04	0.34	0.44	0.49	0.29	0.22

We do the **same** for all other observations.



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

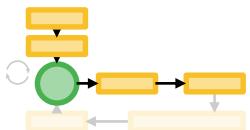


Steps of estimating neural networks

51

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

The general principle of calculating the cost

52

- The **cost function** (a.k.a. loss function) produces a scalar value that quantifies the model's error by comparing the model's response with the correct response.
- It is used to evaluate the possible alternative solutions, i.e. if one combination of parameters provides a better model fit than another combination of parameters.
- We need to define the cost function for our task. We will choose **mean squared error**:

$$C(W) = \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{1}{2}(\hat{y}_i - y_i)^2}_{\text{Cost for observation } i (C_i)}$$

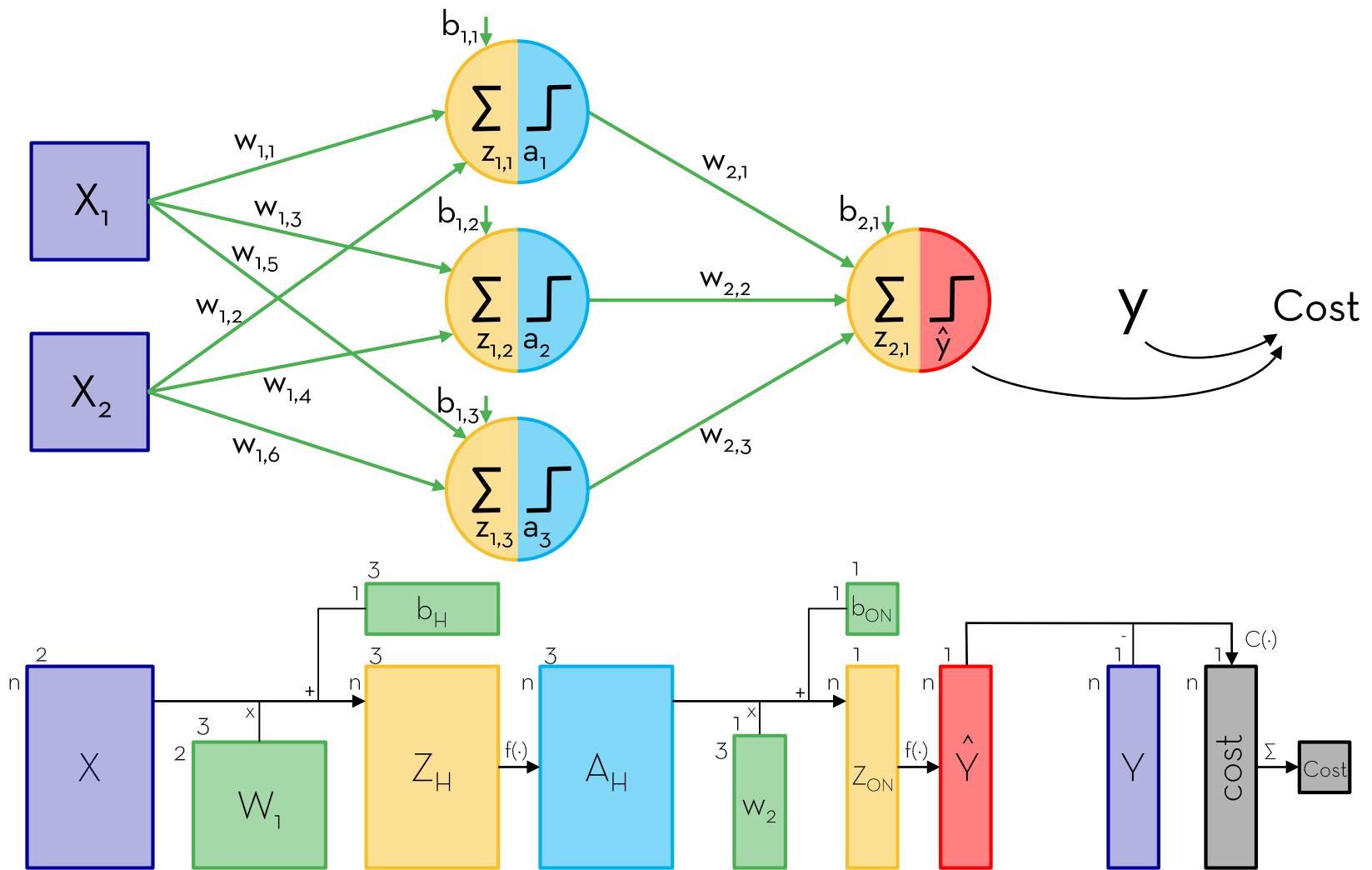
Model's response.

Correct response.

The $\frac{1}{2}$ simplifies the math but can otherwise be ignored.
If $y = \frac{1}{2}x^2 \Rightarrow dy / dx = x$

A summary of the general principle of calculating the cost

53



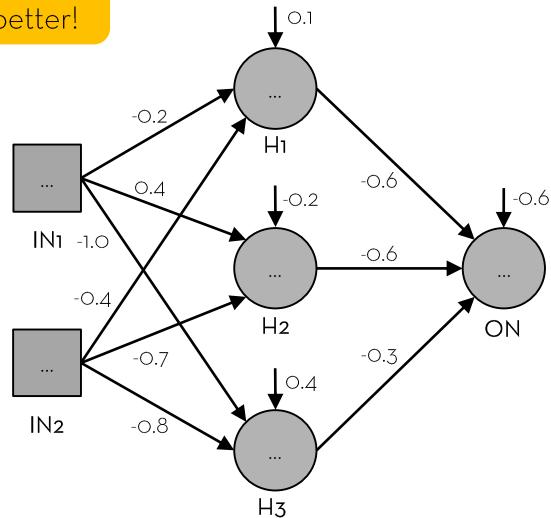
Step 4: Calculate the cost

Is our model good or bad? (1/2)

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.02
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.01
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.01
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.33
0	1.04	0.34	0.44	0.49	0.29	0.22	0.02

Our \hat{y} is labeled "output".

$\text{Cost}_i = \frac{1}{2}(\text{output}_i - y_i)^2$
→ The smaller, the better!



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

$$C(W) = \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{1}{2}(\hat{y}_i - y_i)^2}_{\text{Cost for observation } i (C_i)}$$

Step 4: Calculate the cost

Is our model good or bad? (2/2)

55

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.02
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.01
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.01
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.33
0	1.04	0.34	0.44	0.49	0.29	0.22	0.02

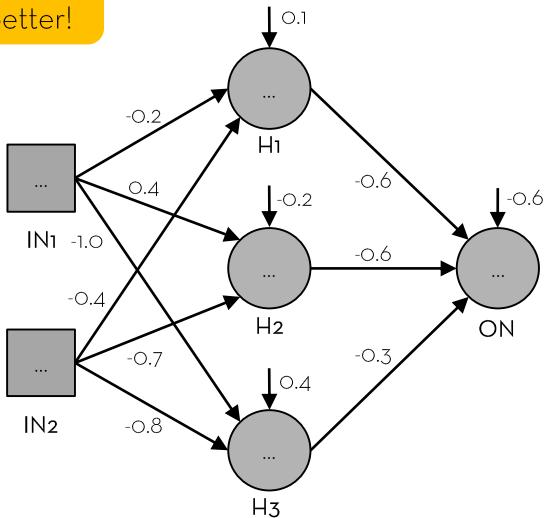
$$\text{Cost}_i = \frac{1}{2}(\text{output}_i - y_i)^2$$

→ The smaller, the better!

$$\Sigma = 0.16$$

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

Our averaged up cost function returns 0.16.



	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

$$C(W) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2}(\hat{y}_i - y_i)^2$$

Estimating a neural network: Backpropagation

Steps of estimating neural networks

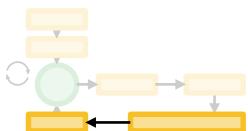
57

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.

5. Backpropagation.
6. Update weights.

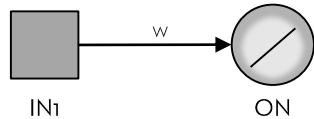
The remaining 2 steps constitute **gradient descent**. Gradient descent is the first-order optimization technique we use to train a neural network.

Repeat 3. - 6. until convergence criteria is met.



If we only have one weight, it is easy to visualize the cost function

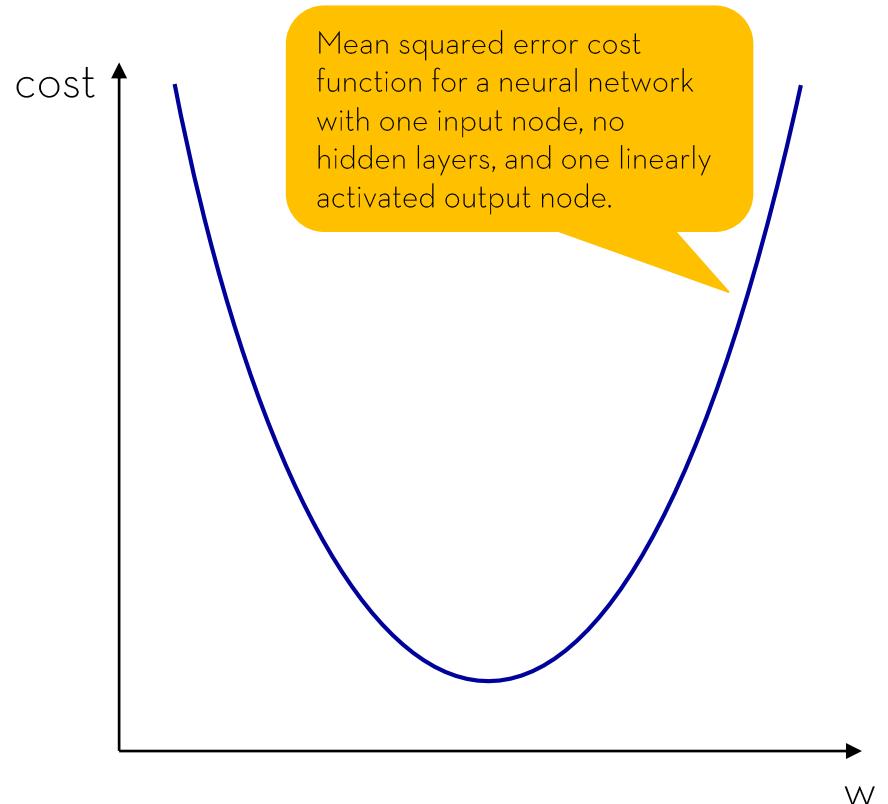
For a very simple neural network:



it is simple to visualize the cost function:

$$C(W) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (x_i w - y_i)^2.$$

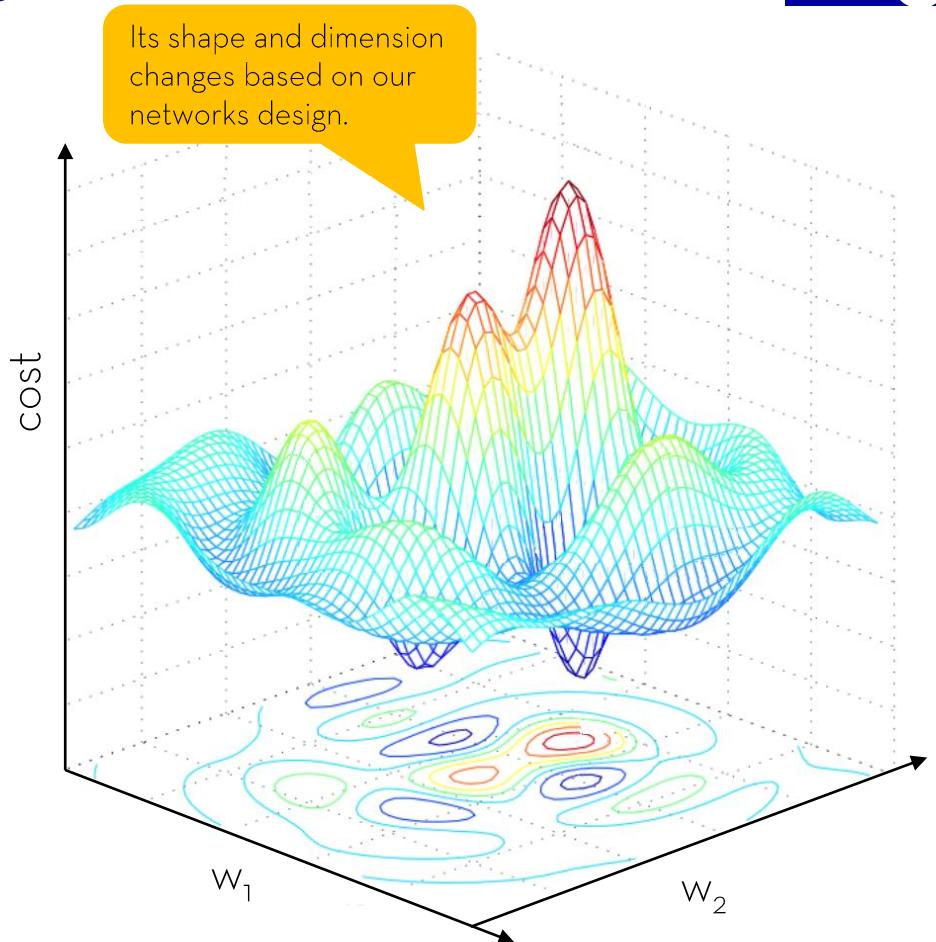
Once we have more than one weight, the cost function is harder to visualize.



Neural network cost functions can be very non-linear and difficult to optimize

The cost is a hyper-dimensional non-linear and non-convex function whose shape we don't know.

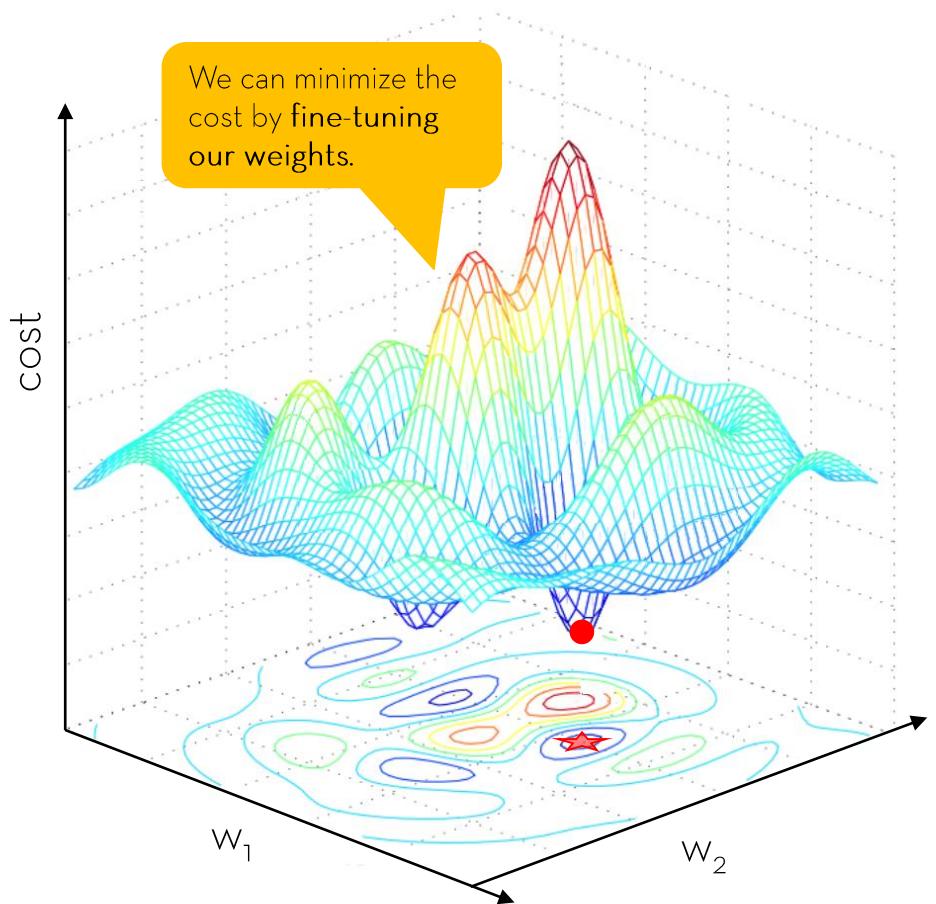
- **Hyper-dimensional** because we are training many weights (in our case 9 weights and 4 biases).
- **Non-linear** and **non-convex** because of the activation functions and cost function.



How do we optimize our weights?

60

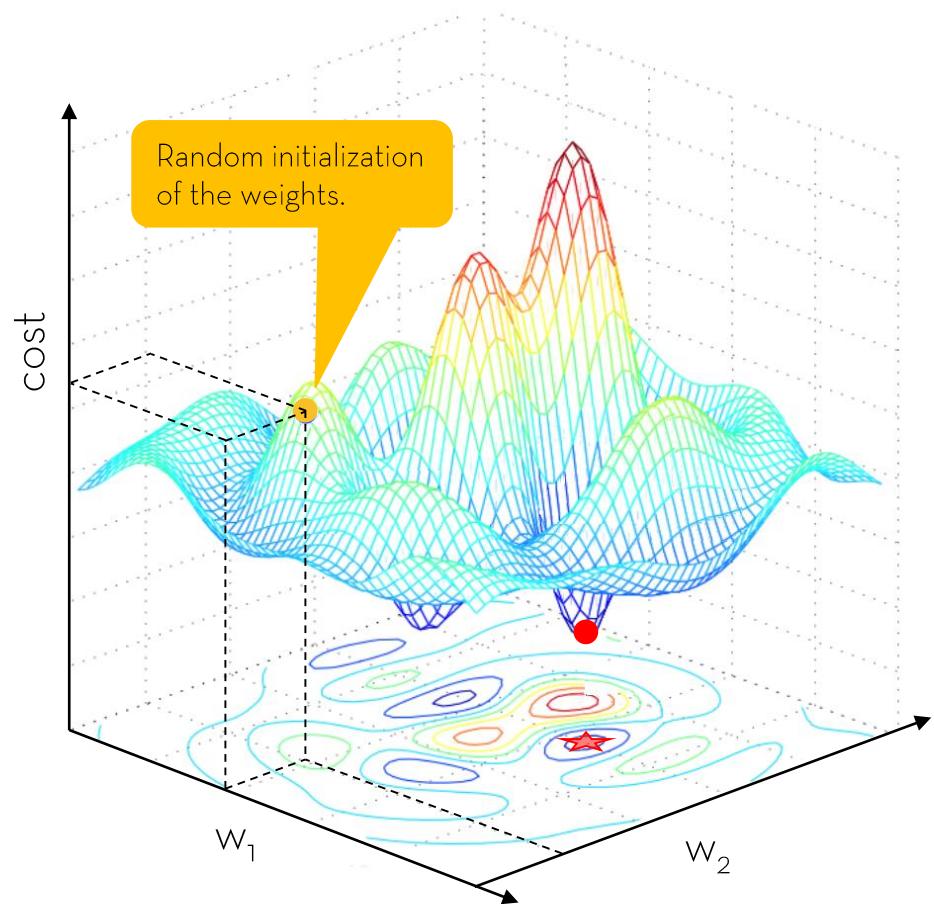
- We are trying to find the **optimal weights** – i.e., the combination of weights that minimizes the cost function.



How do we optimize our weights?

61

- We are trying to find the **optimal weights** – i.e., the combination of weights that minimizes the cost function.
- By randomly initializing the weights (step 2) we start at a random point on the function.

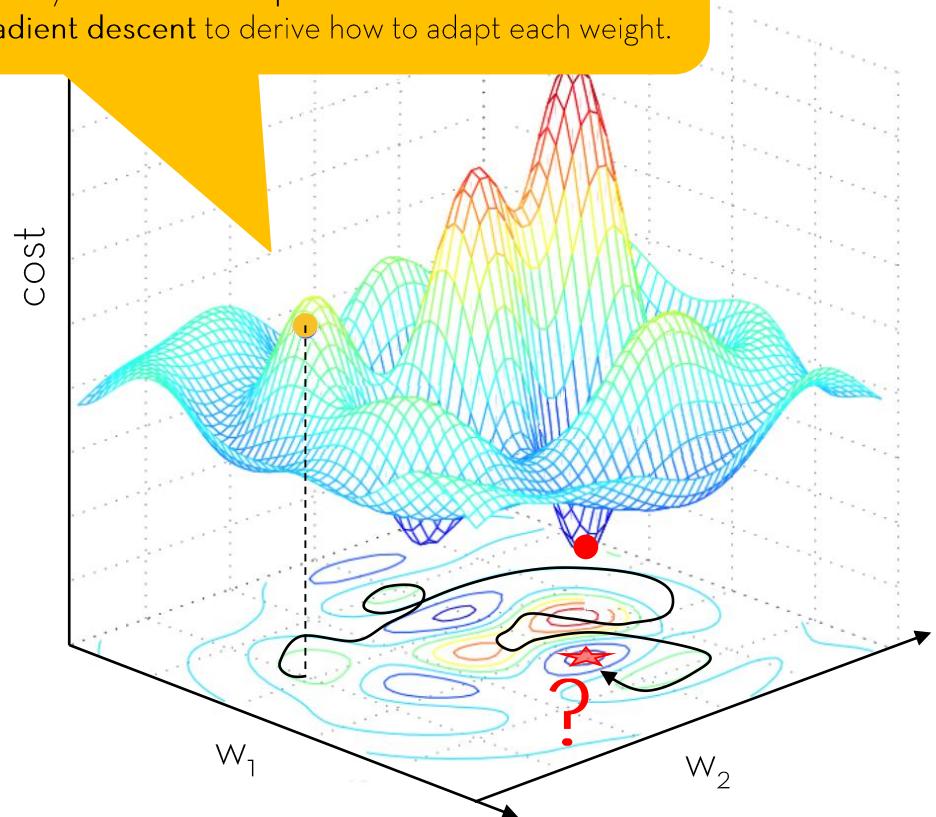


How do we optimize our weights?

62

- We are trying to find the **optimal weights** – i.e., the combination of weights that minimizes the cost function.
- By randomly initializing the weights (step 2) we start at a random point on the function.
- How do we find the best combination of weights? We use **gradient descent** to navigate to the lowest cost solution.

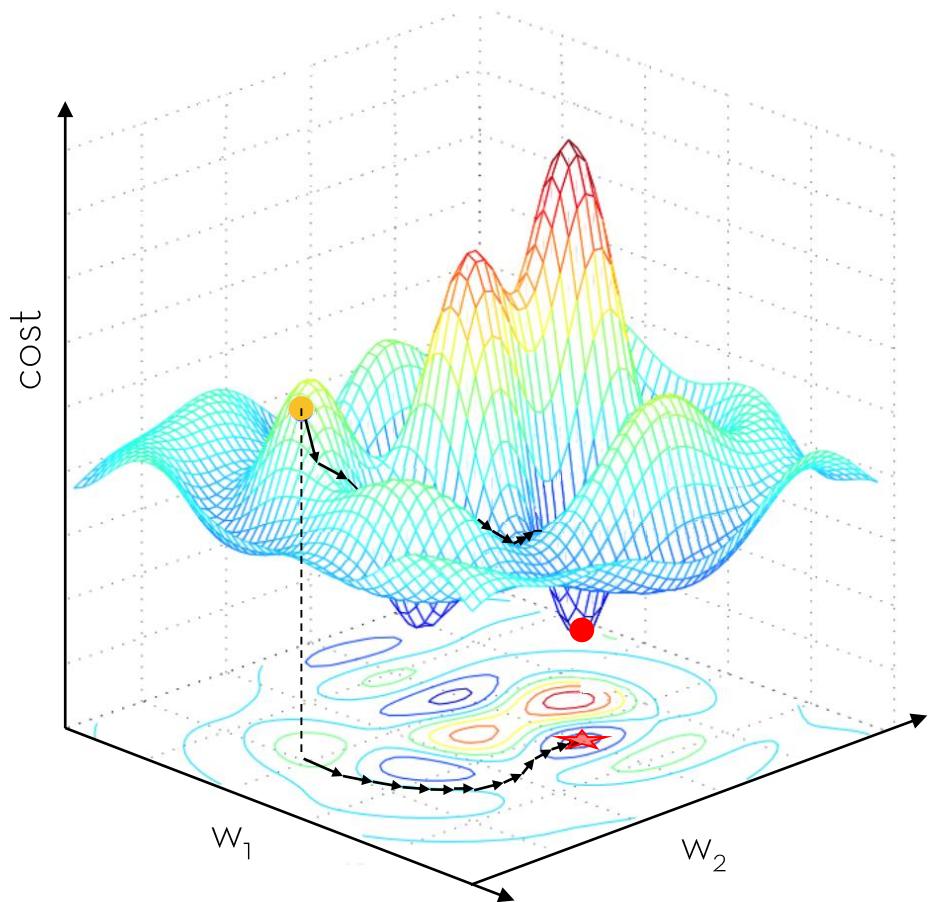
Simply trying all possible weights to find the best combination is not possible due to exponential work. We rely on **numerical optimization methods** such as gradient descent to derive how to adapt each weight.



Why gradient?

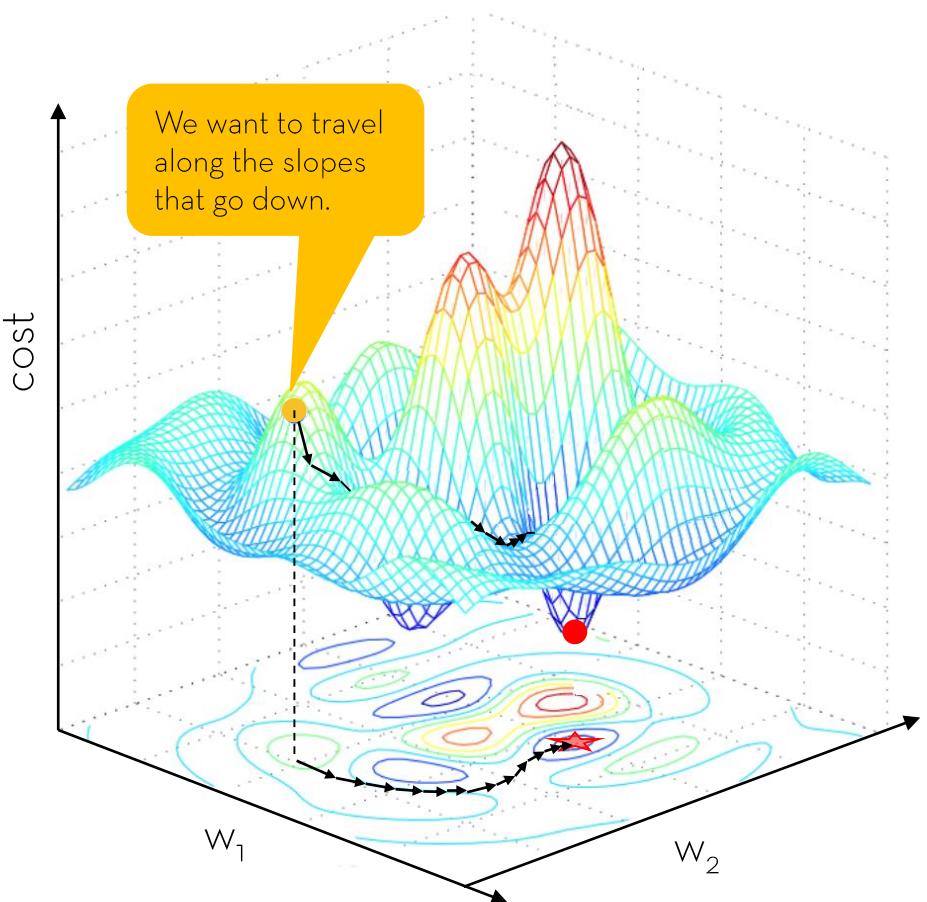
63

- The gradient tells us in which direction and how far we need to go for each weight.



Why descent?

- The gradient tells us in **which direction** and **how far we need to go** for each weight.
- Because we are trying to **minimize the cost**, we need to travel in the direction that **descends the gradient**.
- That means that we want to update the weights such that we travel in the direction in which the gradient is negative.



A quick refresher on gradients

65

The gradient answers the question
“how much does y change if we
change x ? ”

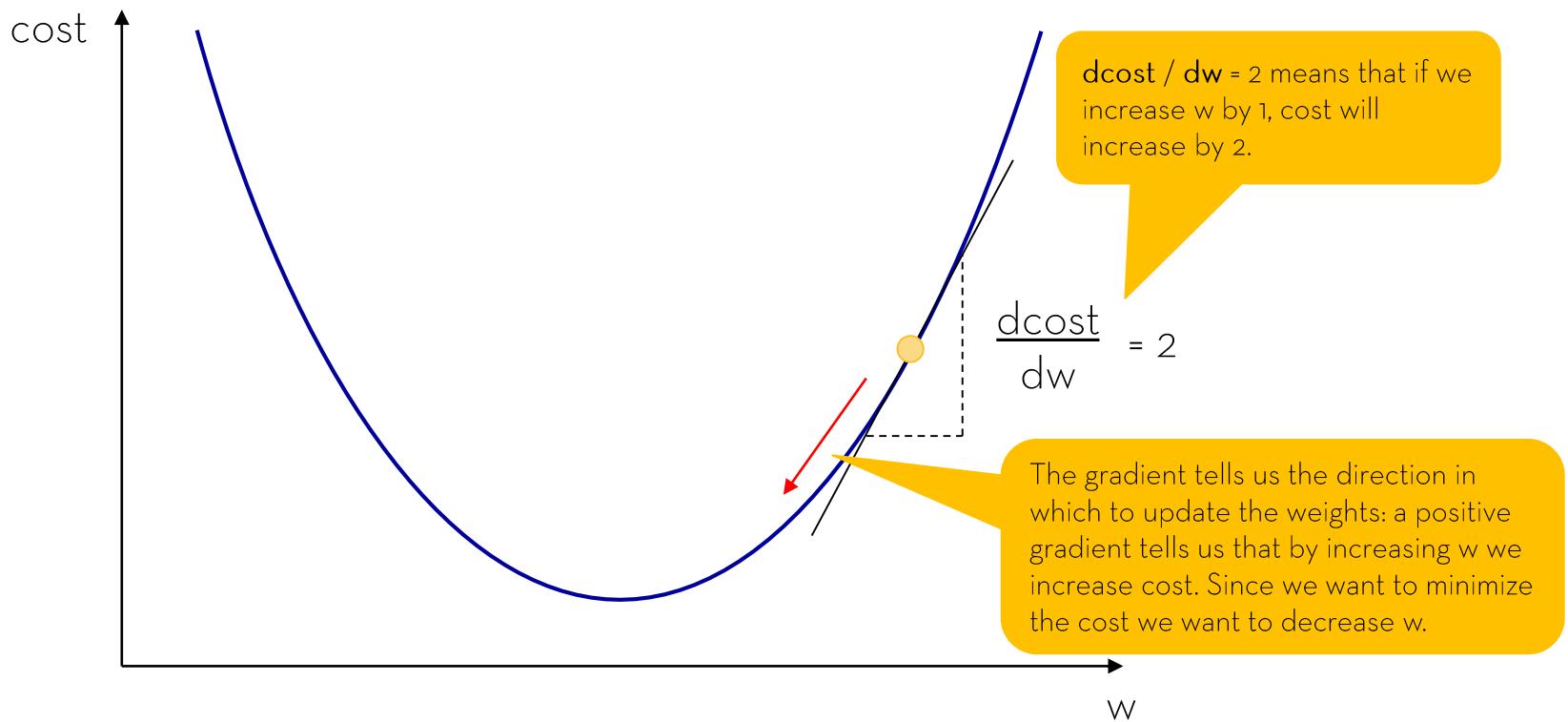
$$\frac{dy}{dx}$$

$$f'(x)$$

In most of the lectures we
will use function notation.
If $y = f(x)$ then $dy/dx = f'(x)$.

A quick refresher on gradients

66



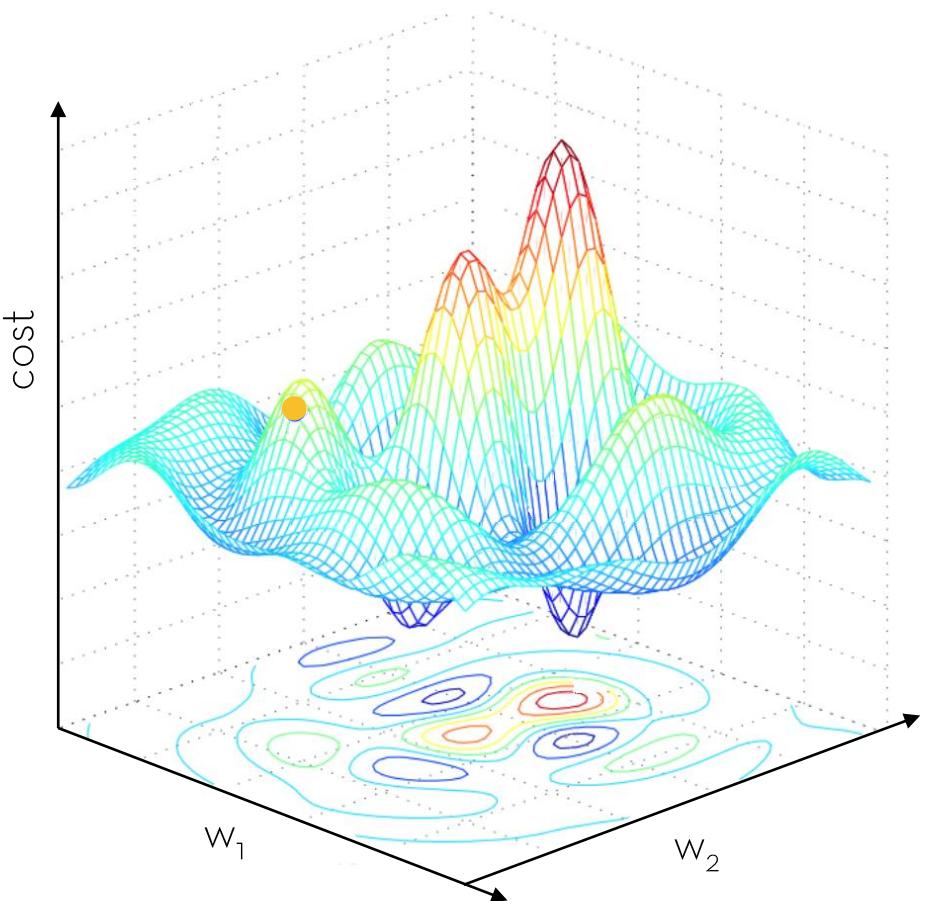
How to update the weights using gradients

67

Since we are optimizing many weights simultaneously, we take partial derivatives – i.e., the derivative of the cost with respect to one weight at a time.

$$\frac{\partial \text{cost}}{\partial w_1} \quad \frac{\partial \text{cost}}{\partial w_2}$$

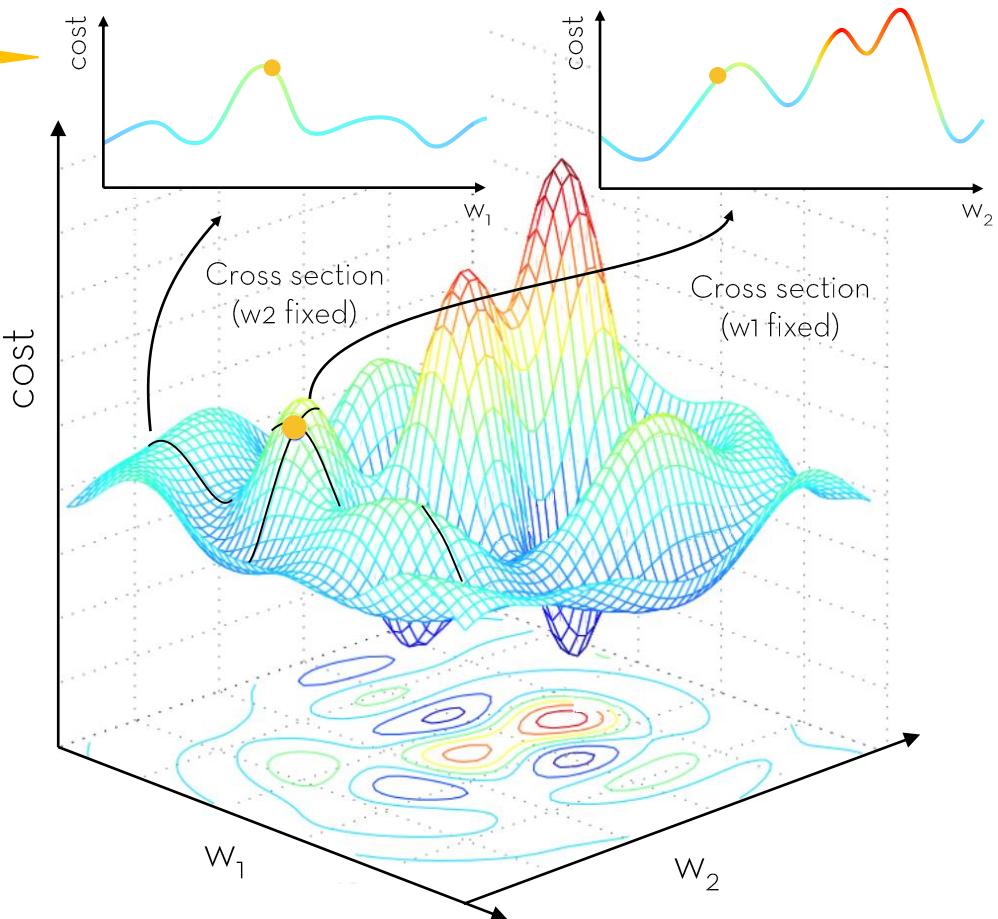
Then, once calculated we use $\partial \text{cost}/\partial w_1$ to update w_1 and $\partial \text{cost}/\partial w_2$ to update w_2 .



How to update the weights using gradients

We can take the cross sections of the loss function at the current point to visualize the gradients better.

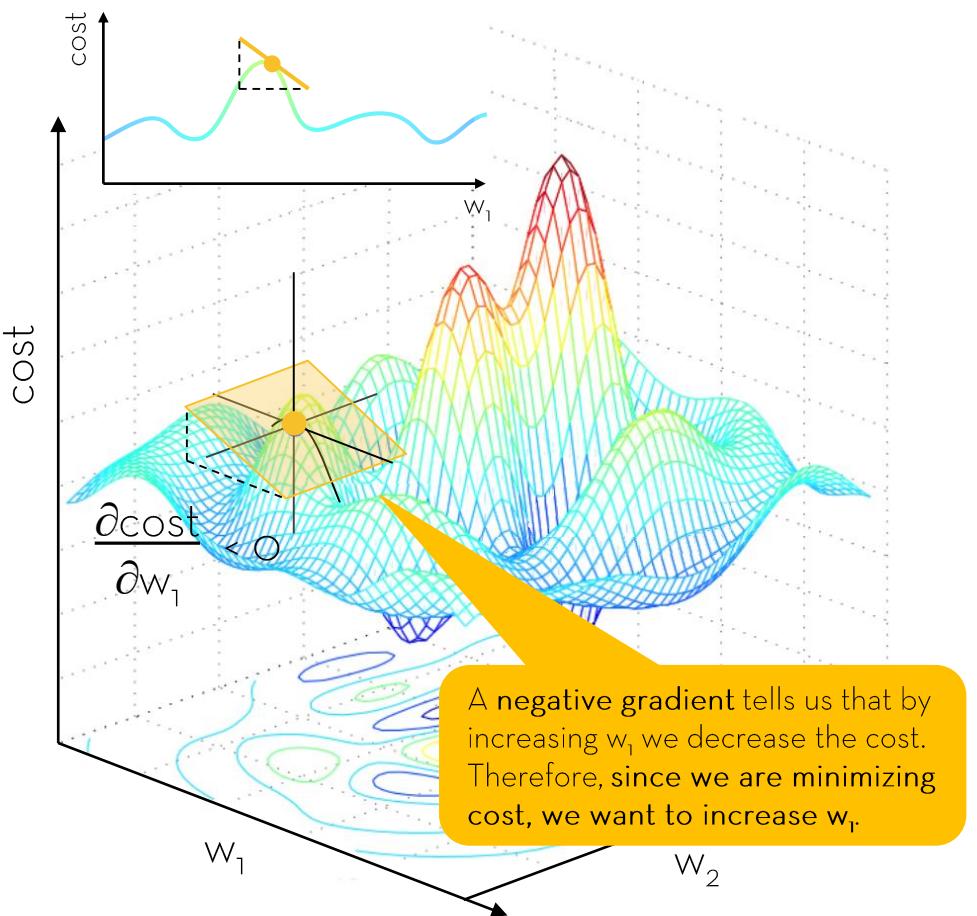
$$\frac{\partial \text{cost}}{\partial w_1} \quad \frac{\partial \text{cost}}{\partial w_2}$$



How to update the weights using gradients

- If $\frac{\partial \text{cost}}{\partial w_1} > 0$, decreasing the w_1 will decrease the cost.
- If $\frac{\partial \text{cost}}{\partial w_1} < 0$, increasing the w_1 will decrease the cost.

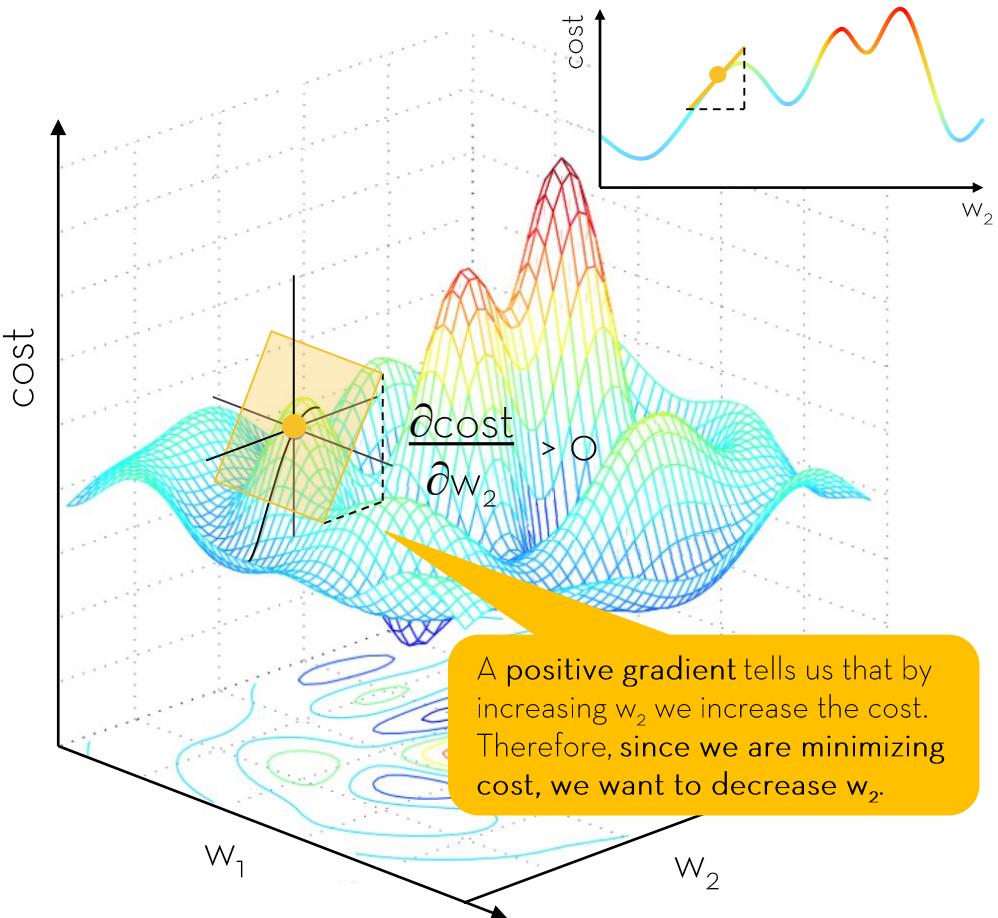
$$\frac{\partial \text{cost}}{\partial w_1} \quad \frac{\partial \text{cost}}{\partial w_2}$$



How to update the weights using gradients

$$\frac{\partial \text{cost}}{\partial w_1} \quad \frac{\partial \text{cost}}{\partial w_2}$$

- If $\frac{\partial \text{cost}}{\partial w_2} > 0$, decreasing the w_2 will decrease the cost.
- If $\frac{\partial \text{cost}}{\partial w_2} < 0$, increasing the w_2 will decrease the cost.



How to update the weights using gradients

Since $\partial \text{cost} / \partial w_1$ is negative we want to increase w_1 .

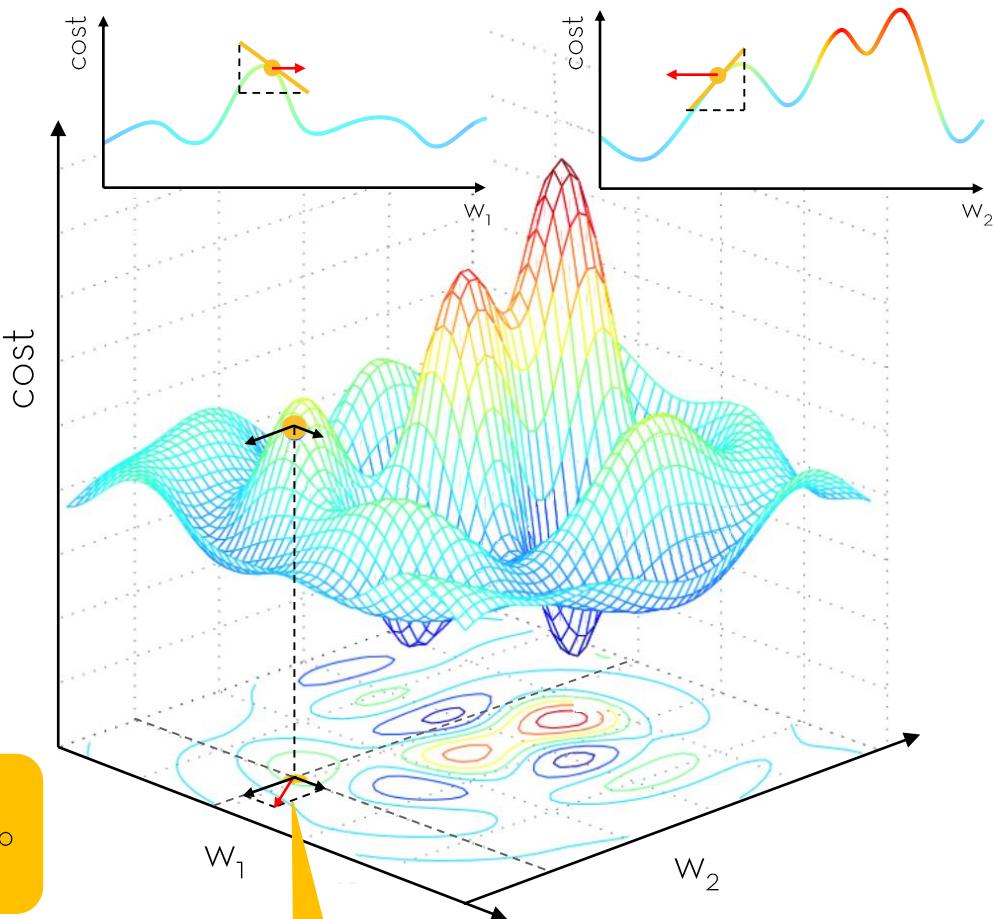
Since $\partial \text{cost} / \partial w_2$ is positive we want to decrease w_2 .

$$\frac{\partial \text{cost}}{\partial w_1}$$

$$\frac{\partial \text{cost}}{\partial w_2}$$

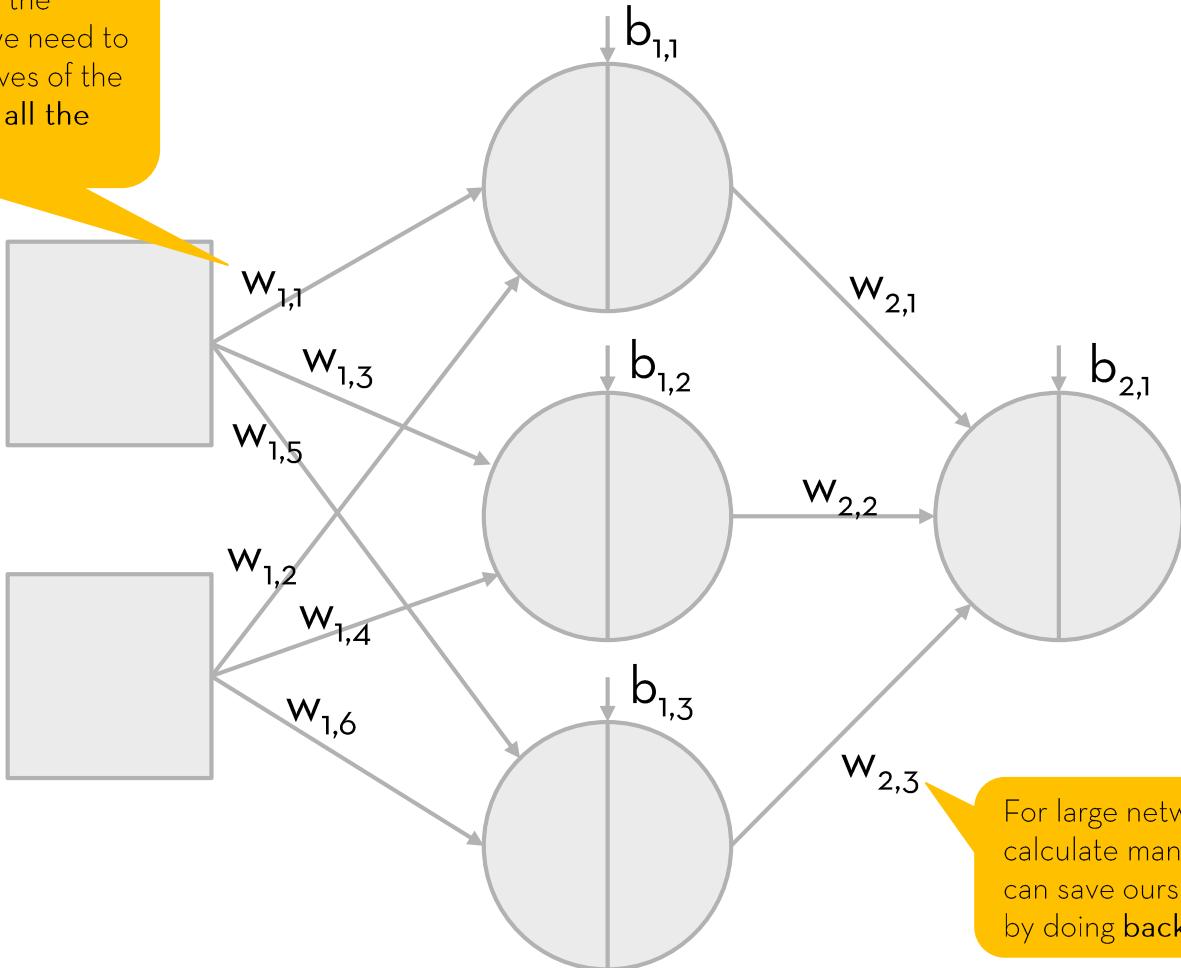
The gradient not only tells us which direction to travel but also how far we need to travel.

Since $\partial \text{cost} / \partial w_2$ is steeper than $\partial \text{cost} / \partial w_1$ we prefer to descend along w_2 since it will have the larger effect on the cost.



How to calculate the gradients

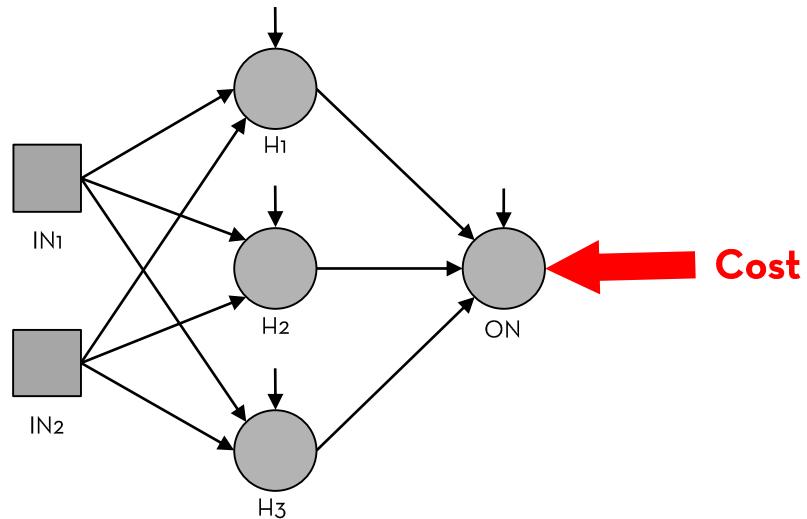
In order to update all the weights and biases, we need to calculate the derivatives of the cost with respect to all the weights and biases.



For large networks, we need to calculate many gradients! We can save ourselves some time by doing backpropagation.

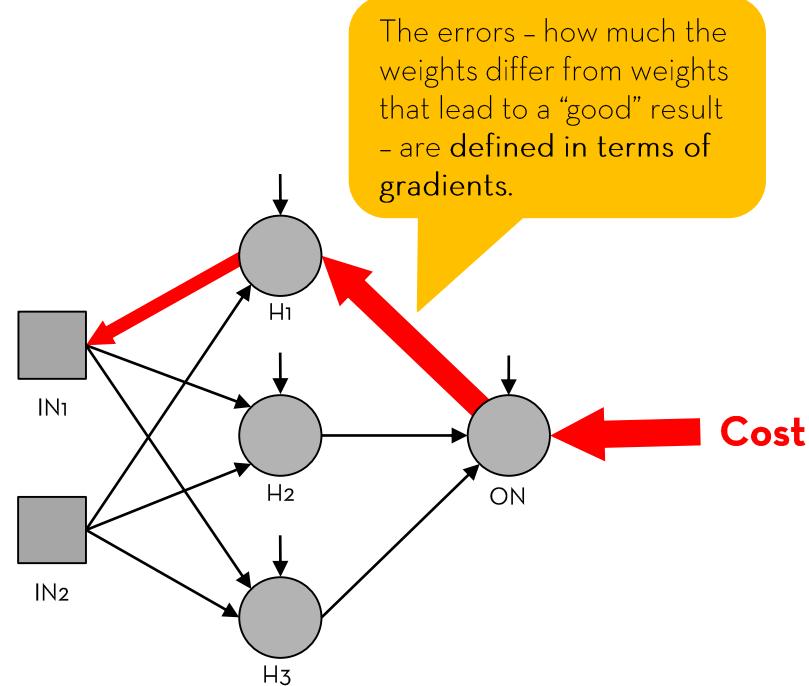
We use backpropagation to get the gradients of the network: Why back?

- We start at the end of the network – i.e., at the back.
- Then, we send the gradients backwards through the network.



We use backpropagation to get the gradients of the network: Why propagation?

- To calculate the gradient of earlier layers we propagate the gradients from later layers backwards through the network.
- Backpropagation saves us from having to recalculate the full gradient from scratch at each step.
- Formally, backpropagation is just the repeated application of the **chain rule**.



A quick refresher on the chain rule

75

A gradient can be decomposed using the chain rule.

We rephrase the question from

- How much does y change if we change x ?
to
1) How much does y change if we change z ? and
2) How much does z change if we change x ?

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

If $y = f(x)$ and $z = g(x)$.

$$= f'(z) g'(x)$$

See a primer in backpropagation for a more in-depth explanation.

We “pass the gradients backwards” by calculating the gradient of the cost with respect to earlier values

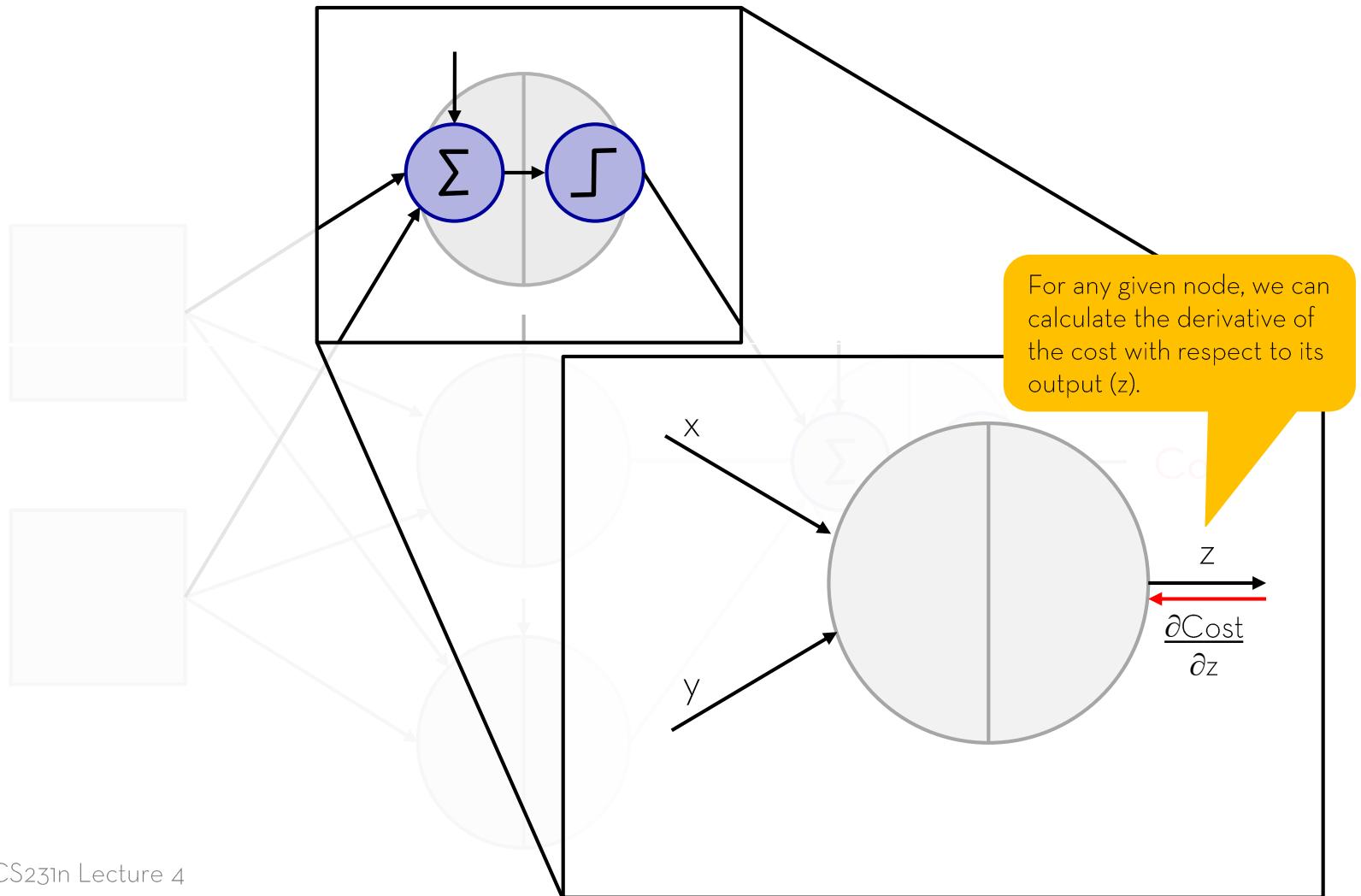


Figure adapted from CS231n Lecture 4

We “pass the gradients backwards” by calculating the gradient of the cost with respect to earlier values

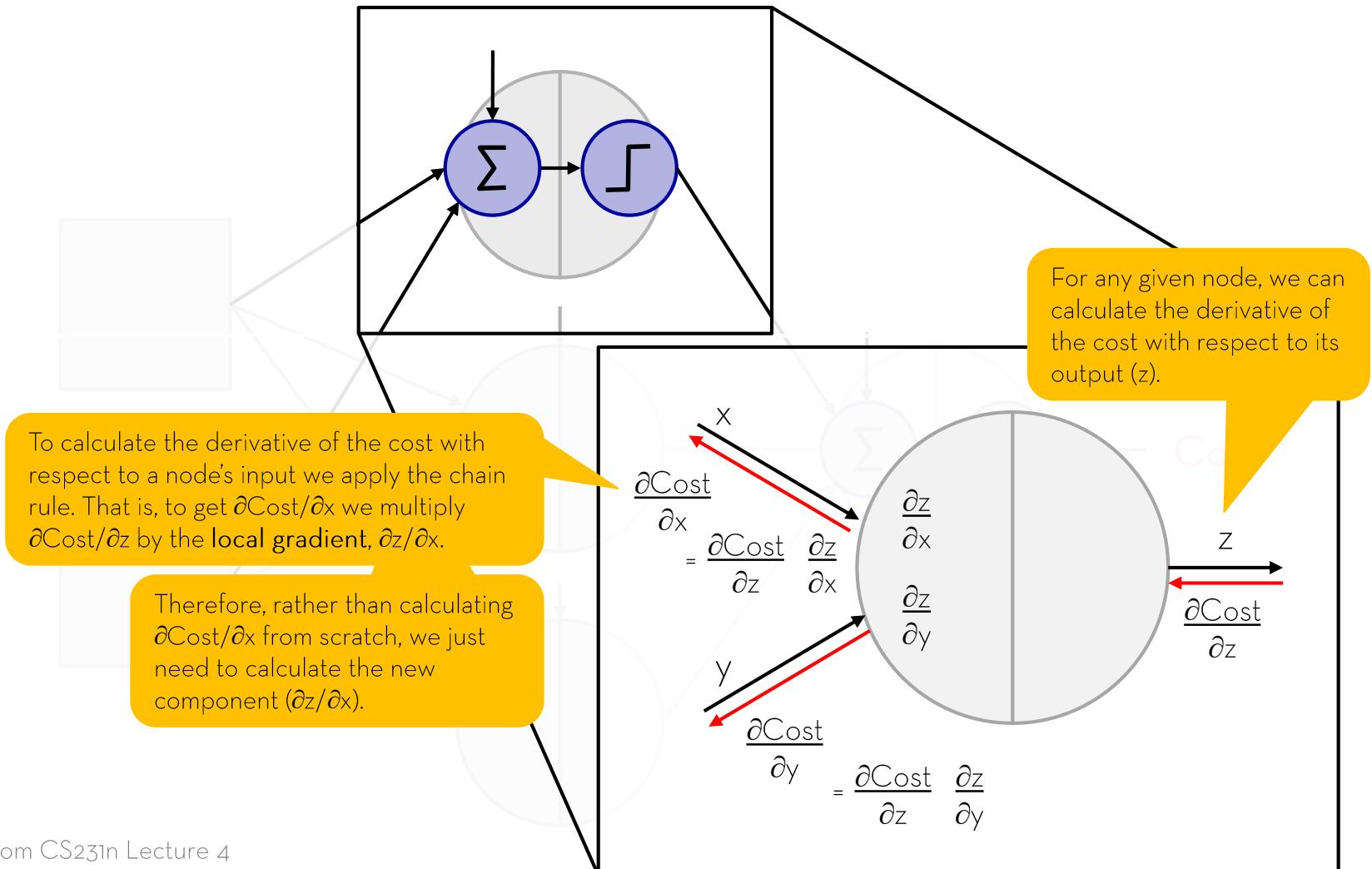
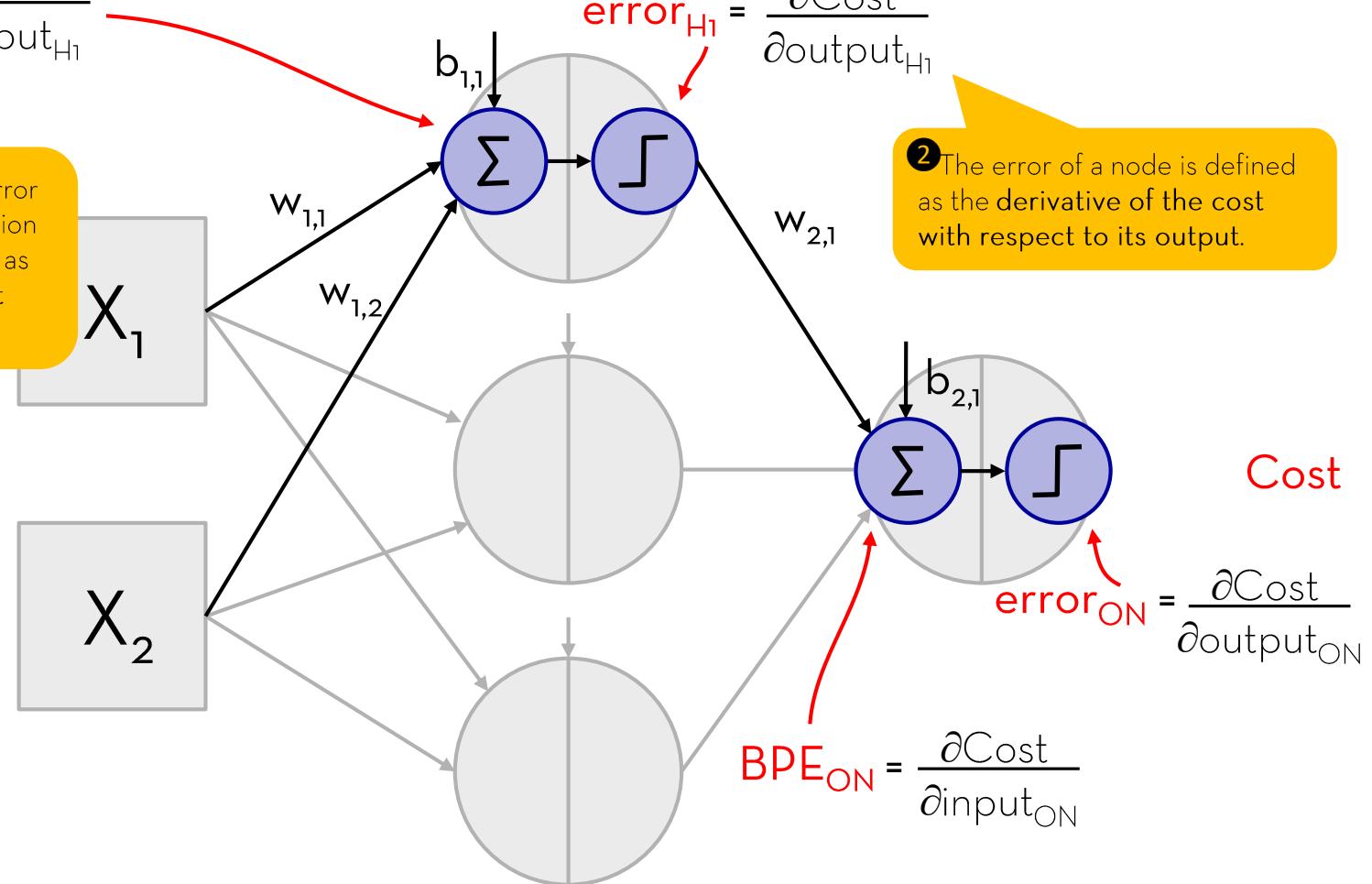


Figure adapted from CS231n Lecture 4

We can define the gradients in terms of errors

$$BPE_{H_1} = \frac{\partial \text{Cost}}{\partial \text{input}_{H_1}}$$

- ① The backpropagation error (BPE) (a.k.a. error information term) of a node is defined as the derivative of the cost with respect to its input.

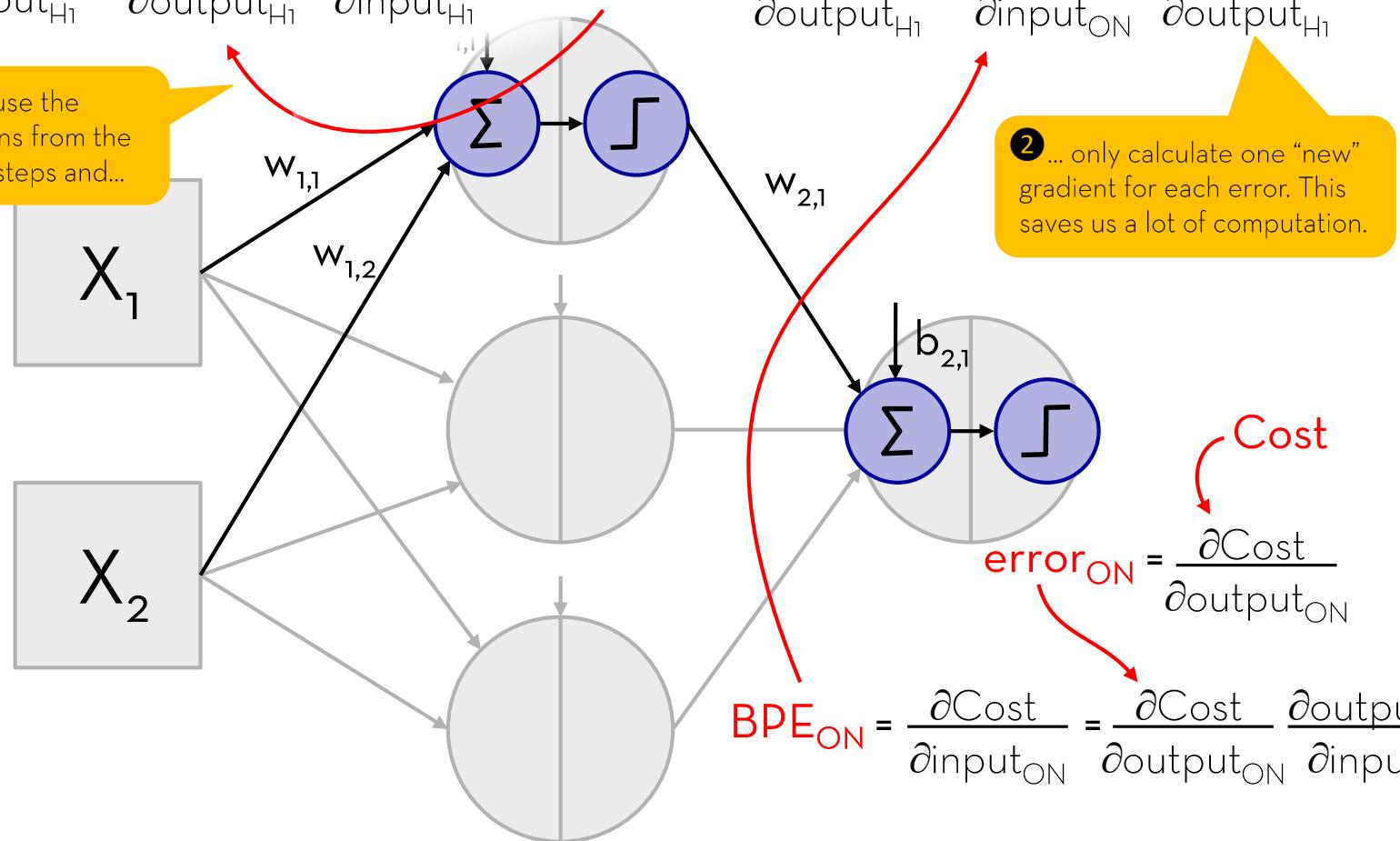


We can calculate all the important gradients using the chain rule

$$\text{BPE}_{H_1} = \frac{\partial \text{Cost}}{\partial \text{input}_{H_1}} = \frac{\partial \text{Cost}}{\partial \text{output}_{H_1}} \frac{\partial \text{output}_{H_1}}{\partial \text{input}_{H_1}}$$

$$\text{error}_{H_1} = \frac{\partial \text{Cost}}{\partial \text{output}_{H_1}} = \frac{\partial \text{Cost}}{\partial \text{input}_{ON}} \frac{\partial \text{input}_{ON}}{\partial \text{output}_{H_1}}$$

① We reuse the calculations from the previous steps and...



$$\text{BPE}_{ON} = \frac{\partial \text{Cost}}{\partial \text{input}_{ON}} = \frac{\partial \text{Cost}}{\partial \text{output}_{ON}} \frac{\partial \text{output}_{ON}}{\partial \text{input}_{ON}}$$

$$\text{error}_{ON} = \frac{\partial \text{Cost}}{\partial \text{output}_{ON}}$$

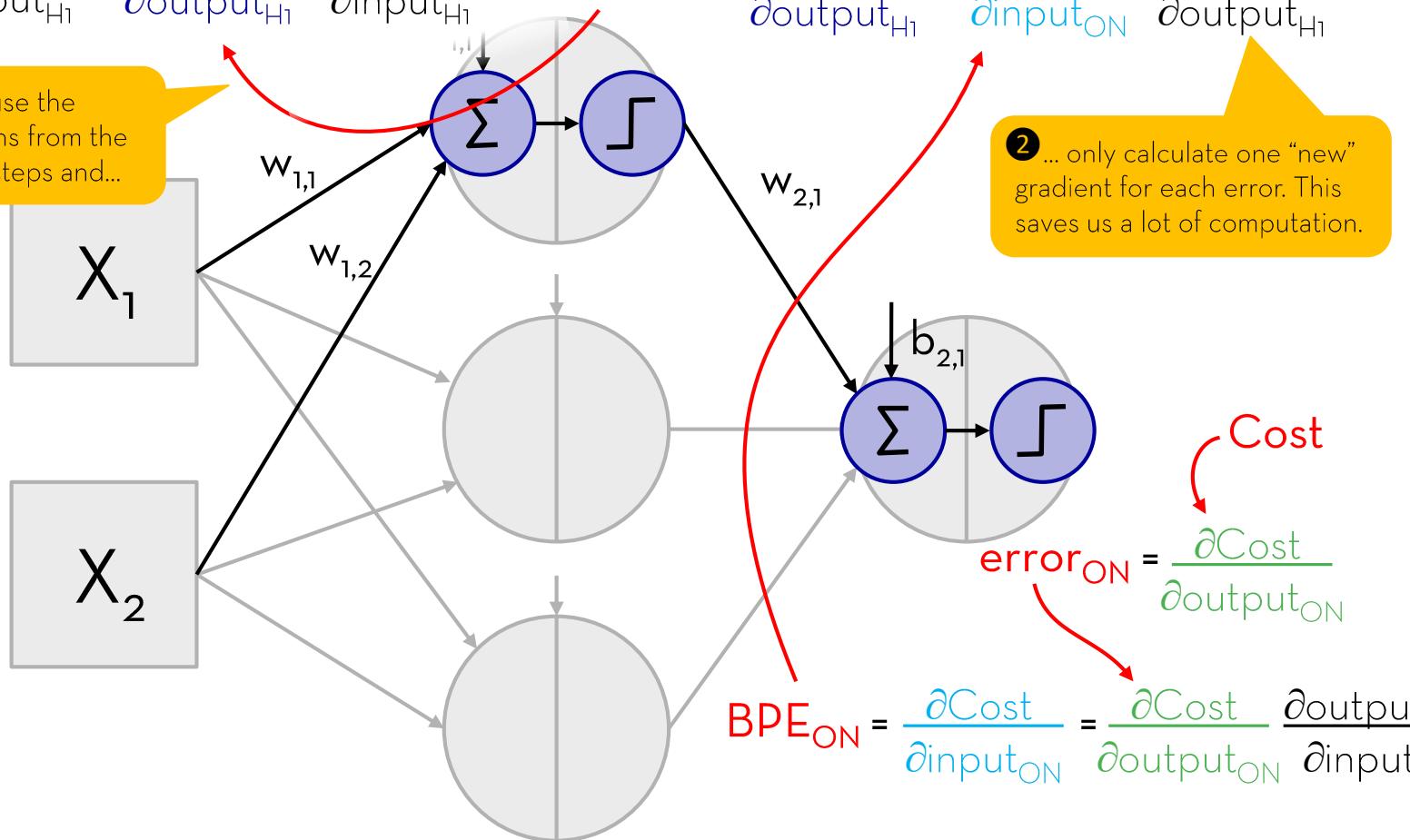
We can calculate all the important gradients using the chain rule

$$\text{BPE}_{H_1} = \frac{\partial \text{Cost}}{\partial \text{input}_{H_1}} = \frac{\partial \text{Cost}}{\partial \text{output}_{H_1}} \frac{\partial \text{output}_{H_1}}{\partial \text{input}_{H_1}}$$

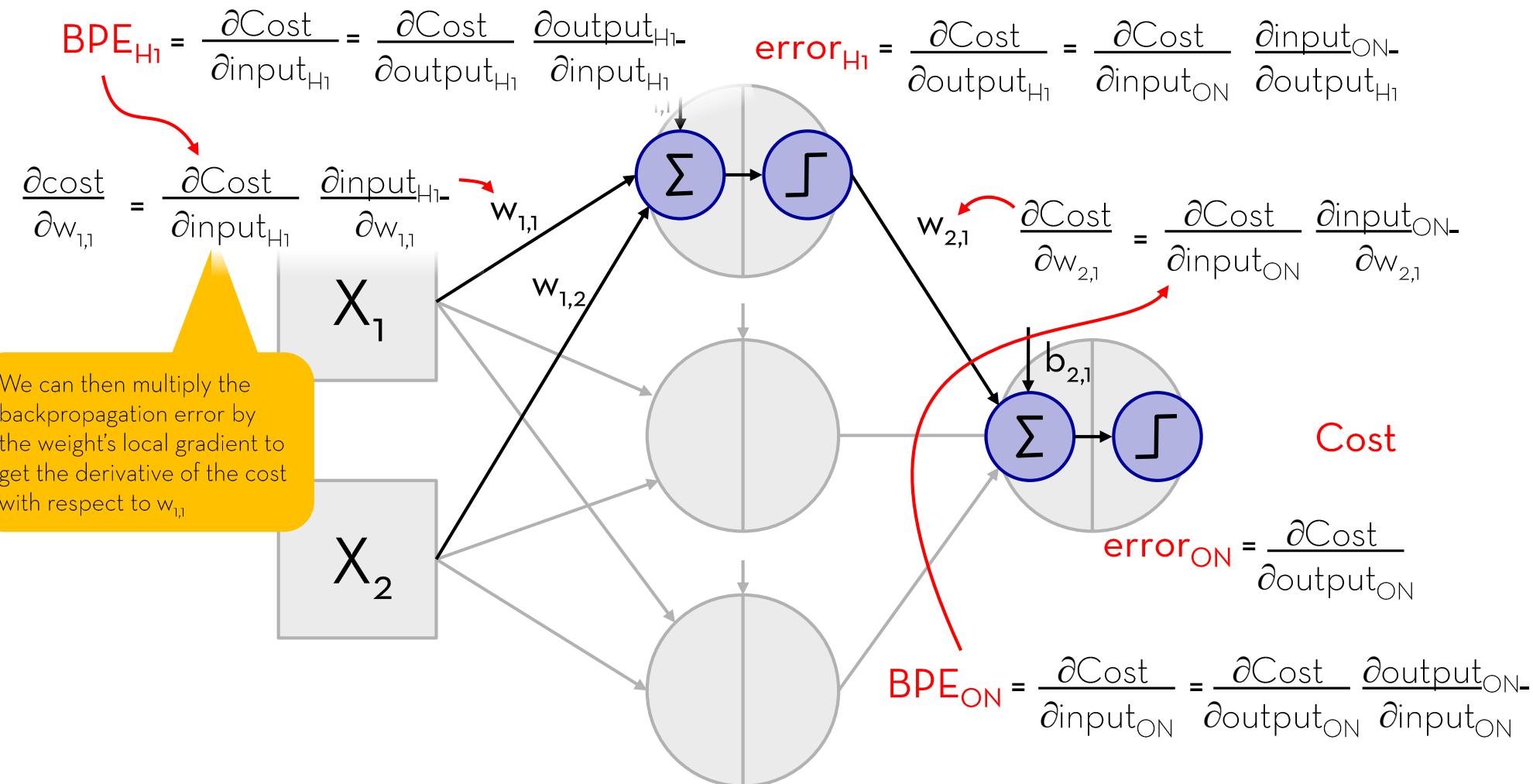
$$\text{error}_{H_1} = \frac{\partial \text{Cost}}{\partial \text{output}_{H_1}} = \frac{\partial \text{Cost}}{\partial \text{input}_{ON}} \frac{\partial \text{input}_{ON}}{\partial \text{output}_{H_1}}$$

1 We reuse the calculations from the previous steps and...

2 ... only calculate one “new” gradient for each error. This saves us a lot of computation.



We can calculate all the important gradients using the chain rule

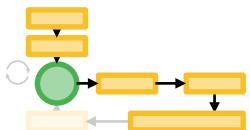


Steps of estimating neural networks

82

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

Step 5: Backpropagation

How to calculate the gradients?

a) Output layer:

1. Calculate error.
2. Calculate the backpropagation error.

b) Hidden layer:

1. Calculate error.
2. Calculate the backpropagation error.



We work our way
backward through the
network.

Step 5: Backpropagation

How to calculate the gradients?

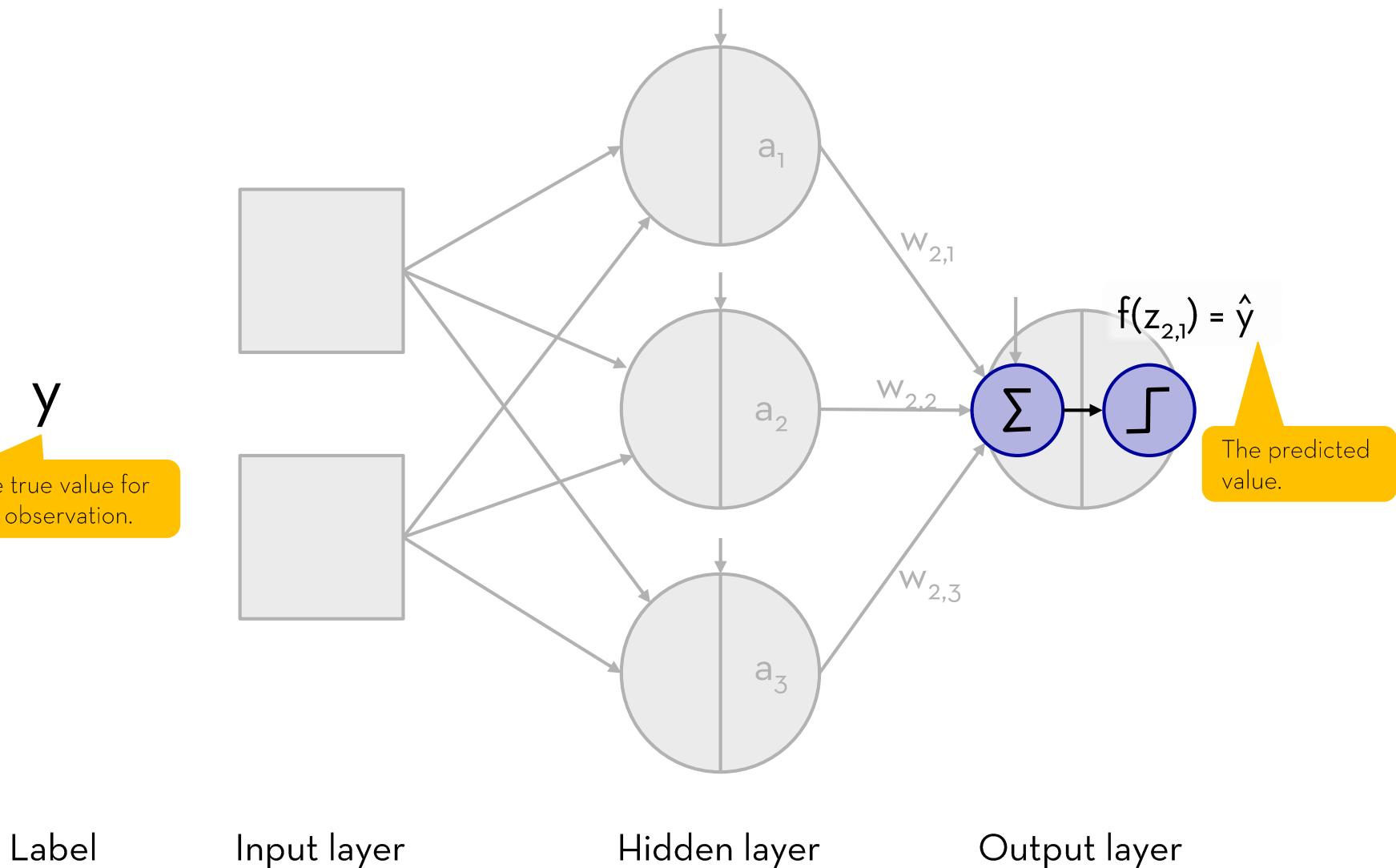
84

- a) Output layer:
 - 1. Calculate error.
 - 2. Calculate the backpropagation error.
- b) Hidden layer:
 - 1. Calculate error.
 - 2. Calculate the backpropagation error.

Step 5: Backpropagation

a1) Calculate the error for the output layer (1/2)

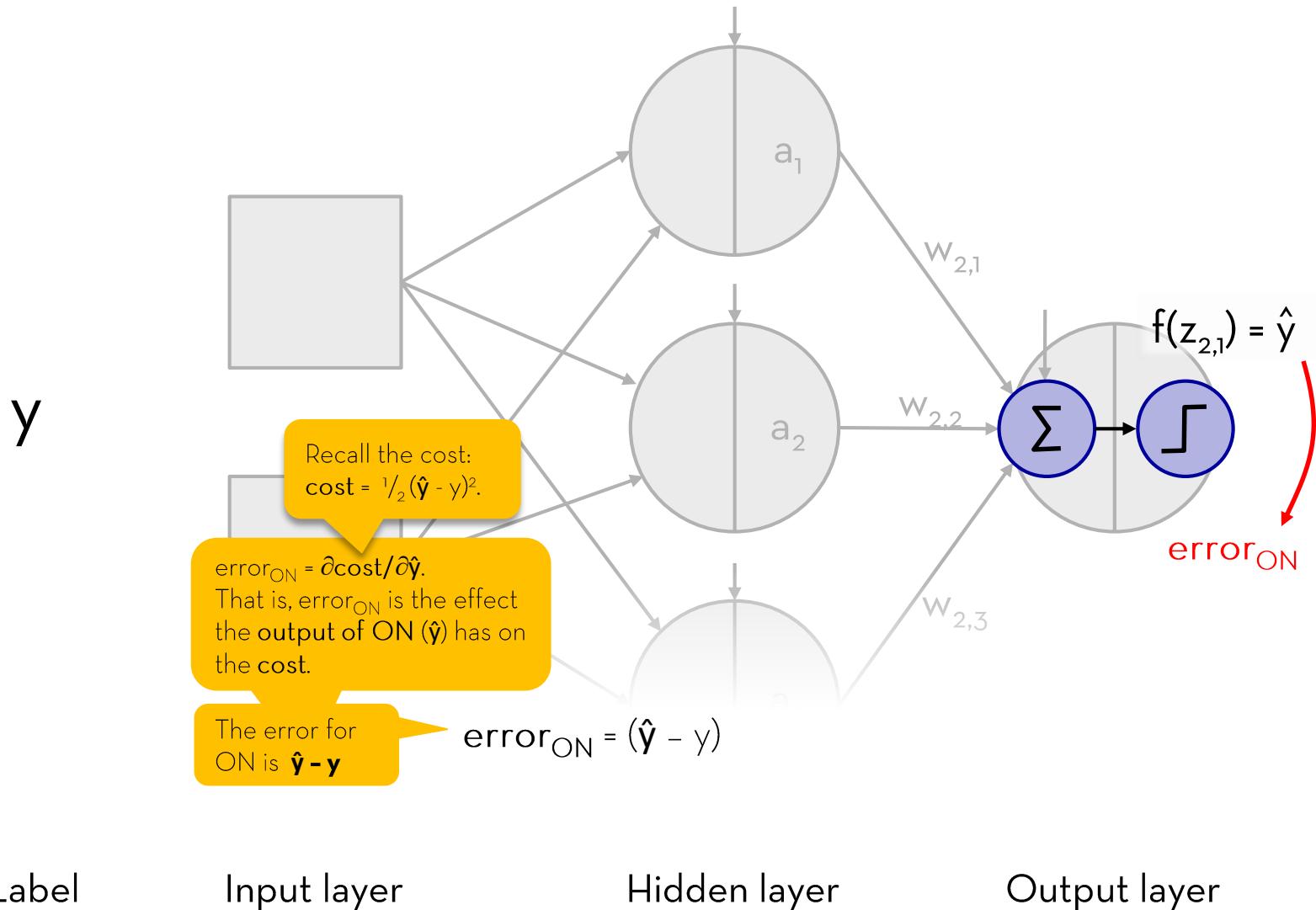
85



Step 5: Backpropagation

a1) Calculate the error for the output layer (2/2)

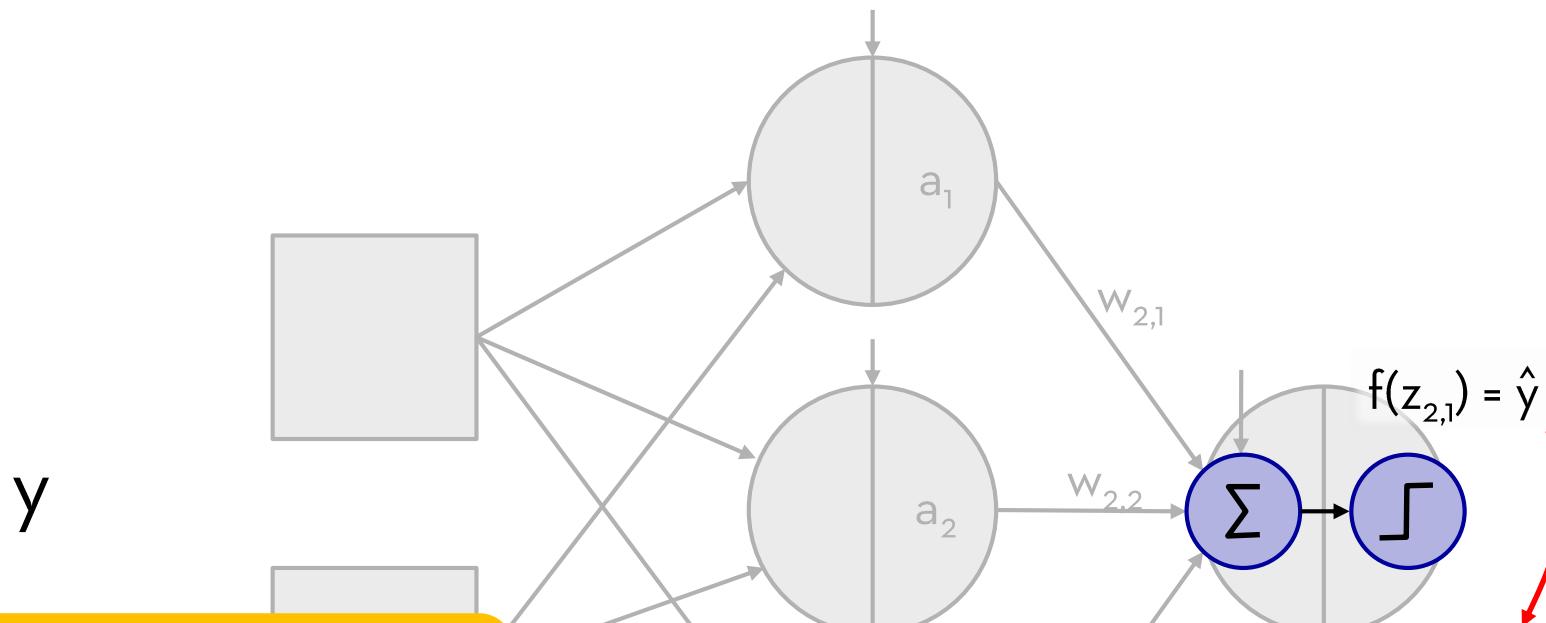
86



Step 5: Backpropagation

a) Calculate the backpropagation error for ON (1/2)

87



The backpropagation error is defined as the error of the node multiplied by the derivative of the activation of the node at the input of the node.

backprop error_{ON} = $\partial \text{cost} / \partial \text{input}_{\text{ON}}$. That is, backprop error_{ON} is the effect of the input of ON on the cost.

$$\text{error}_{\text{ON}} = (\hat{y} - y)$$

$$\text{backprop error}_{\text{ON}} = \text{error}_{\text{ON}} * f'(\text{input}_{\text{ON}})$$

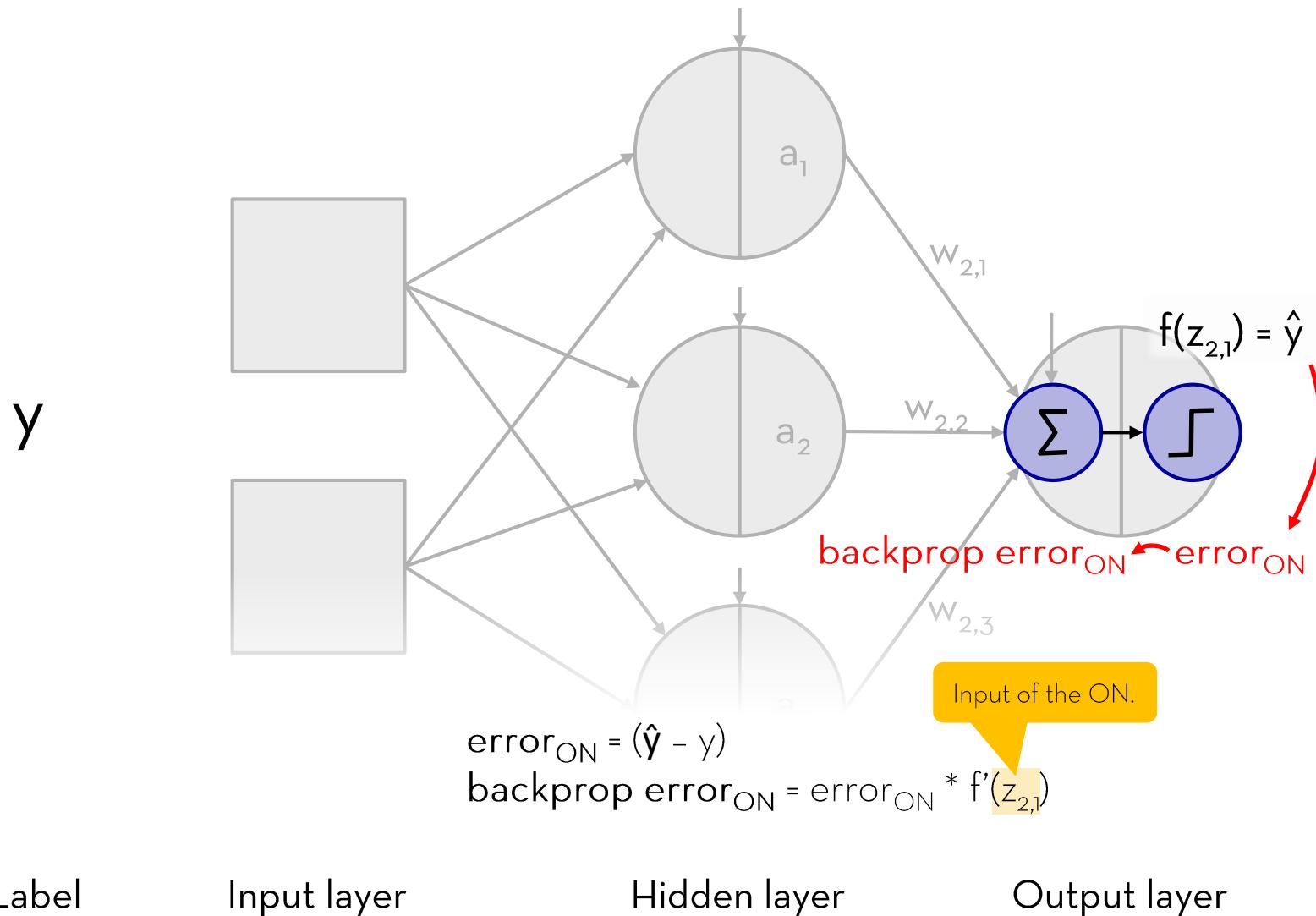
$f'(\cdot)$ is the derivative of the sigmoid function.

$$f(x) = 1 / (1 + e^{-x}), \\ f'(x) = f(x)(1 - f(x)).$$

Step 5: Backpropagation

a) Calculate the backpropagation error for ON (2/2)

88

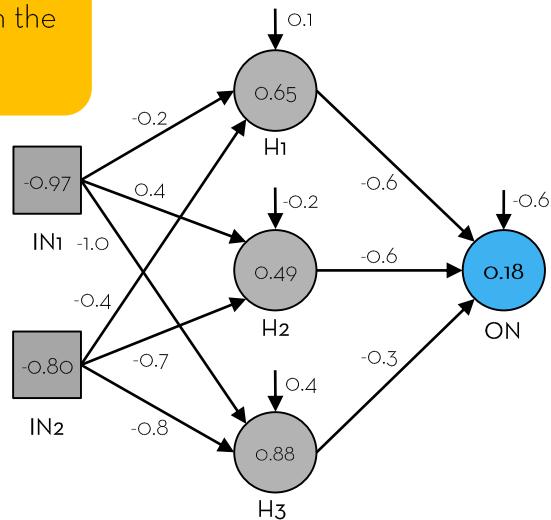


Step 5: Backpropagation

a1) Calculate the error for output node (1/2)

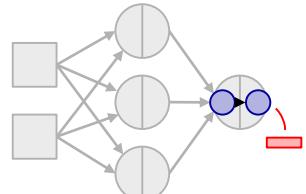
Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

The **error** is the difference between the estimated and true value: $\hat{y} - y$.



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6



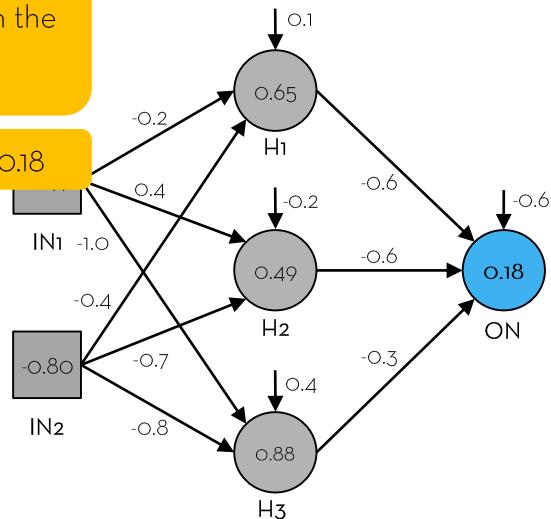
Step 5: Backpropagation

a1) Calculate the error for output node (2/2)

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

The **error** is the difference between the estimated and true value: $\hat{y} - y$.

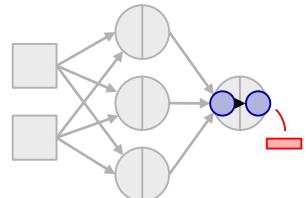
$$\text{error}_{\text{ON}} = 0.18 - 0 = 0.18$$



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error				0.18

$$\text{error}_{\text{ON}} = \partial \text{cost} / \partial \hat{y}$$



Step 5: Backpropagation

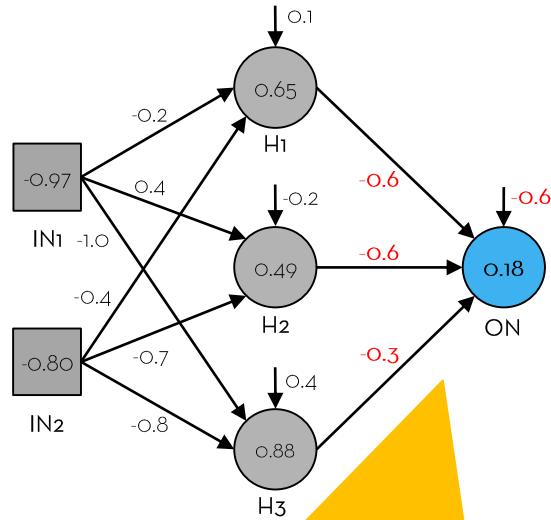
a) Calculate the backpropagation error for ON

91

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error				0.18
Backprop. Err				0.025

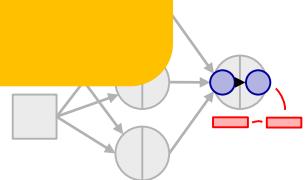


The **backpropagation error for ON**:

backprop error_{ON}

$$\begin{aligned}
 &= \text{error}_{\text{ON}} * f'(a_1w_1 + a_2w_2 + a_3w_3 + b) \\
 &= \text{error}_{\text{ON}} * f'(-0.6*0.65 + -0.6*0.49 + -0.3*0.88 + -0.6) \\
 &= 0.18 * f'(-1.55) \\
 &= 0.025
 \end{aligned}$$

backprop error_{ON} = $\partial \text{cost} / \partial \text{input}_{\text{ON}}$



Step 5: Backpropagation

How to calculate the gradients?

a) Output layer:

1. Calculate error.
2. Calculate the backpropagation error.

b) Hidden layer:

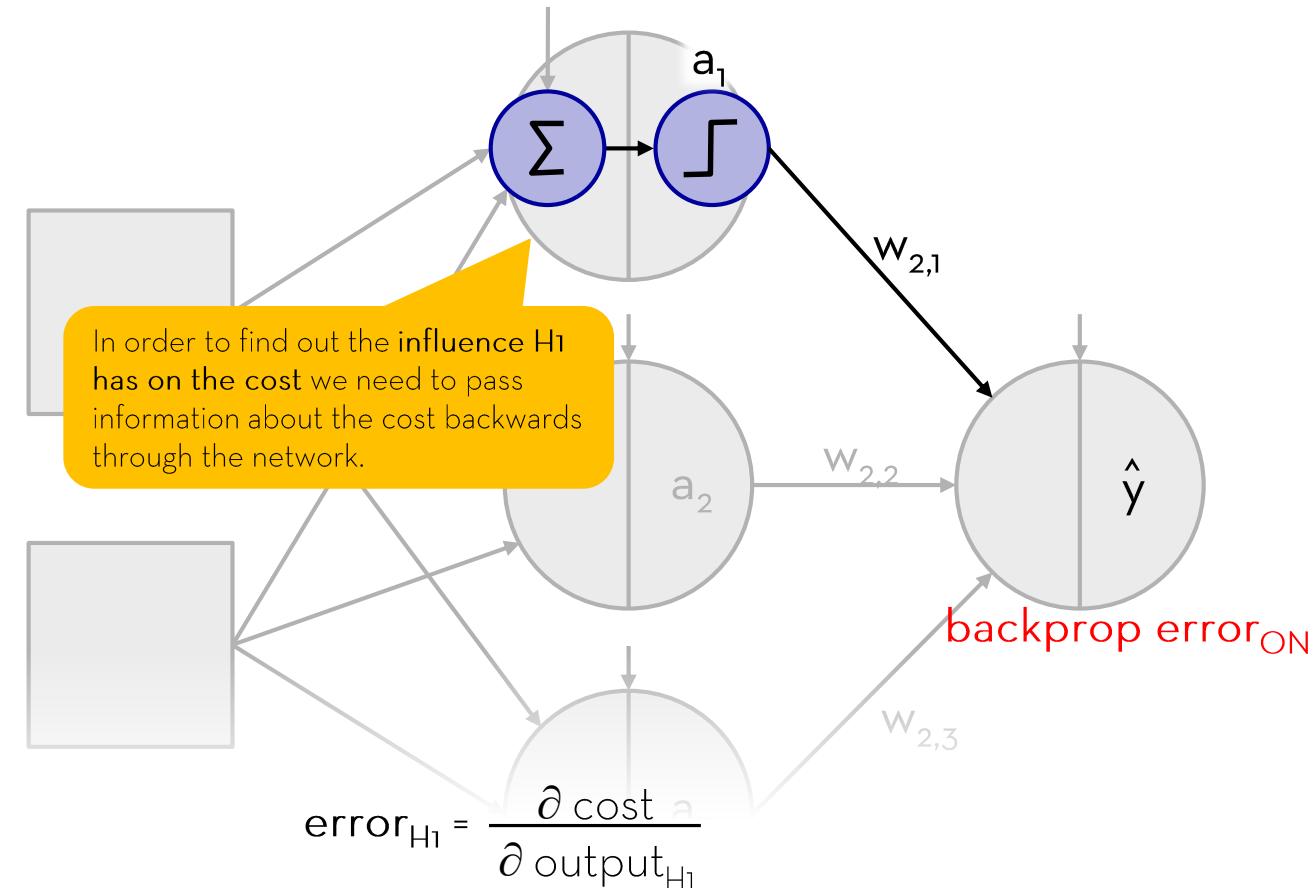
1. Calculate error.
2. Calculate the backpropagation error.

Step 5: Backpropagation

b1) Calculate the error for the hidden layer ($H_{1,1}$) (1/4)

93

y



Label

Input layer

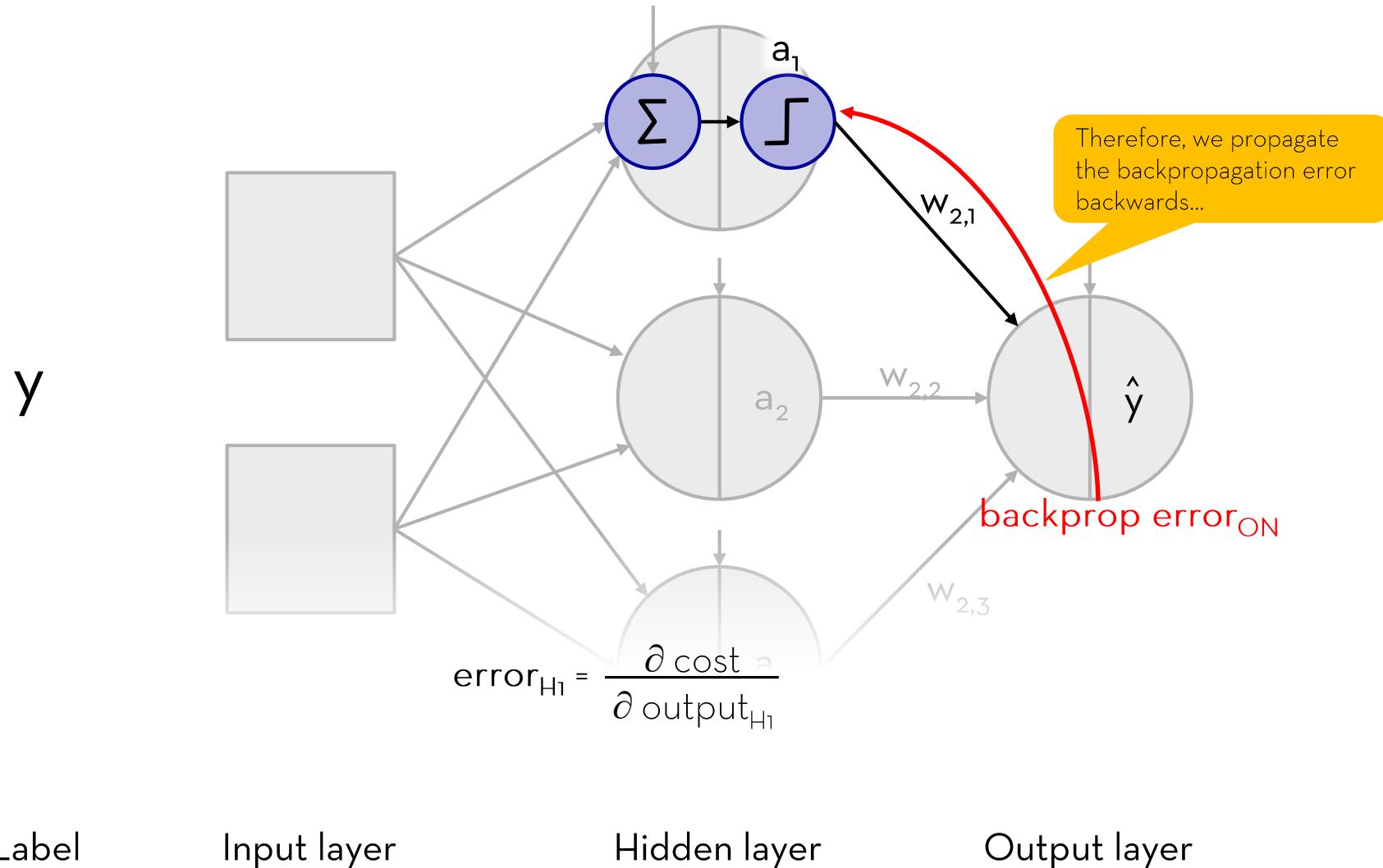
Hidden layer

Output layer

Step 5: Backpropagation

b1) Calculate the error for the hidden layer ($H_{1,1}$) (2/4)

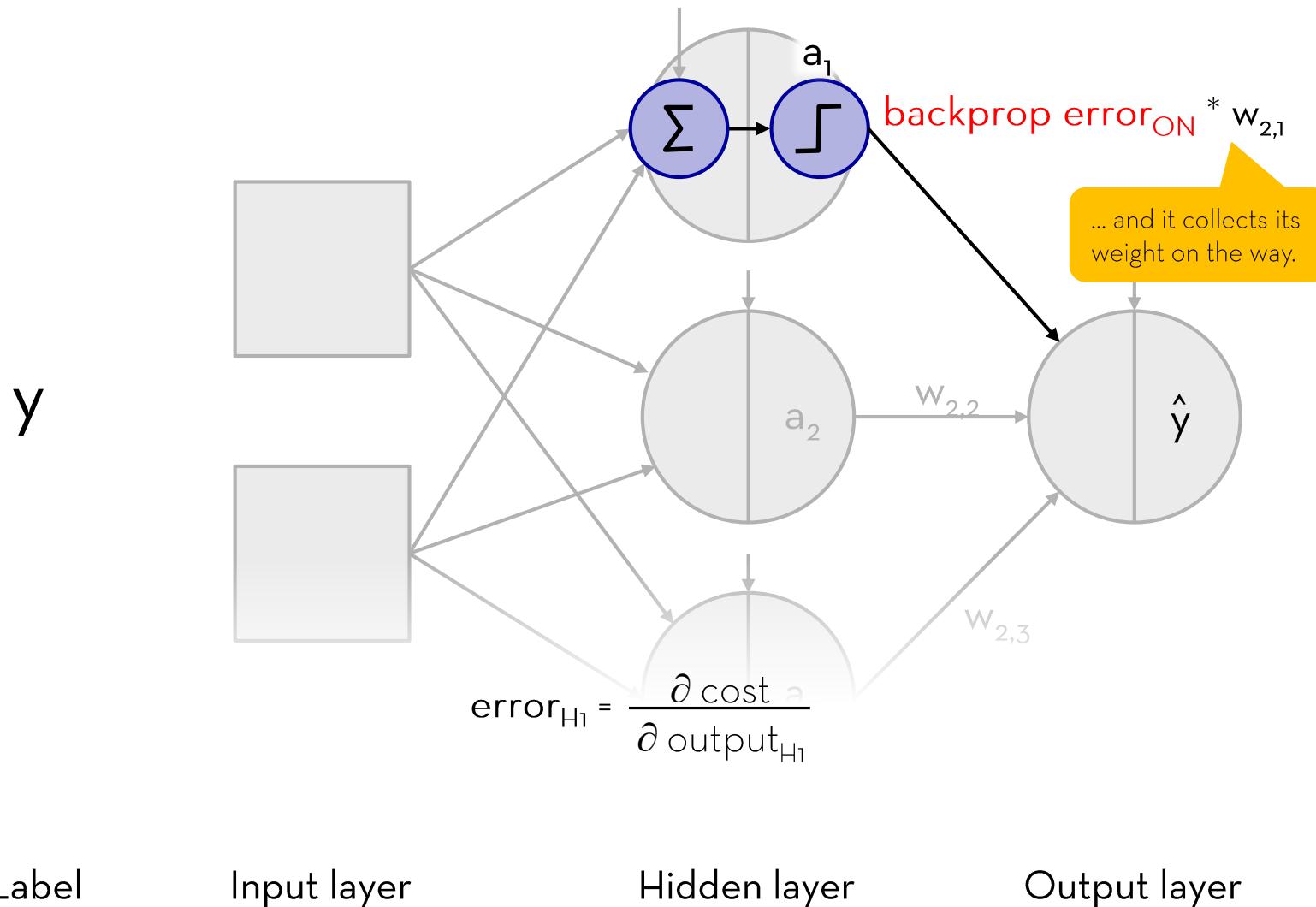
94



Step 5: Backpropagation

b1) Calculate the error for the hidden layer ($H_{1,1}$) (3/4)

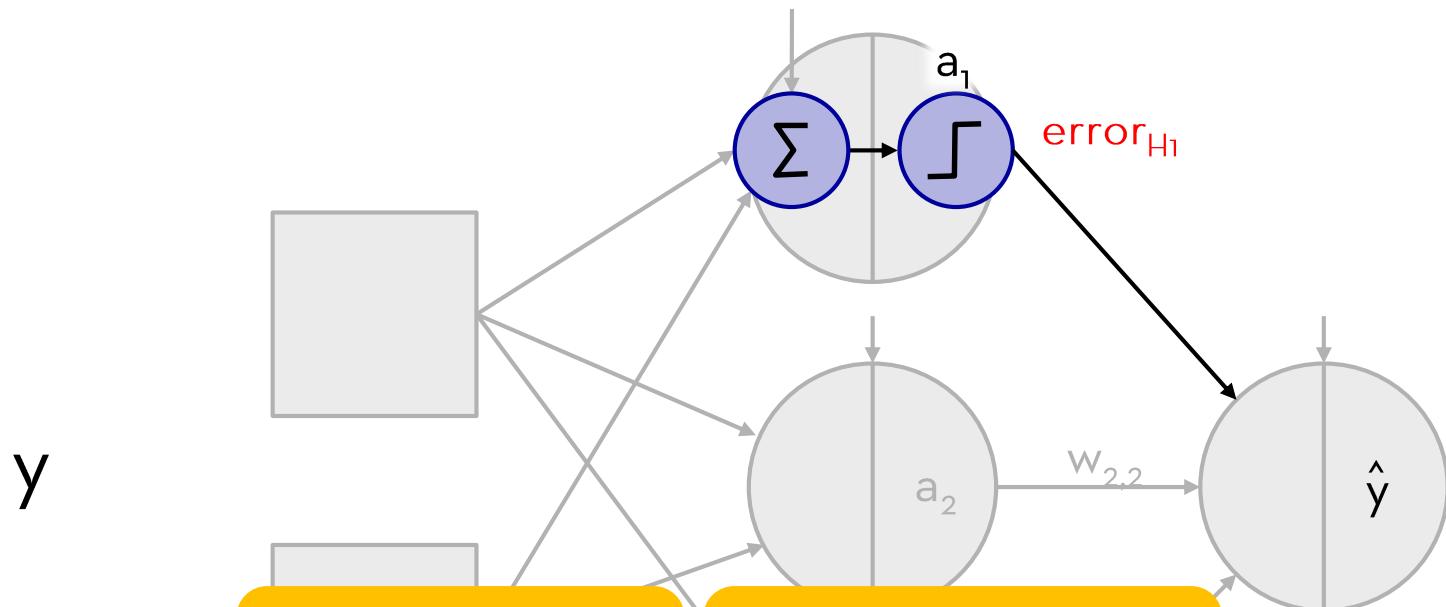
95



Step 5: Backpropagation

b1) Calculate the error for the hidden layer ($H_{1,1}$) (4/4)

96



$\text{error}_{H_1} = \frac{\partial \text{cost}}{\partial a_1}$
That is, error_{H_1} is the effect
the output of H_1 (a_1) has on
the cost.

Why? Because by the chain rule

$$\frac{\partial \text{cost}}{\partial a_1} = \frac{\partial \text{cost}}{\partial \text{input}_{ON}} * \frac{\partial \text{input}_{ON}}{\partial a_1}$$

Recall:
 $\text{input}_{ON} = a_1 w_{2,1} + a_2 w_{2,2} + a_3 w_{2,3} + b_{2,1}$

$$\text{error}_{H_1} = \text{backprop error}_{ON} * w_{2,1}$$

Label

Input layer

Hidden layer

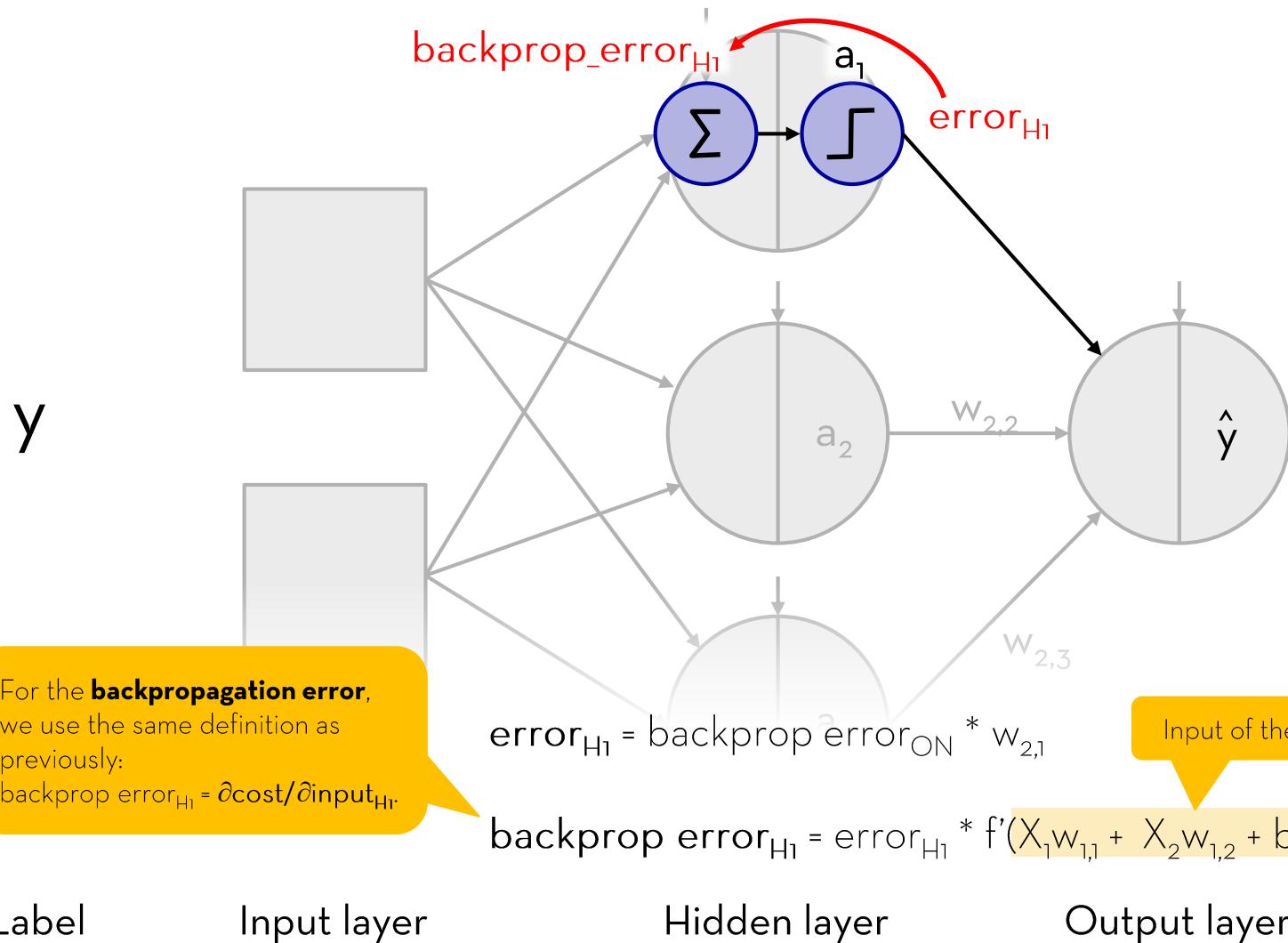
Output layer

See a primer in
backpropagation
for a more in-depth
explanation.

Step 5: Backpropagation

b2) Calculate the backpropagation error for the HL

97



Step 5: Backpropagation

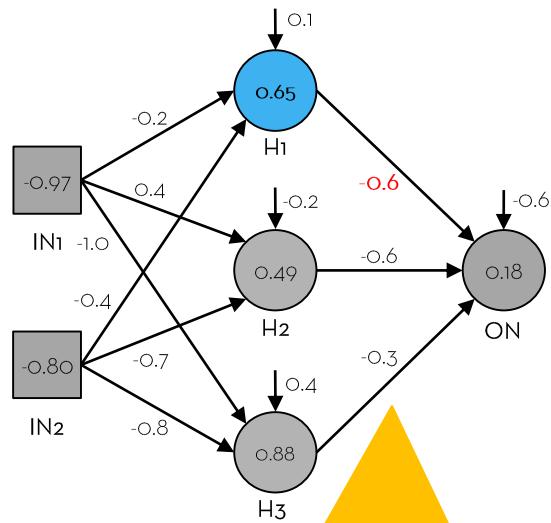
b1) Calculate error for the HL nodes (1/3)

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

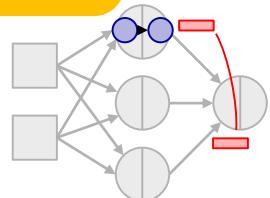
	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015			0.18
Backprop. Err				0.025

$$\text{error}_{H1} = \partial \text{cost} / \partial \text{output}_{H1}$$



Compute the error:

$$\begin{aligned}
 \text{error}_{H1} &= \text{backprop error}_{ON} * w_{H1 \rightarrow ON} \\
 &= 0.025 * -0.6 \\
 &= -0.015.
 \end{aligned}$$



Step 5: Backpropagation

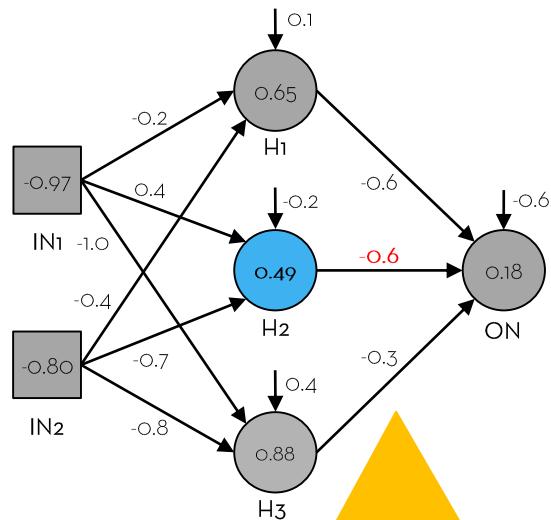
b1) Calculate error for the HL nodes (2/3)

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

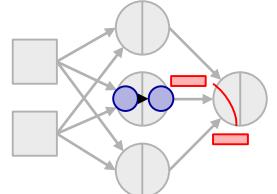
	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015		0.18
Backprop. Err				0.025

$$\text{error}_{H2} = \frac{\partial \text{cost}}{\partial \text{output}_{H2}}$$



$\text{error}_{H2} = \text{backprop error}_{ON} * w_{H2 \rightarrow ON}$
 $= 0.025 * -0.6$
 $= -0.015.$



Step 5: Backpropagation

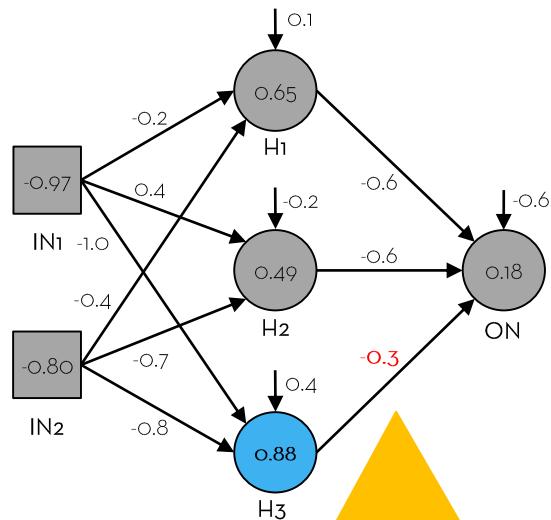
b1) Calculate error for the HL nodes (3/3)

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

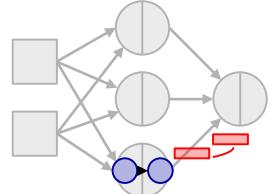
	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err				0.025

$$\text{error}_{H3} = \partial \text{cost} / \partial \text{output}_{H3}$$



$$\begin{aligned}\text{error}_{H3} &= \text{backprop error}_{ON} * w_{H3 \rightarrow ON} \\ &= 0.025 * -0.3 \\ &= -0.008\end{aligned}$$



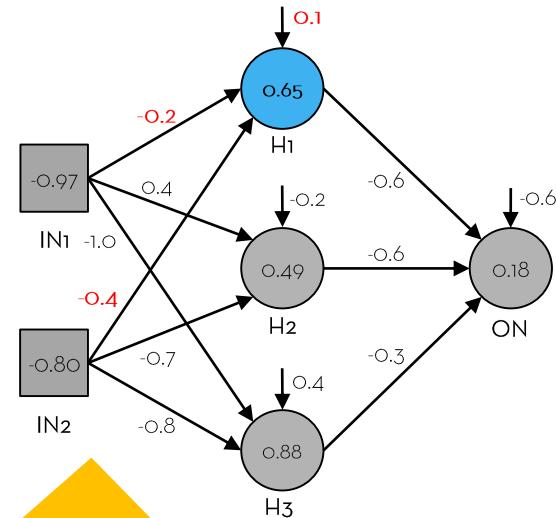
Step 5: Backpropagation

b2) Calculate backpropagation error for HL nodes

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003			0.025



Then, calculate the **backpropagation error**:

backprop error_{H1}

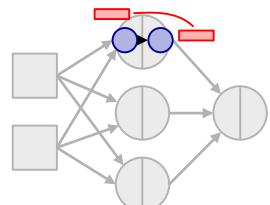
$$= \text{error}_{H1} * f'(\text{input}_{H1})$$

$$= \text{error}_{H1} * f'(-0.97 * -0.2 + -0.80 * -0.4 + 0.1)$$

$$= -0.015 * f'(0.614)$$

$$= -0.003$$

backprop error_{H1} = $\frac{\partial \text{cost}}{\partial \text{input}_{H1}}$



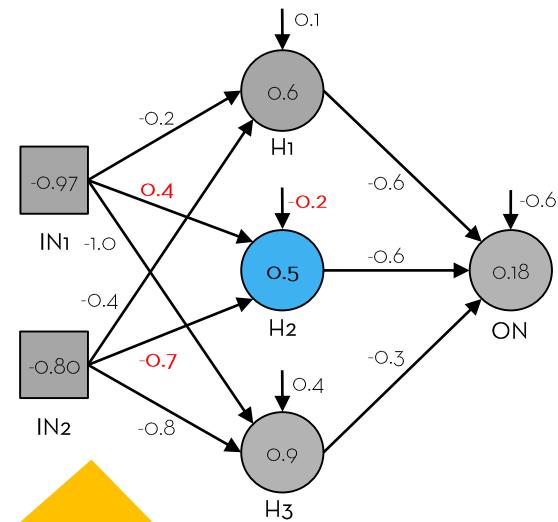
Step 5: Backpropagation

b2) Calculate backpropagation error for HL nodes

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

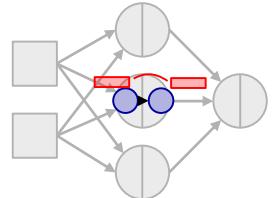
	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004		0.025



Then, calculate the **backpropagation error**:
 $\text{backprop error}_{H2}$

$$\begin{aligned}
 &= \text{error}_{H2} * f'(\text{input}_{H2}) \\
 &= \text{error}_{H2} * f'(-0.97 * 0.4 + -0.80 * -0.7 + -0.2) \\
 &= -0.015 * f'(-0.026) \\
 &= -0.004
 \end{aligned}$$

$\text{backprop error}_{H2} = \frac{\partial \text{cost}}{\partial \text{input}_{H2}}$



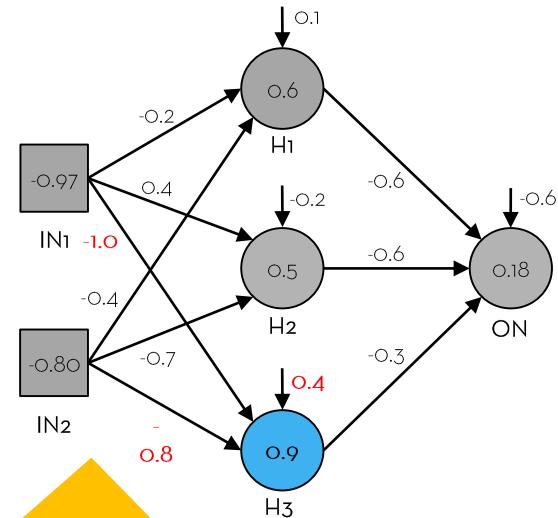
Step 5: Backpropagation

b2) Calculate backpropagation error for HL nodes

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

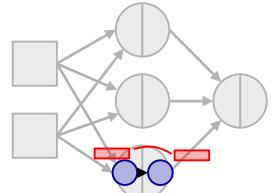
	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



Then, calculate the **backpropagation error**:
 $\text{backprop error}_{H3}$

$$\begin{aligned}
 &= \text{error}_{H3} * f'(\text{input}_{H3}) \\
 &= \text{error}_{H3} * f'(-0.97 * -1.0 + -0.80 * -0.7 + 0.4) \\
 &= -0.008 * f'(2.01) \\
 &= -0.001
 \end{aligned}$$

$\text{backprop error}_{H3} = \frac{\partial \text{cost}}{\partial \text{input}_{H3}}$



Step 5: Backpropagation

Repeat steps a) and b) for all observations

Y	X1	X2	H1	H2	H3	Output	Cost	Err H1	Err H2	Err H3	Err ON	BPE H1	BPE H2	BPE H3	BPE ON
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.015	-0.015	-0.008	0.18	-0.003	-0.004	-0.001	0.025
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03								
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03								
...							
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67								
O	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>								
							$\Sigma = 1.55$								

So far, we have calculated the errors and backpropagation errors for the output and hidden layer for the first observation.

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

Step 5: Backpropagation

Repeat steps a) and b) for all observations

105

Y	X1	X2	H1	H2	H3	Output	Cost	Err H1	Err H2	Err H3	Err ON	BPE H1	BPE H2	BPE H3	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.015	-0.015	-0.008	0.18	-0.003	-0.004	-0.001	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.014	-0.014	-0.007	0.17	-0.003	-0.004	-0.001	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.014	-0.014	-0.007	0.17	-0.003	-0.004	-0.001	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.073	0.073	0.037	-0.82	0.017	0.018	0.005	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	-0.023	-0.023	-0.012	0.22	-0.006	-0.006	-0.002	0.039

$$\Sigma = 1.55$$

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

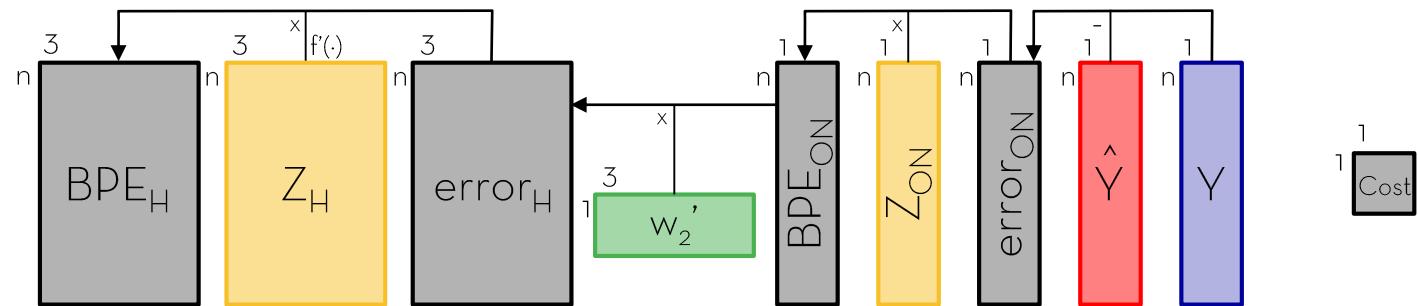
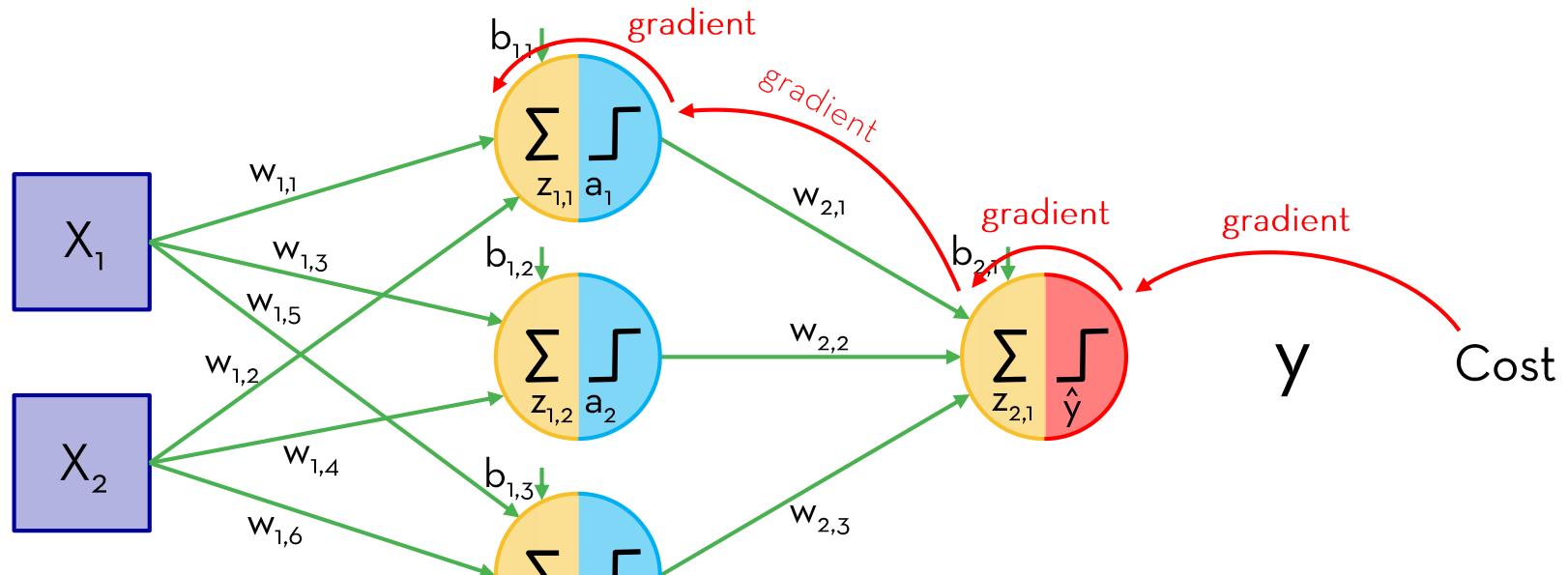
We repeat the same procedure for all observations.

We will need these values to update the weights.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6

A summary of the general principle of backpropagation

106

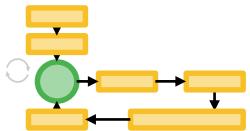


Steps of estimating neural networks

107

1. Design the neural network.
2. Initialize the neural network.
3. Forward propagation.
4. Calculate the cost.
5. Backpropagation.
6. Update weights.

Repeat 3. - 6. until convergence criteria is met.



<http://pages.cs.wisc.edu/~bolo/shipyard/neural/local.html>

Step 6: Update weights

How to optimize the weights?

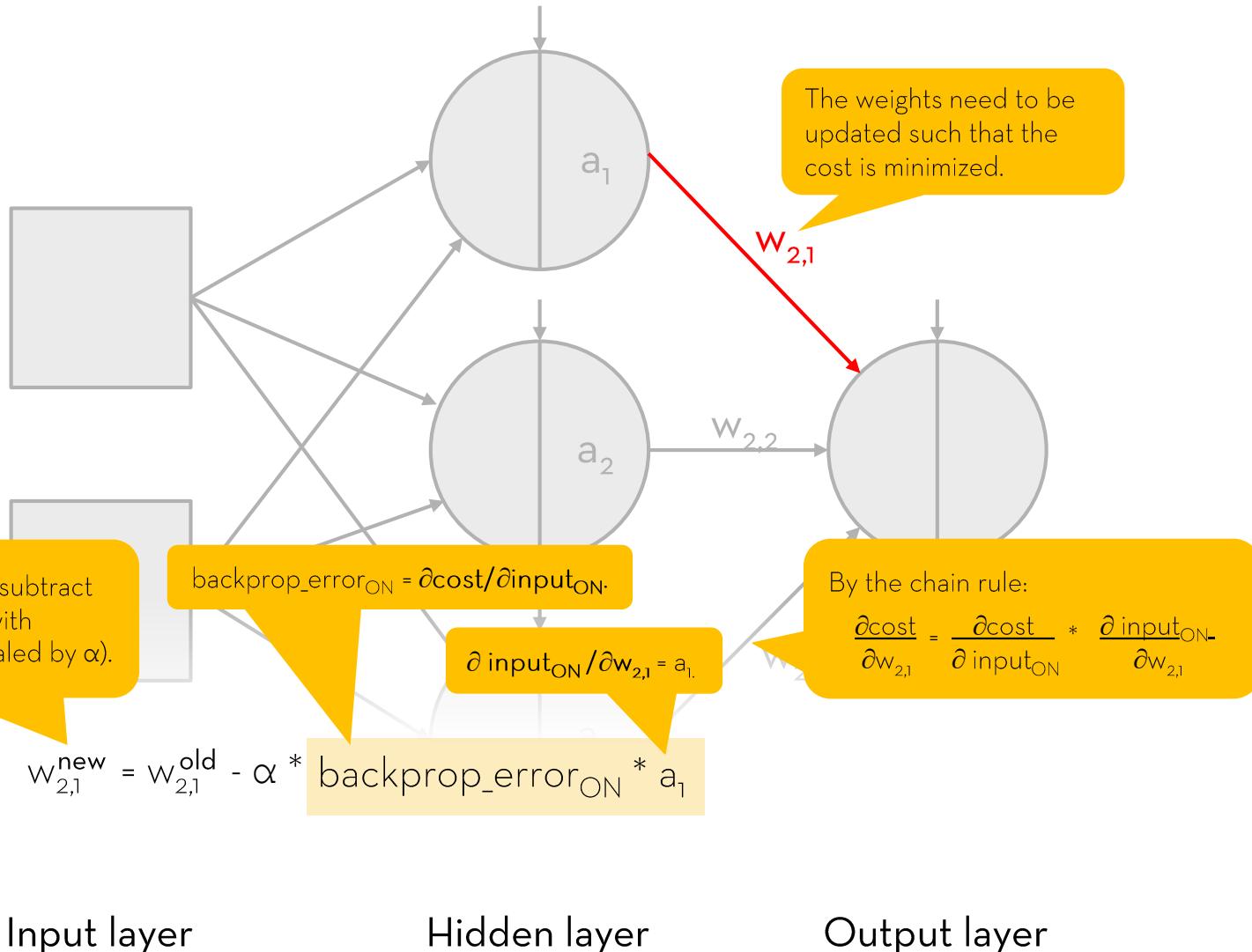
108

- a) Update the weights and biases leading to the output node.
- b) Update the weights and biases leading to the hidden layer.

Step 6: Update weights

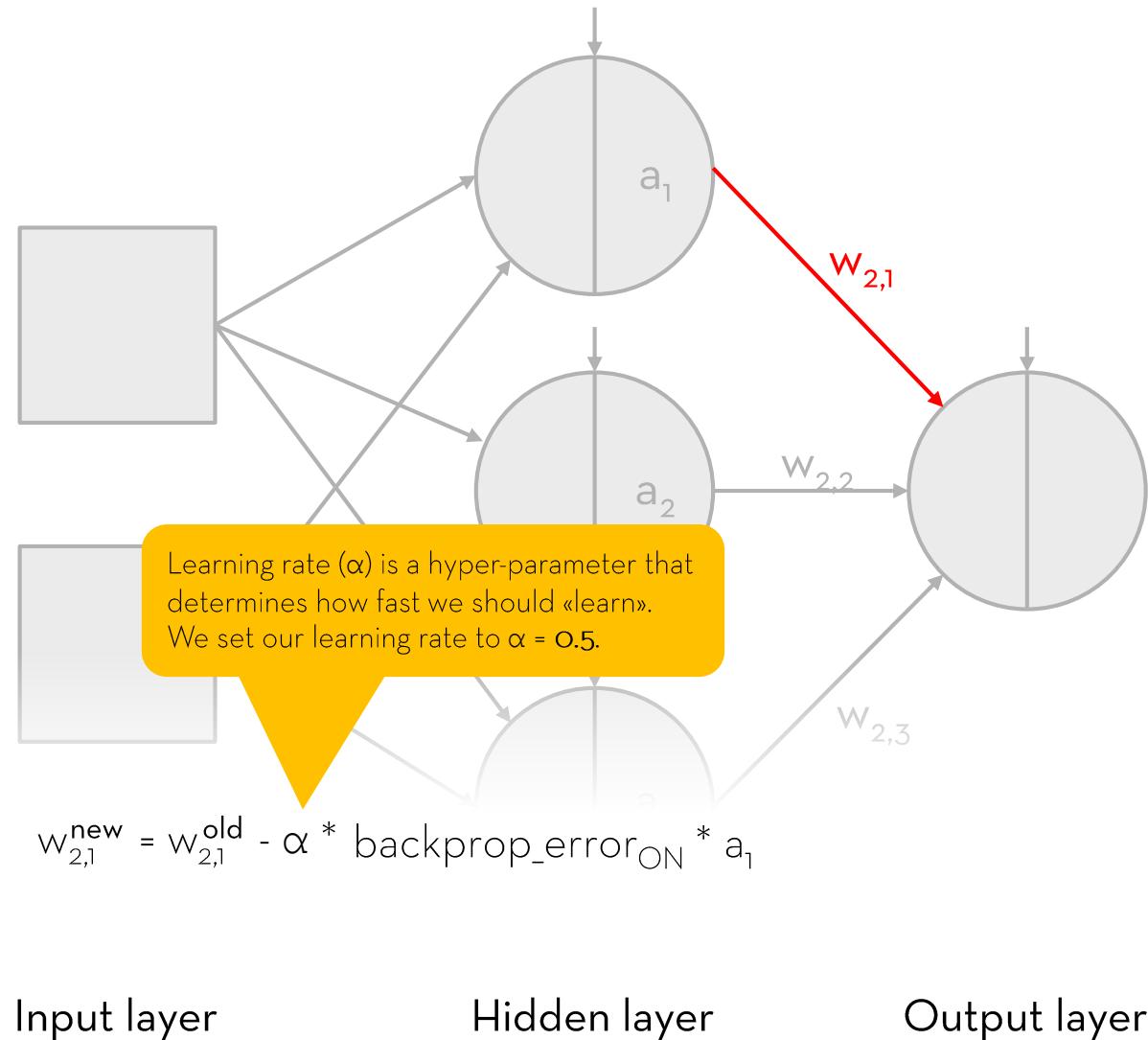
How to update the weights and biases? (1/4)

109



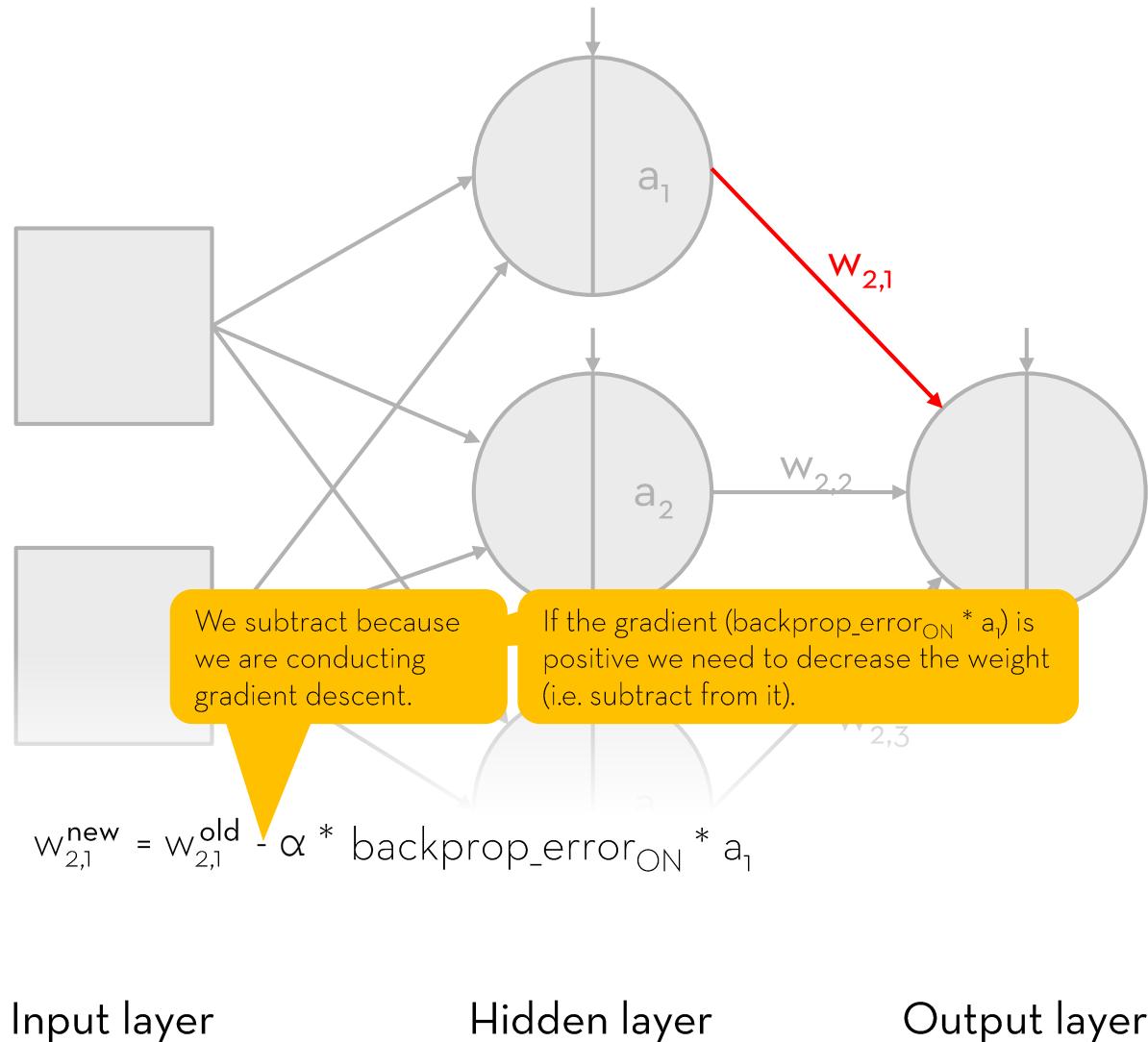
Step 6: Update weights

How to update the weights and biases? (2/4)



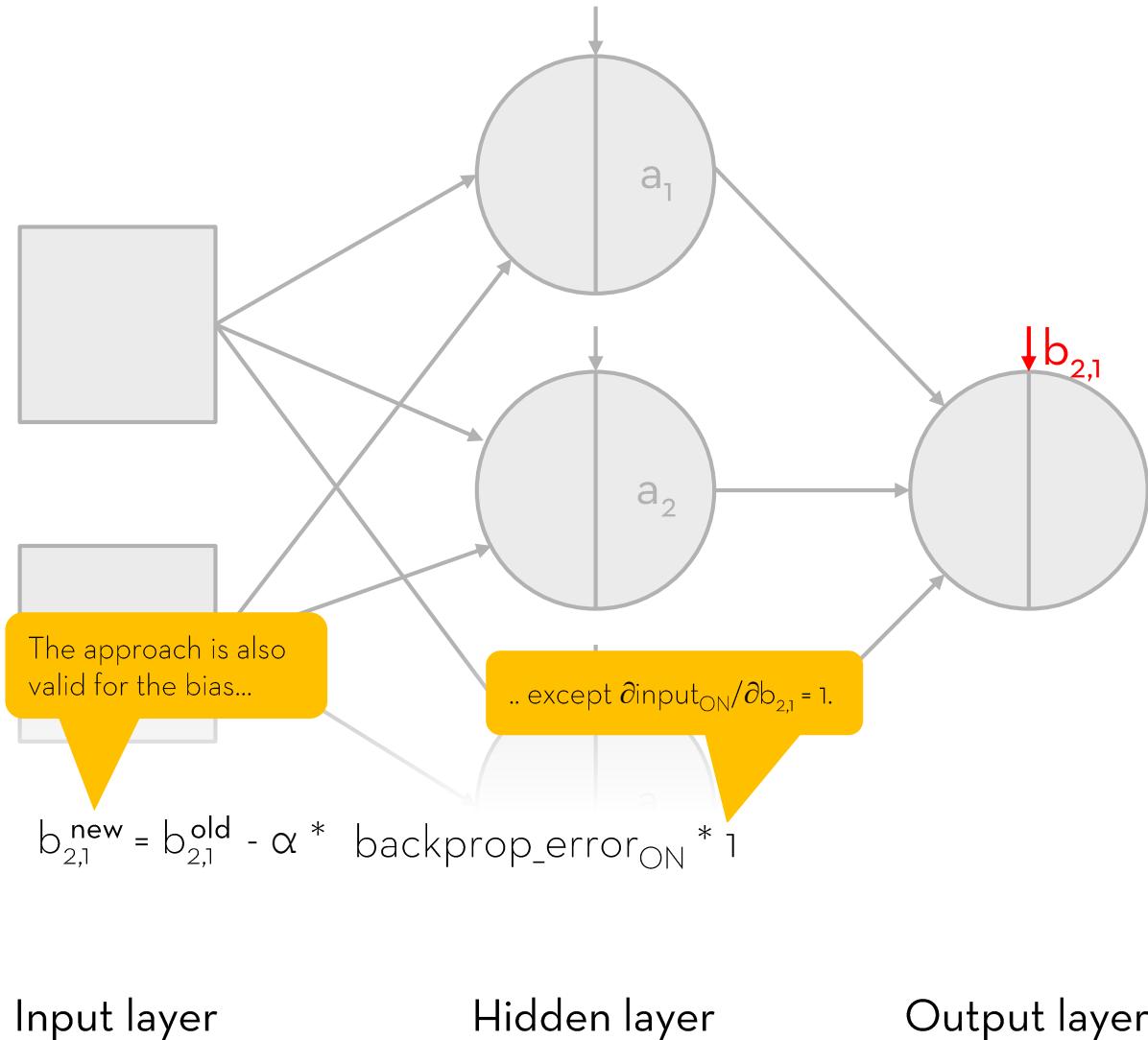
Step 6: Update weights

How to update the weights and biases? (3/4)



Step 6: Update weights

How to update the weights and biases? (4/4)



Step 6: Update weights

a) Update the weights and biases leading to the ON

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039

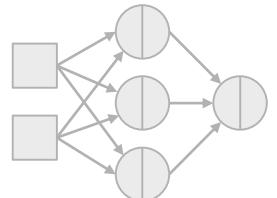
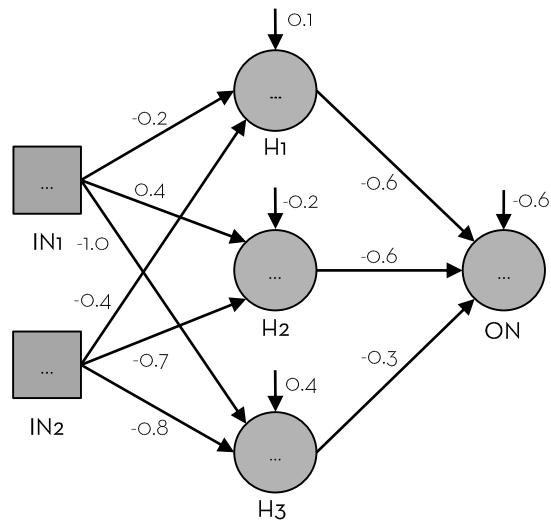
$\Sigma = 1.55$

	W1	W2	W3
IN1 \rightarrow HL	-0.2	0.4	-1.0
IN2 \rightarrow HL	-0.4	-0.7	-0.8
HL \rightarrow ON	-0.6	-0.6	-0.3

We import the output BPEs for every observation from the previous section to facilitate understanding.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

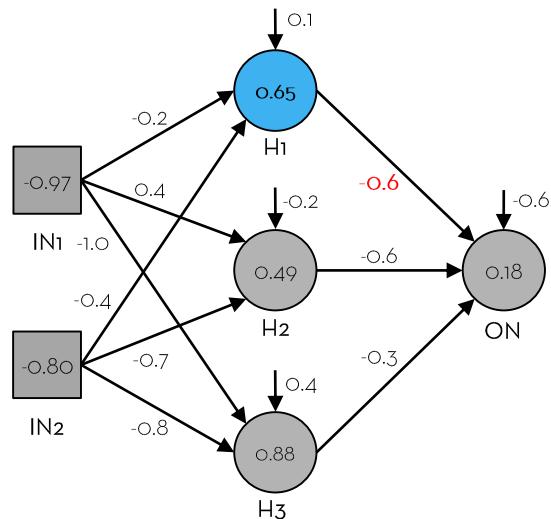
For familiarity, we will keep the backpropagation here for the first observation.



Step 6: Update weights

a) Update the weights leading to the ON (1/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039
$\Sigma = 1.55$								

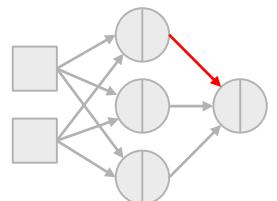


	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.6	-0.6	-0.3

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

To get the **new weight** we take the old weight, subtract the backpropagation error multiplied by the **learning rate 0.5** and the value of the **H1**.

$$\begin{aligned}
 \text{weight}_{\text{H1} \rightarrow \text{ON}} &= \text{weight}_{\text{H1} \rightarrow \text{ON}} - \alpha * \text{backprop_error}_{\text{ON}} * a_{\text{H1}} \\
 &= -0.6 - 0.5 * 0.025 * 0.65 \\
 &= -0.608
 \end{aligned}$$

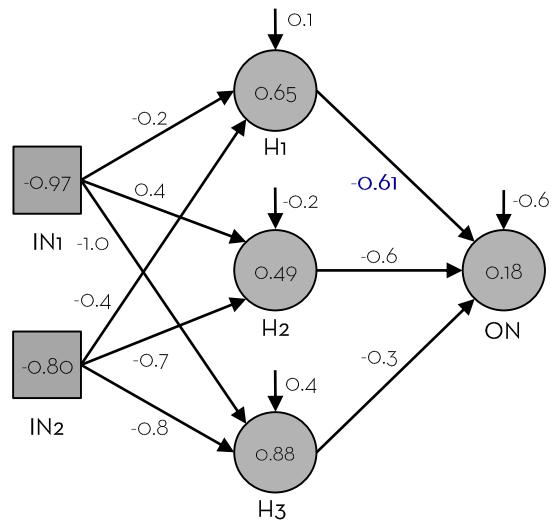


Step 6: Update weights

a) Update the weights leading to the ON (2/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039

$\Sigma = 1.55$

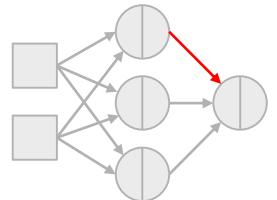


	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.608	-0.6	-0.3

The weight changes from -0.6 to -0.608.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

Notice: the weight decreased from the original value since the gradient is positive.

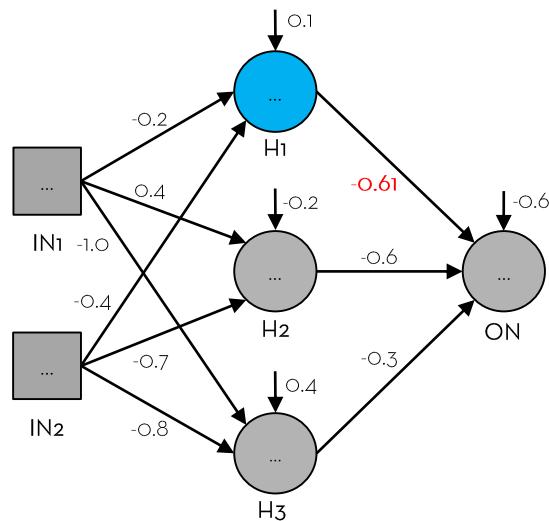


Step 6: Update weights

a) Update the weights leading to the ON (3/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039

$\Sigma = 1.55$



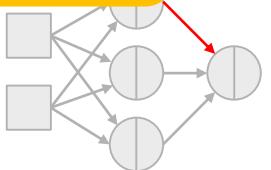
	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.608	-0.6	-0.3

We repeat this for every observation:

$$\begin{aligned}
 \text{weight}_{\text{H}1 \rightarrow \text{ON}} &= \text{weight}_{\text{H}1 \rightarrow \text{ON}} - \alpha * (\text{backprop_error}_{\text{ON}_2} * a_{\text{H}1,2} + \dots + \text{backprop_error}_{\text{ON}_{10}} * a_{\text{H}1,10}) \\
 &= -0.608 - 0.5 * (0.024 * 0.65 + \dots + 0.039 * 0.44) \\
 &= -0.496
 \end{aligned}$$

	H1	H2	H3
Bias	0.1	-0.2	0.4
Error	-0.015	-0.015	-0.008
Backprop. Err	-0.003	-0.004	-0.001

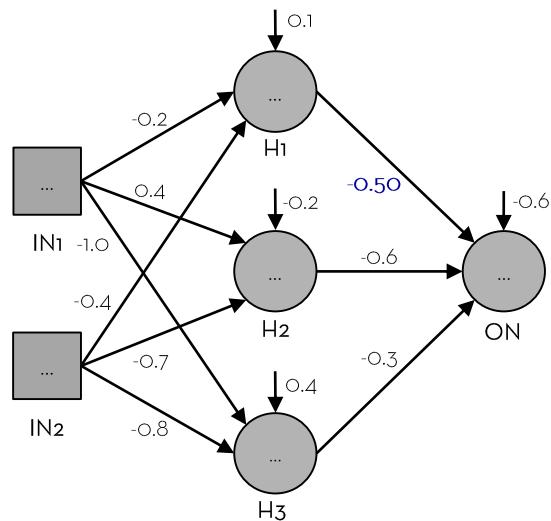
Since updating the weight is linear, you can update the weight one observation at a time, or all together.



Step 6: Update weights

a) Update the weights leading to the ON (4/9)

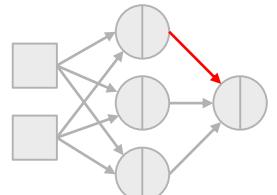
Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039
$\Sigma = 1.55$								



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.5	

The final weight after updating the weight with every observation is -0.496.

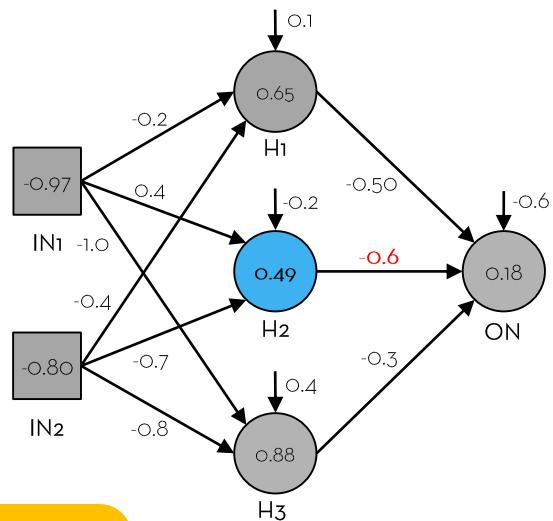
	H1	H2	H3	
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



Step 6: Update weights

a) Update the weights leading to the ON (5/9)

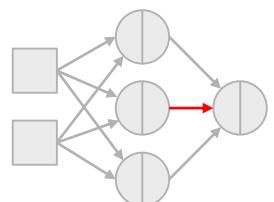
Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039
$\Sigma = 1.55$								



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	0.496	-0.6	-0.3

weight_{H2→ON}
 $= \text{weight}_{H2 \rightarrow ON} - \alpha * \text{backprop_error}_{ON} * a_{H2,1}$
 $= -0.6 - 0.5 * 0.025 * 0.49$
 $= -0.606$

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



Step 6: Update weights

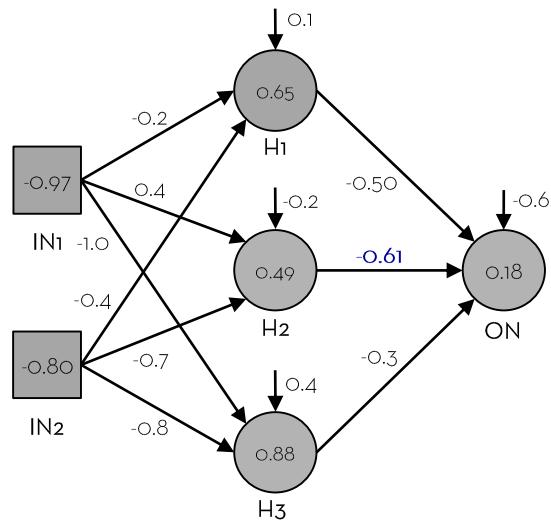
a) Update the weights leading to the ON (6/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039

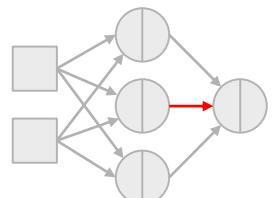
$\Sigma = 1.55$

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.606	-0.7

The weight changes from -0.6 to -0.606.



	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

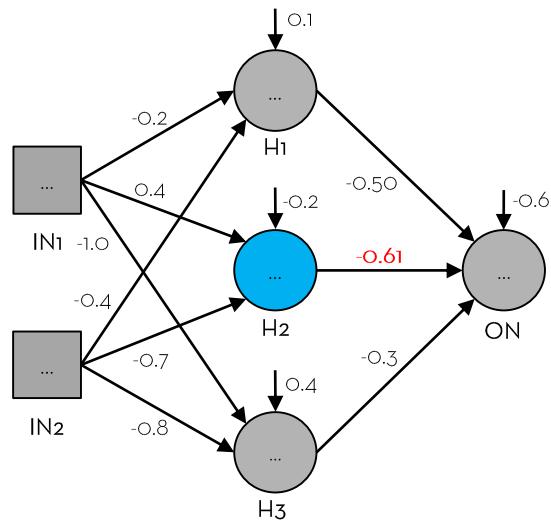


Step 6: Update weights

a) Update the weights leading to the ON (7/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039

$\Sigma = 1.55$



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.606	-0.3

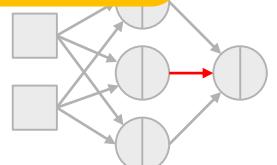
We repeat this for every observation:

$$\begin{aligned}
 \text{weight}_{H2 \rightarrow ON} &= \text{weight}_{H2 \rightarrow ON} - \alpha * (\text{backprop_error}_{ON2} * a_{H1,2} + \dots + \text{backprop_error}_{ON10} * a_{H1,10}) \\
 &= -0.606 - 0.5 * (0.024 * 0.55 + \dots + 0.039 * 0.49) \\
 &= -0.506
 \end{aligned}$$

	H1	H2	H3
Bias	0.1	-0.2	0.4
Error	-0.015	-0.015	-0.008
Backprop. Err	-0.003	-0.004	-0.001

0.18

0.025



Step 6: Update weights

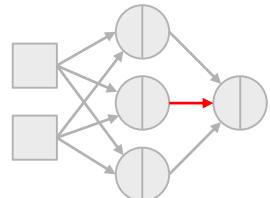
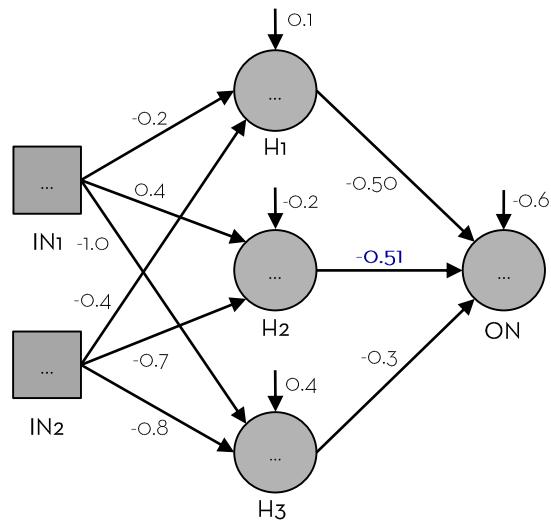
a) Update the weights leading to the ON (8/9)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039
$\Sigma = 1.55$								

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.3

The final weight after updating the weight with every observation is -0.506.

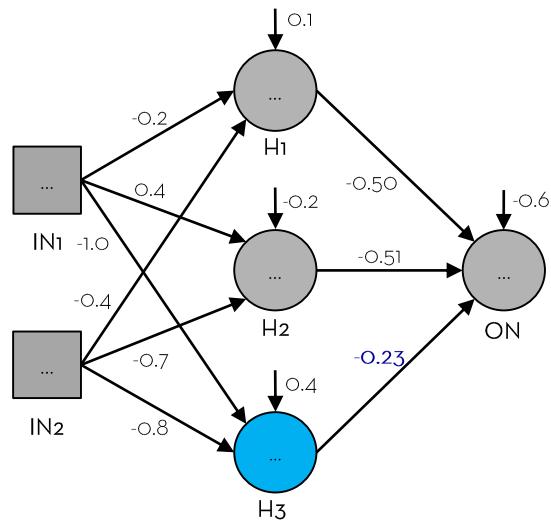
	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



Step 6: Update weights

a) Update the weights leading to the ON (9/9)

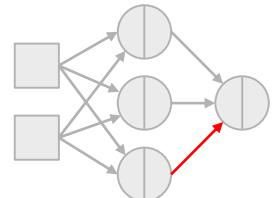
Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039
$\Sigma = 1.55$								



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

Repeat for the final weight.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

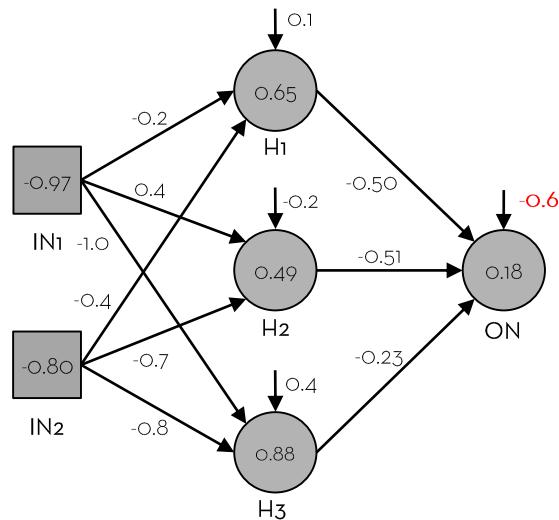


Step 6: Update weights

a) Update the bias of the ON (1/3)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039

$\Sigma = 1.55$

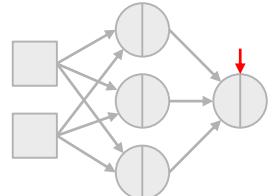


	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.6
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

We **update the bias** as follows:

$$\begin{aligned} \text{bias}_{\text{ON}} &= \text{bias}_{\text{ON}} - \alpha * \text{backprop_error}_{\text{ON}_i} \\ &= \textcolor{red}{-0.6} - 0.5 * \textcolor{teal}{0.025} \\ &= -0.613 \end{aligned}$$

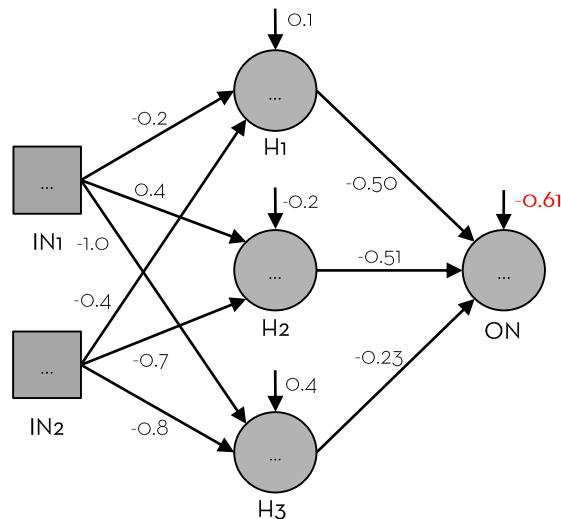


Step 6: Update weights

a) Update the bias of the ON (2/3)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	0.039

$\Sigma = 1.55$

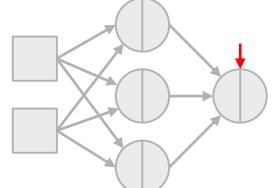


	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

And repeat for every observation.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.613
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

$$\begin{aligned}
 \text{bias}_{\text{ON}} &= \text{bias}_{\text{ON}} - \alpha * (\text{backprop_error}_{\text{ON}2} + \dots + \text{backprop_error}_{\text{ON}10}) \\
 &= -0.613 - 0.5 * (0.024 + \dots + 0.039) \\
 &= -0.327
 \end{aligned}$$

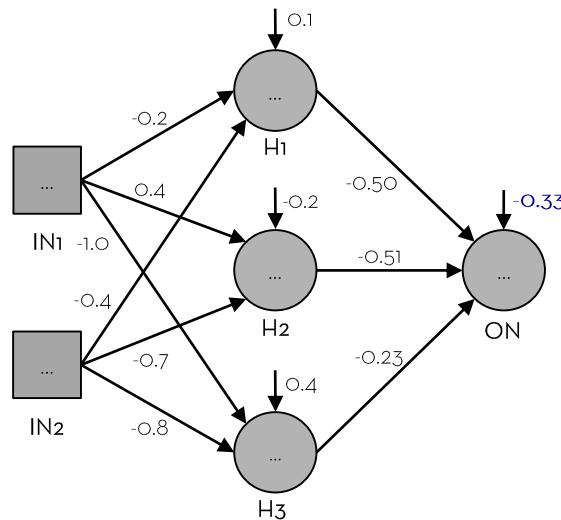


Step 6: Update weights

a) Update the bias of the ON (3/3)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE ON
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	0.025
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	0.024
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	0.023
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	-0.122
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	0.039

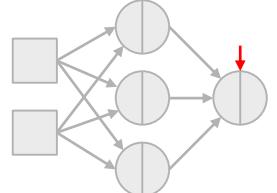
$\Sigma = 1.55$



	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

The bias of ON
changes to -0.327.



Step 6: Update weights

b) Update the weights and biases leading to the HL

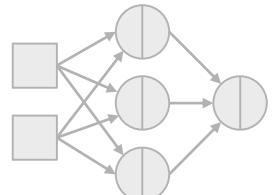
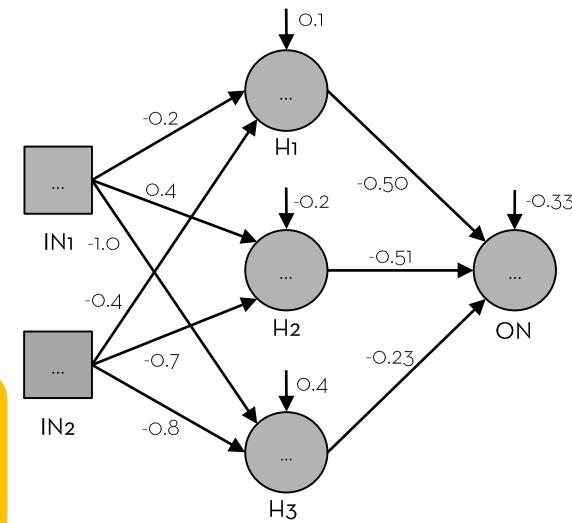
Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	-0.006	-0.006	-0.002

$$\Sigma = 1.55$$

	W1	W2	W3
IN1 → HL	-0.2	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

Now, we import the hidden layer BPEs for every observation from the previous section to facilitate understanding.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



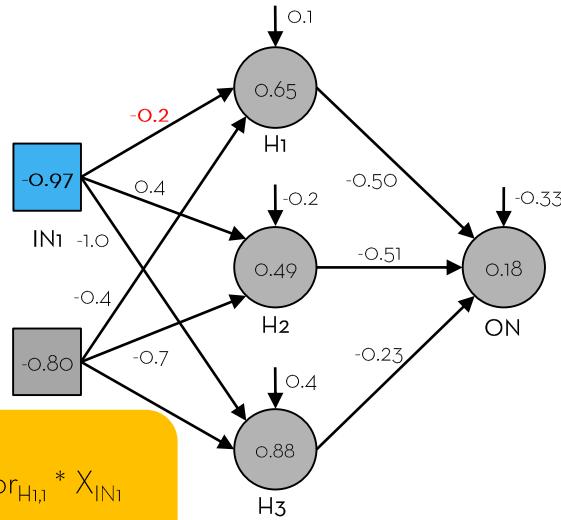
Step 6: Update weights

b) Update the weights leading to the HL nodes (1/5)

127

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
O	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	-0.006	-0.006	-0.002

$$\Sigma = 1.55$$

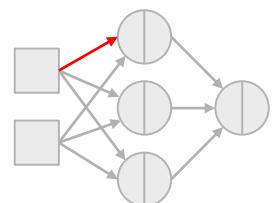


	W1	W2	W3
IN1 → HL	-0.2	-0.1	
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

Update the weight from IN1 to H1:

$$\begin{aligned} \text{weight}_{\text{IN1} \rightarrow \text{H1}} &= \text{weight}_{\text{IN1} \rightarrow \text{H1}} - \alpha * \text{backprop_error}_{\text{H1},1} * X_{\text{IN1}} \\ &= -0.2 - 0.5 * -0.003 * -0.97 \\ &= -0.202 \end{aligned}$$

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



Step 6: Update weights

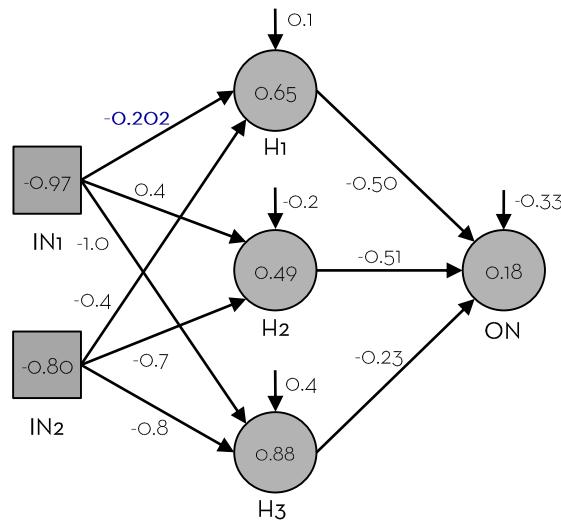
b) Update the weights leading to the HL nodes (2/5)

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	-0.006	-0.006	-0.002

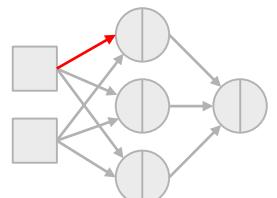
$\Sigma = 1.55$

	W1	W2	W3
IN1 → HL	-0.202	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

The weight changes from -0.2 to -0.202.



	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



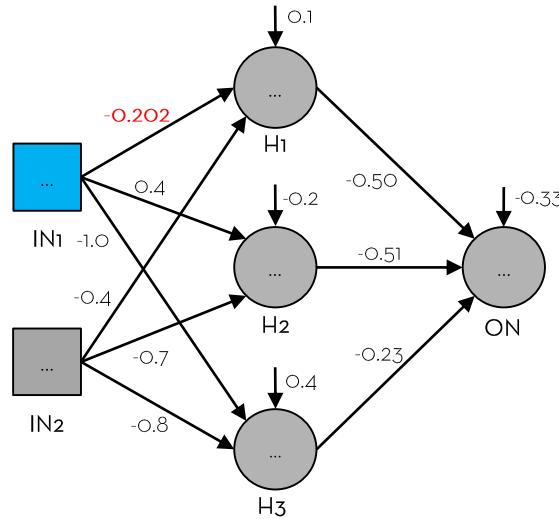
Step 6: Update weights

b) Update the weights leading to the HL nodes (3/5)

129

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05	-0.006	-0.006	-0.002

$$\Sigma = 1.55$$

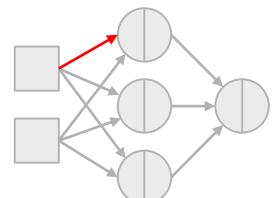


	W1	W2	W3
IN1 → HL	-0.202	0.4	-1.0
IN2 → HL	-0.4	-0.7	-0.8
HL → ON	-0.496	-0.506	

	H1
Bias	0.1
Error	-0.015
Backprop. Err	-0.003 -0.004 -0.001 0.025

We repeat this for every observation:

$$\begin{aligned} \text{weight}_{\text{IN} \rightarrow \text{H}1} &= \text{weight}_{\text{IN} \rightarrow \text{H}1} - \alpha * (\text{backprop_error}_{\text{HL}1,2} * X_{\text{IN}2} + \dots + \text{backprop_error}_{\text{HL}1,10} * X_{\text{IN}10}) \\ &= -0.202 - 0.5 * (-0.003 * -0.66 + \dots + -0.006 * 1.04) \\ &= -0.230 \end{aligned}$$

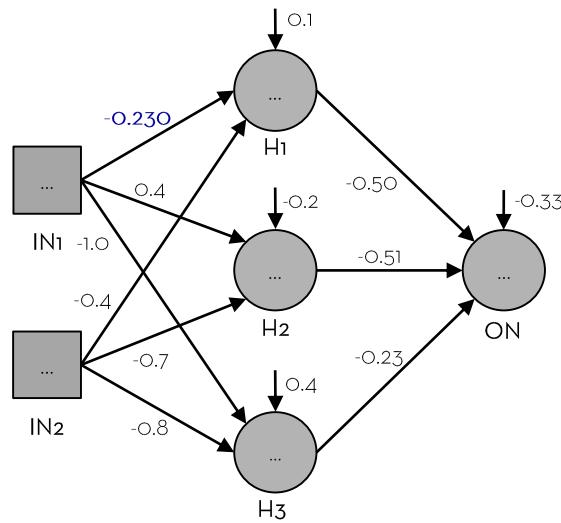


Step 6: Update weights

b) Update the weights leading to the HL nodes (4/5)

130

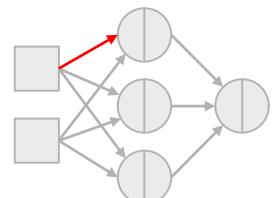
Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	-0.006	-0.006	-0.002
$\Sigma = 1.55$										



	W1	W2	W3
IN1 → HL	-0.230	0.4	
IN2 → HL	-0.4	-0.7	-0.1
HL → ON	-0.496	-0.506	-0.2

The final weight after updating the weight with every observation is -0.230.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



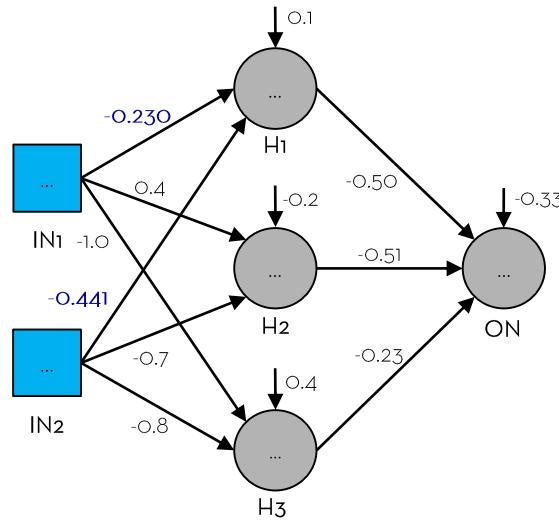
Step 6: Update weights

b) Update the weights leading to the HL nodes (5/5)

131

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	0.05	-0.006	-0.006	-0.002

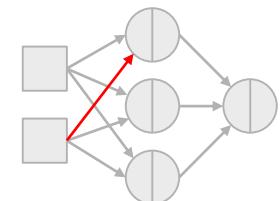
$$\Sigma = 1.55$$



	W1	W2	W3
IN1 → HL	-0.230	0.4	-1.0
IN2 → HL	-0.441	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

We update both weights of this hidden layer.

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



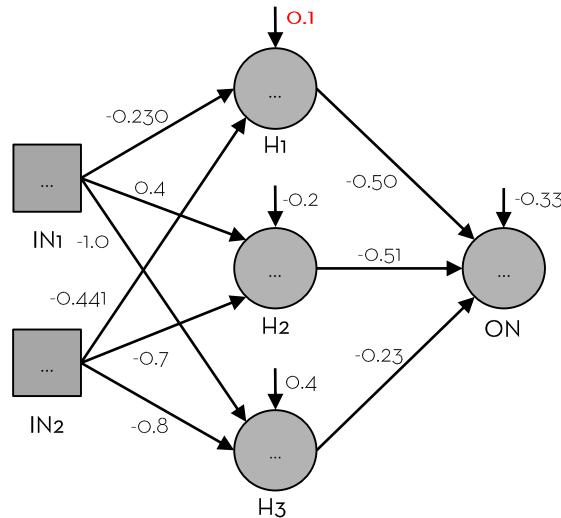
Step 6: Update weights

b) Update the bias of the HL nodes (1/2)

132

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05	-0.006	-0.006	-0.002

$$\Sigma = 1.55$$

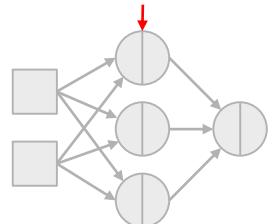


	W1	W2	W3
IN1 → HL	-0.230	0.4	-1.0
IN2 → HL	-0.441	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

	H1	H2	H3	ON
Bias	0.1	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

We **update the bias** as follows:

$$\begin{aligned}
 \text{bias}_{H1} &= \text{bias}_{H1} - \alpha * (\text{backprop_error}_{H1,1,1} + \dots + \text{backprop_error}_{H1,1,10}) \\
 &= 0.1 - 0.5 * (-0.003 + \dots + -0.006) \\
 &= 0.061
 \end{aligned}$$



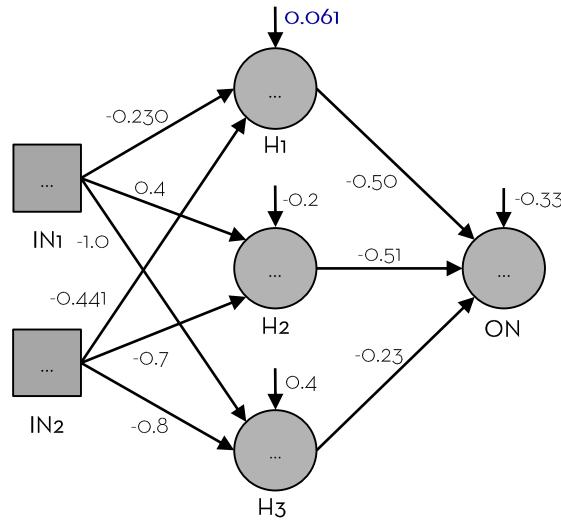
Step 6: Update weights

b) Update the bias of the HL nodes (2/2)

133

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	-0.006	-0.006	-0.002

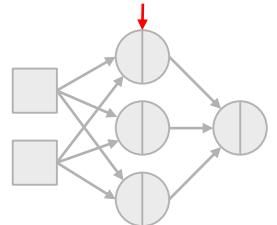
$\Sigma = 1.55$



	W1	W2	W3
IN1 → HL	-0.230	0.4	-1.0
IN2 → HL	-0.441	-0.7	-0.8
HL → ON	-0.496	-0.506	-0.233

	H1	H2	H3	ON
Bias	0.061	-0.2	0.4	-0.327
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025

The bias changes from 0.1 to 0.061.



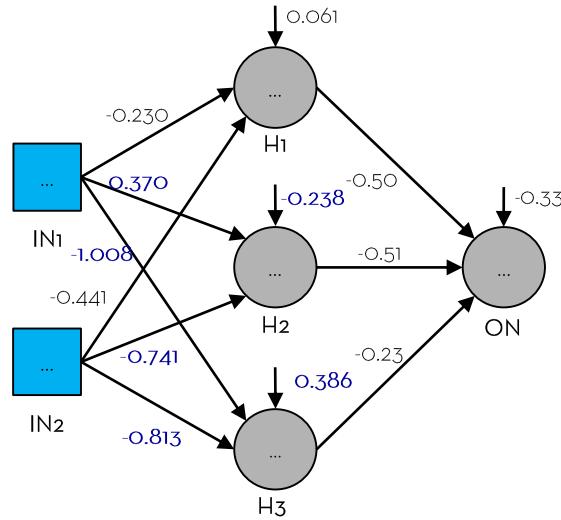
Step 6: Update weights

b) Update the rest of the weights and biases

134

Y	X1	X2	H1	H2	H3	Output	Cost	BPE H1	BPE H2	BPE H3
0	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03	-0.003	-0.004	-0.001
0	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03	-0.003	-0.004	-0.001
0	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03	-0.003	-0.004	-0.001
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67	0.017	0.018	0.005
0	1.04	0.34	0.44	0.49	0.29	0.22	<u>0.05</u>	-0.006	-0.006	-0.002

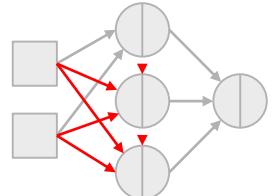
$$\Sigma = 1.55$$



	W1	W2	W3
IN1 → HL	-0.230	0.370	-1.008
IN2 → HL	-0.441	-0.741	-0.813
HL → ON	-0.496	-0.506	-0.233

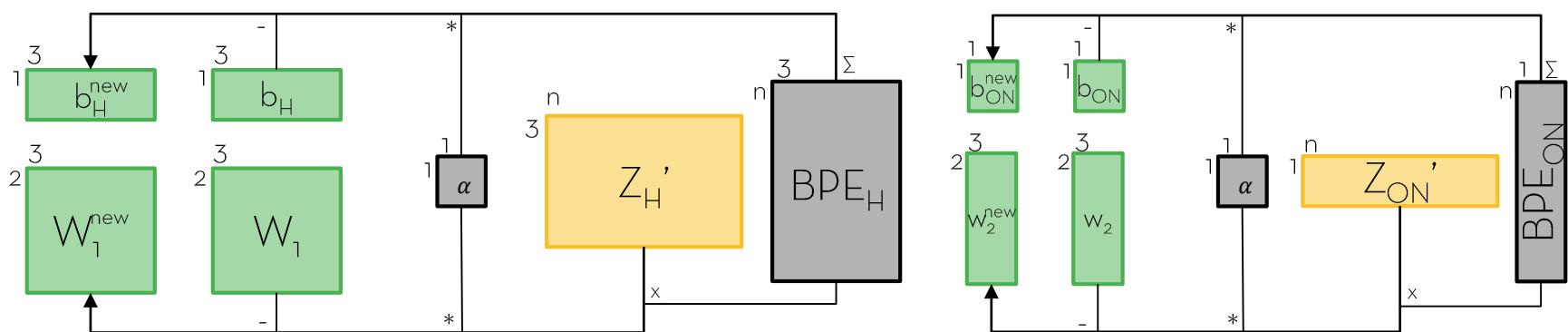
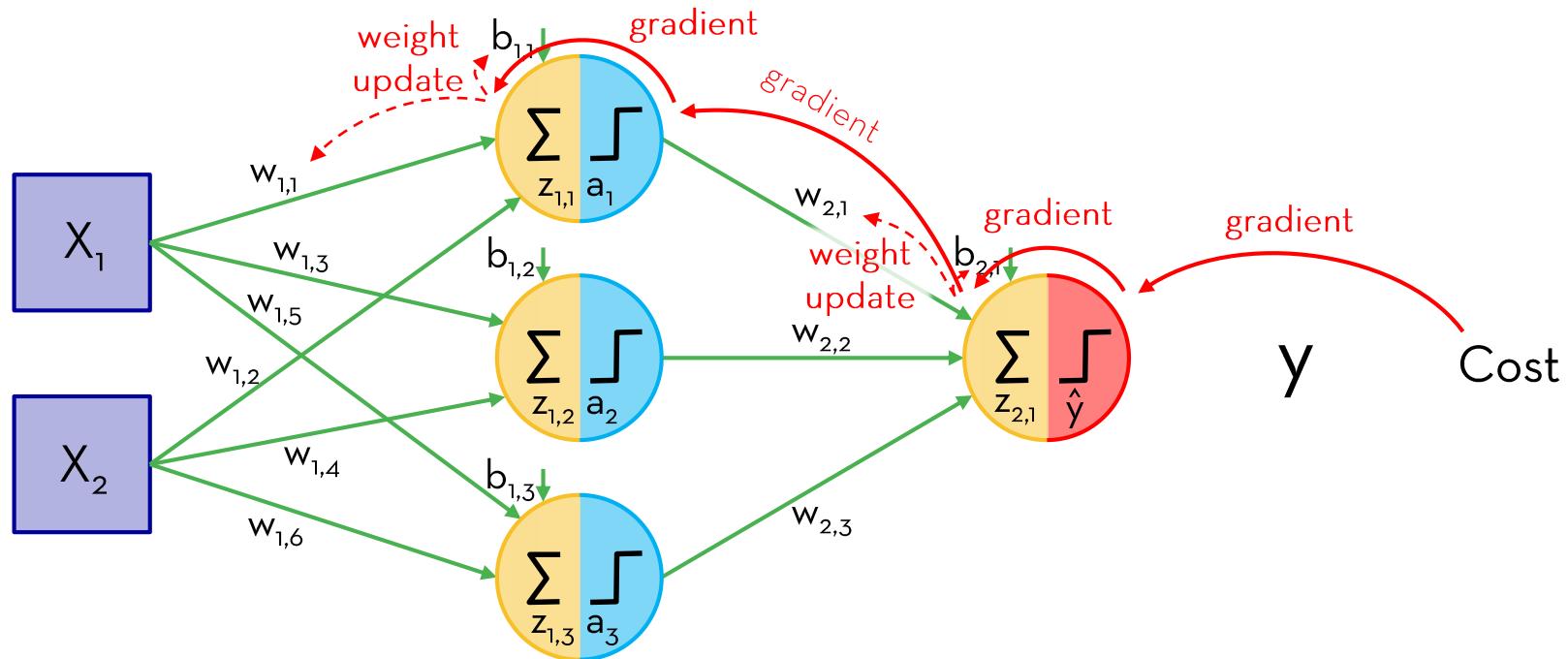
We proceed in the **same way** for the rest of the HL and update all weights and biases.

	H1	H2	H3	ON
Bias	0.061	-0.238	0.386	0.386
Error	-0.015	-0.015	-0.008	0.18
Backprop. Err	-0.003	-0.004	-0.001	0.025



A summary of the general principle of updating the weights using gradient descent

135



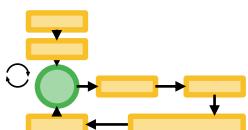
Repeat steps 3. - 6. until convergence criteria is met

136

- We now **repeat** the forward propagation and backpropagation process.
- One **iteration** (i.e., estimating once forward propagation and backpropagation for the entire dataset) **is called one cycle or epoch**.
- **Training of neural networks requires many cycles** (10'000s of cycles) until the neural network is good enough (the summed-up error doesn't change anymore).
- We set the following training goal: **summed-up error < 0.001**.

There are **many possible training goals**, for example:

- 1) Maximum number of epochs reached.
- 2) Maximum amount of time exceeded/performance minimized to a threshold.
- 3) Performance gradient falls below a minimum.
- 4) Validation performance has increased x times since the last time it decreased (when using validation).



<https://ch.mathworks.com/help/nnet/ref/traingd.html>

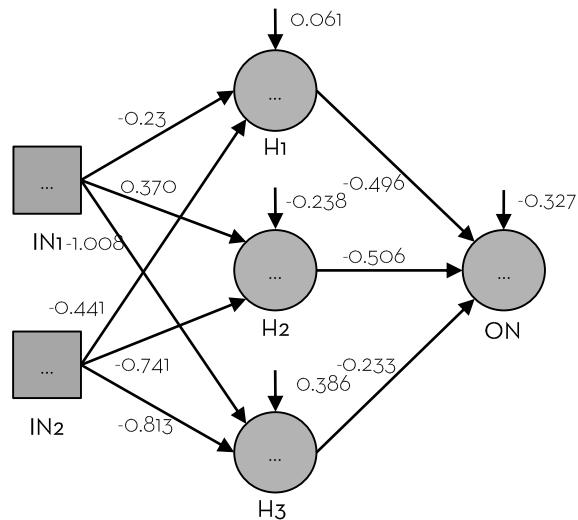
After 1 epoch

During iterative training of a neural network , an **epoch is a single pass through the entire training set**, followed by testing of the verification set.
When training in mini-batches, then an epoch consist of multiple iterations.

Y	X1	X2	H1	H2	H3	Output	Cost
O	-0.97	-0.80	0.65	0.49	0.88	0.18	0.03
O	-0.69	-0.97	0.65	0.55	0.87	0.17	0.03
O	-1.14	-1.09	0.68	0.53	0.92	0.17	0.03
...
1	-0.79	-0.52	0.61	0.46	0.83	0.18	0.67
O	1.04	0.34	0.44	0.49	0.29	0.22	0.05
$\Sigma = 1.55$							

	W1	W2	W3
IN1 → HL	-0.230	0.370	-1.008
IN2 → HL	-0.441	-0.741	-0.813
HL → ON	-0.496	-0.506	-0.233

	H1	H2	H3	ON
Bias	0.061	-0.238	0.386	-0.327



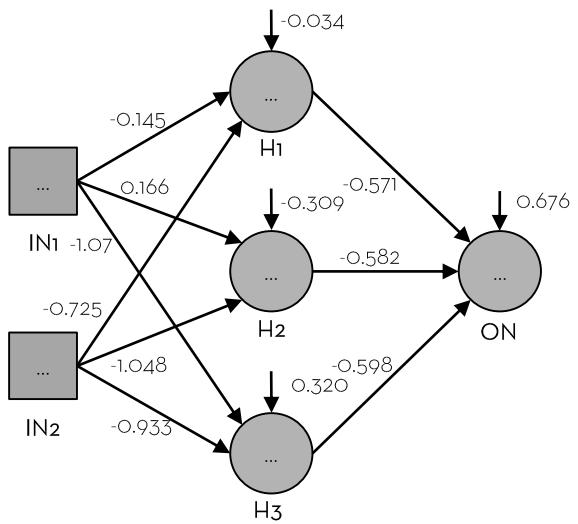
After 10 epochs

138

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.71	0.58	0.89	0.36	0.13
0	-0.69	-0.97	0.72	0.63	0.88	0.36	0.13
0	-1.14	-1.09	0.77	0.64	0.93	0.34	0.12
...
1	-0.79	-0.52	0.66	0.52	0.84	0.38	0.38
0	1.04	0.34	0.34	0.39	0.25	0.52	0.27
$\Sigma = 0.96$							

	W1	W2	W3
IN1 → HL	-0.415	0.166	-1.071
IN2 → HL	-0.725	-1.048	-0.933
HL → ON	-0.571	-0.582	-0.598

	H1	H2	H3	ON
Bias	-0.034	-0.309	0.320	0.676



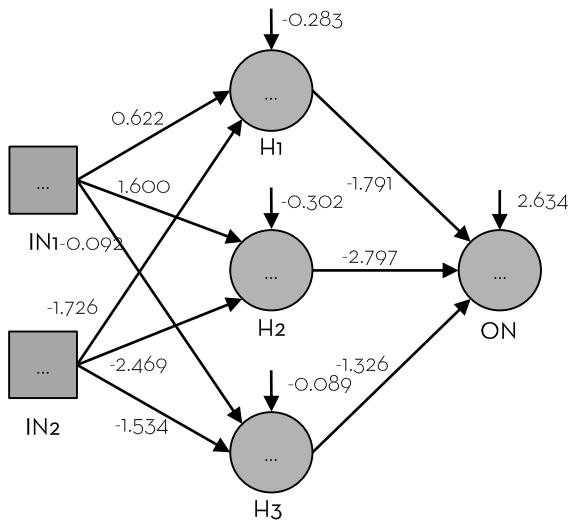
After 100 epochs

139

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.62	0.53	0.78	0.27	0.07
0	-0.69	-0.97	0.72	0.73	0.78	0.15	0.02
0	-1.14	-1.09	0.71	0.64	0.85	0.18	0.03
...
1	-0.79	-0.52	0.53	0.43	0.69	0.39	0.37
0	1.04	0.34	0.44	0.63	0.33	0.42	0.17
$\Sigma = 0.41$							

	W1	W2	W3
IN1 → HL	0.622	1.600	-0.092
IN2 → HL	-1.726	-2.469	-1.534
HL → ON	-1.791	-2.797	-1.326

	H1	H2	H3	ON
Bias	-0.283	-0.302	-0.089	2.634



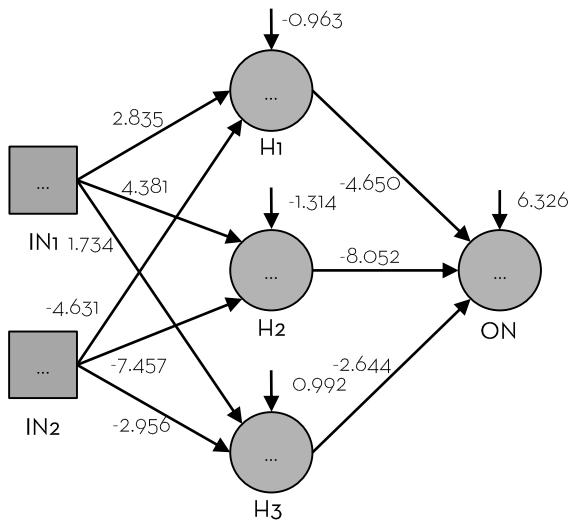
After 1'000 epochs

140

Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.50	0.60	0.42	0.127	0.016
0	-0.69	-0.97	0.83	0.95	0.66	0.001	0.000
0	-1.14	-1.09	0.70	0.86	0.56	0.005	0.000
...
1	-0.79	-0.52	0.31	0.29	0.30	0.851	0.022
0	1.04	0.34	0.60	0.67	0.45	0.045	0.002
$\Sigma = 0.02$							

	W1	W2	W3
IN1 → HL	2.835	4.381	1.734
IN2 → HL	-4.631	-7.457	-2.956
HL → ON	-4.650	-8.052	-2.644

	H1	H2	H3	ON
Bias	-0.963	-1.314	0.992	6.326



After 8'000 epochs

141

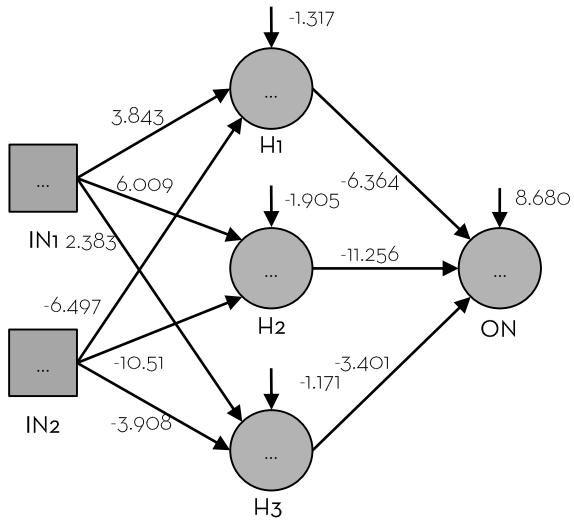
Y	X1	X2	H1	H2	H3	Output	Cost
0	-0.97	-0.80	0.54	0.66	0.41	0.027	0.001
0	-0.69	-0.97	0.91	0.98	0.73	0.000	0.000
0	-1.14	-1.09	0.80	0.94	0.59	0.000	0.000
...
1	-0.79	-0.52	0.27	0.23	0.26	0.968	0.001
0	1.04	0.34	0.62	0.68	0.49	0.010	<u>0.000</u>

$\Sigma < 0.001$

Goal reached!

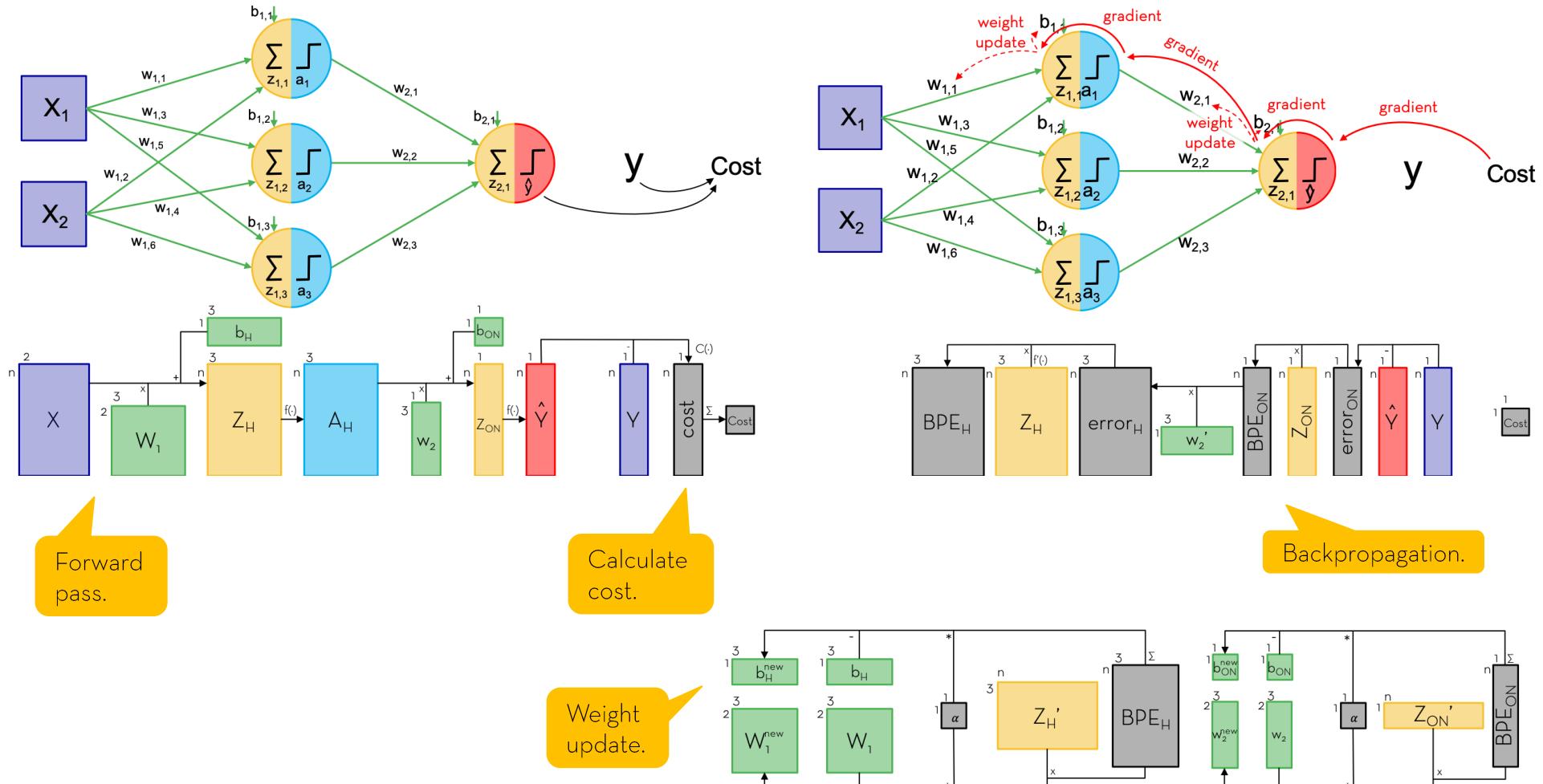
	W1	W2	W3
IN1 → HL	3.843	6.009	2.383
IN2 → HL	-6.497	-10.508	-3.908
HL → ON	-6.364	-11.256	-3.401

	H1	H2	H3	ON
Bias	-1.317	-1.905	-1.171	8.680



Summary of the neural network training algorithm

142



Summary of the neural network training algorithm

143

3. Forward propagation:

- Calculate the values of the hidden layer: $z_H = f(\sum_{i \in IN} X_i w_{i \rightarrow H} + b_H)$ for all $H \in HL$.
- Calculate the values of the output layer: $\hat{y} = f(\underbrace{\sum_{i \in HL} a_i w_{i \rightarrow ON} + b_{ON}}_{z_{ON}})$.

IN: input nodes.
ON: output node.
HL: hidden layer nodes.
 $f(\cdot)$: sigmoid function.

4. Backpropagation:

- For each layer (in reverse):

- Calculate error:

- Calculate the backpropagation error (BPE):

Output node (ON)	Hidden node $H \in HL$.
$(\hat{y} - y)$	$BPE_{ON} * w_{H \rightarrow ON}$
$error_{ON} * f'(z_{ON})$	$error_H * f'(z_H)$

5. Weight updates:

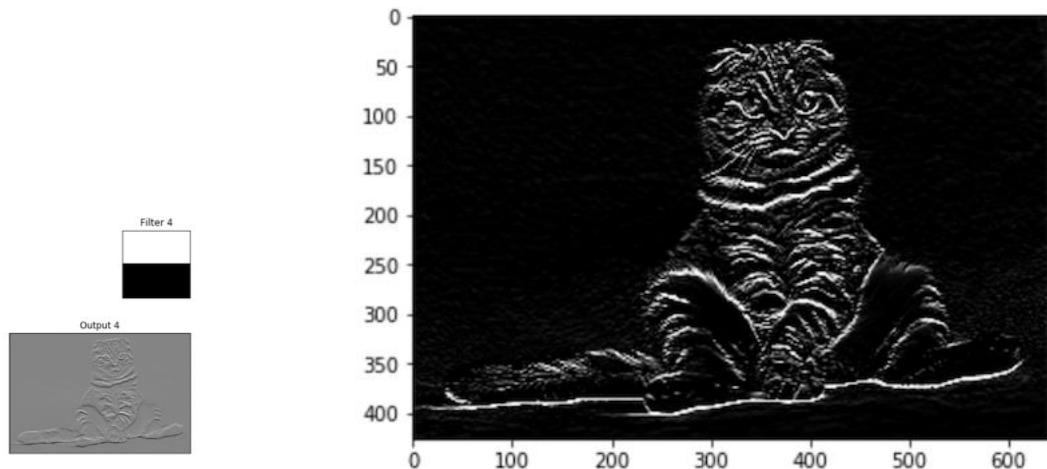
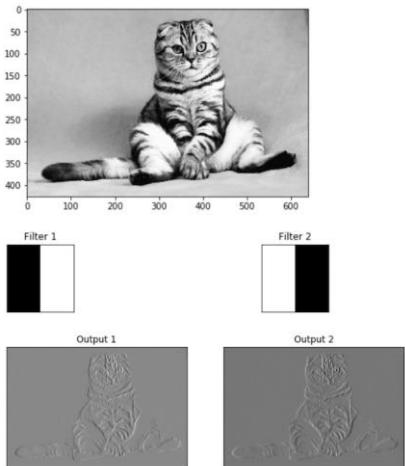
- Update the weights: $w_{from \rightarrow to}^{new} = w_{from \rightarrow to}^{old} - \alpha * bpe_{to} * a_{from}$.
- Update the biases: $b_{node}^{new} = b_{node}^{old} - \alpha * bpe_{node}$.

Advantages and disadvantages of neural networks

144

Advantages:

- Relatively simple learning algorithm.
- Can approximate any function, regardless of its linearity.
- Great for complex / abstract problems like image recognition.



<https://www.quora.com/What-are-the-advantages-disadvantages-of-Artificial-Neural-networks>

<http://neuralnetworksanddeeplearning.com/chap4.html>

<https://towardsdatascience.com/wtf-is-image-classification-8e78a8235acb>

Advantages and disadvantages of neural networks

145

Disadvantages:

- Difficult to interpret what a neural network has learned and why.
- Cannot deal with uncertainties.
- Computationally intensive to train.
- Mimics the brain – machine cannot completely copy a complex organism.

Breaking Linear Classifiers on ImageNet

Mar 30, 2015

You've probably heard that Convolutional Networks work very well in practice and across a wide range of visual recognition problems. You may have also read articles and papers that claim to reach a near "*human-level performance*". There are all kinds of caveats to that (e.g. see my G+ post on [Human Accuracy is not a point, it lives on a tradeoff curve](#)), but that is not the point of this post. I do think that these systems now work extremely well across many visual recognition tasks, especially ones that can be posed as simple classification.

Yet, a second group of seemingly baffling results has emerged that brings up an apparent contradiction. I'm referring to several people who have noticed that it is possible to take an image that a state-of-the-art Convolutional Network thinks is one class (e.g. "panda"), and it is possible to change it almost imperceptibly to the human eye in such a way that the Convolutional Network suddenly classifies the image as any other class of choice (e.g. "gibbon"). We say that we *break*, or *fool* ConvNets. See the image below for an illustration:

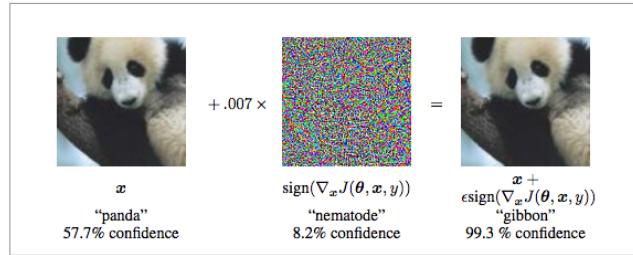


Figure from [Explaining and Harnessing Adversarial Examples](#) by Goodfellow et al.

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

<https://www.quora.com/What-are-the-advantages-disadvantages-of-Artificial-Neural-networks>

Making machines think like humans.

How (Shallow) Neural Networks work