

Indian Institute of Information Technology, Allahabad



Neural Networks Feed Forward Networks

By

Dr. Shiv Ram Dubey

Assistant Professor

Computer Vision And Biometrics Lab (CVBL)

Department Of Information Technology

Indian Institute Of Information Technology, Allahabad

Email: srdubey@iiita.ac.in

Web: <https://profile.iiita.ac.in/srdubey/>



ABOUT IIIT ALLAHABAD (A RESEARCH LED INSTITUTE OF NATIONAL IMPORTANCE)



DISCLAIMER

The content (text, image, and graphics) used in this slide are adopted from many sources for Academic purposes. Broadly, the sources have been given due credit appropriately. However, there is a chance of missing out some original primary sources. The authors of this material do not claim any copyright of such material.

IMAGE CATEGORIZATION / CLASSIFICATION



THE STATISTICAL LEARNING FRAMEWORK

- Apply a prediction function to a feature representation of the image to get the desired output:

$f(\text{apple image}) = \text{"apple"}$

$f(\text{tomato image}) = \text{"tomato"}$

$f(\text{cow image}) = \text{"cow"}$

THE STATISTICAL LEARNING FRAMEWORK

$$y = f(\mathbf{x})$$

output prediction
 function Image
 feature

The diagram illustrates the statistical learning framework equation $y = f(\mathbf{x})$. Three red arrows point from the labels below to the components of the equation: one from 'output' to 'y', one from 'prediction function' to 'f', and one from 'Image feature' to 'x'.

THE STATISTICAL LEARNING FRAMEWORK

$$y = f(\mathbf{x})$$

output prediction Image
 function feature

The diagram illustrates the statistical learning framework equation $y = f(\mathbf{x})$. Three red arrows point from the labels below to the components of the equation: one from 'output' to y , one from 'prediction function' to f , and one from 'Image feature' to \mathbf{x} .

- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, estimate the prediction function f by minimizing the prediction error on the training set

THE STATISTICAL LEARNING FRAMEWORK

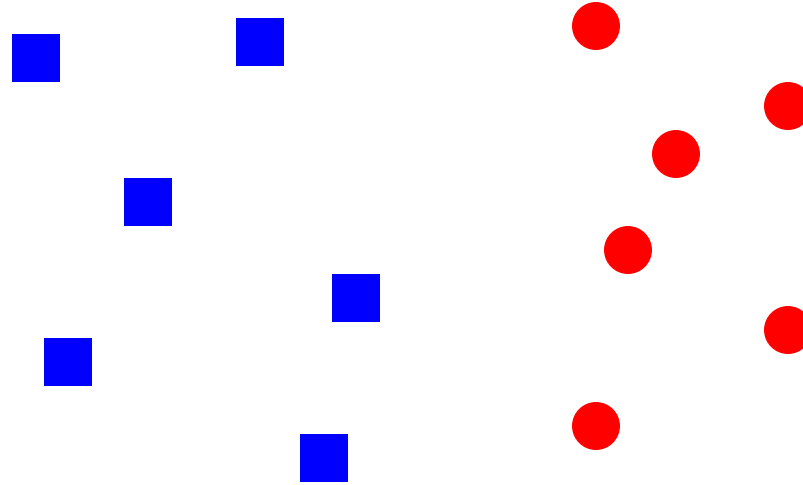
$$y = f(x)$$

output prediction Image
 function feature

The diagram illustrates the statistical learning framework equation $y = f(x)$. Three red arrows point from labels below to the equation: one from 'output' to y , one from 'prediction function' to f , and one from 'Image feature' to x .

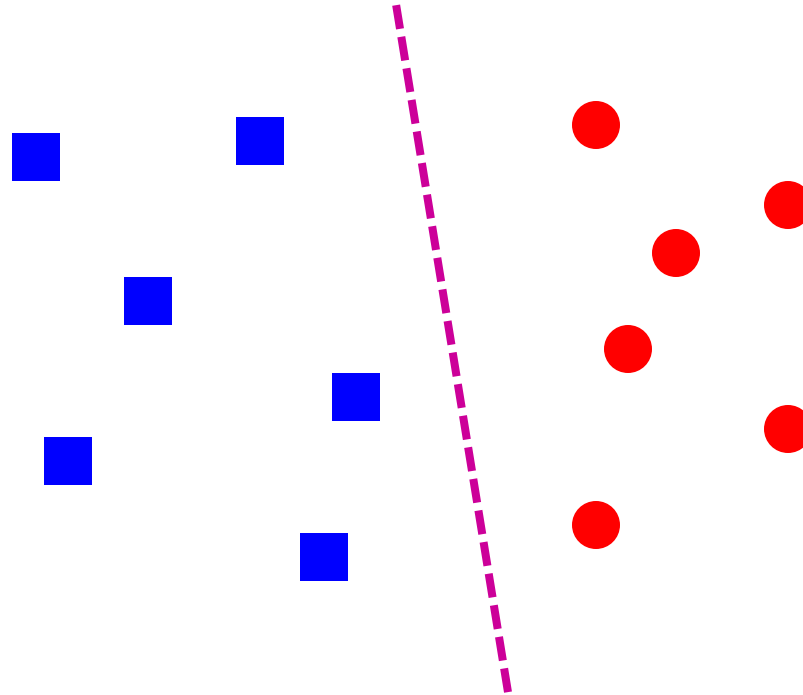
- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, estimate the prediction function f by minimizing the prediction error on the training set
- **Testing:** apply f to a never before seen *test example* \mathbf{x} and output the predicted value $y = f(\mathbf{x})$

LINEAR CLASSIFIERS — 2 CLASS PROBLEM



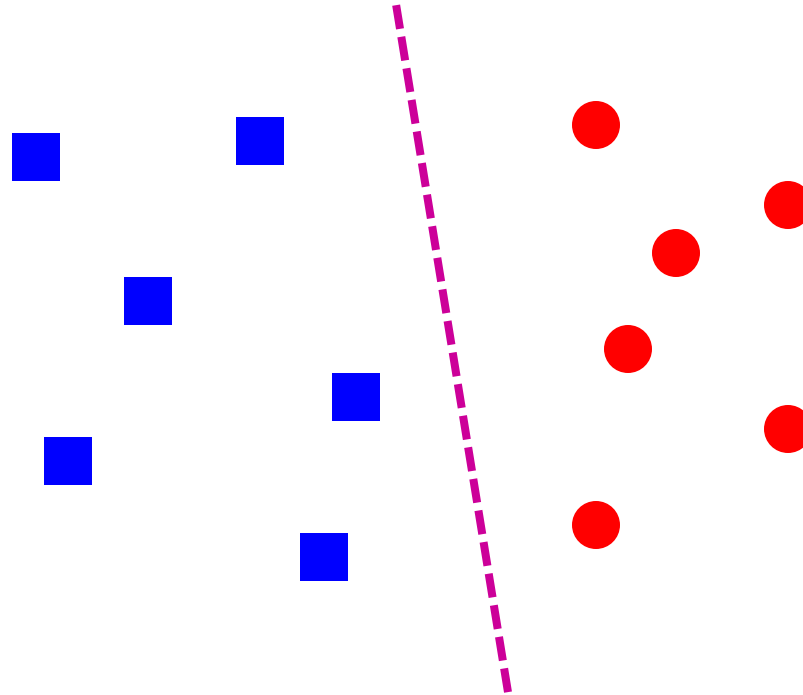
- Find a *linear function* to separate the classes:

LINEAR CLASSIFIERS — 2 CLASS PROBLEM



- Find a *linear function* to separate the classes:

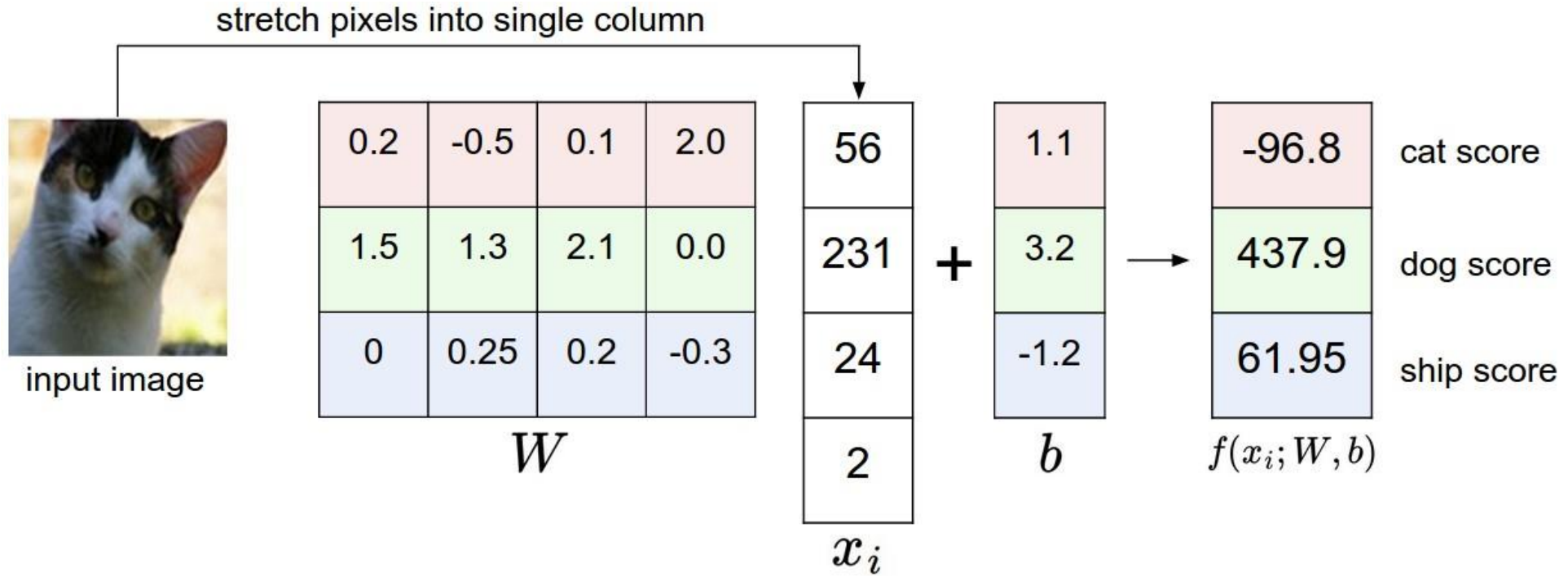
LINEAR CLASSIFIERS – 2 CLASS PROBLEM



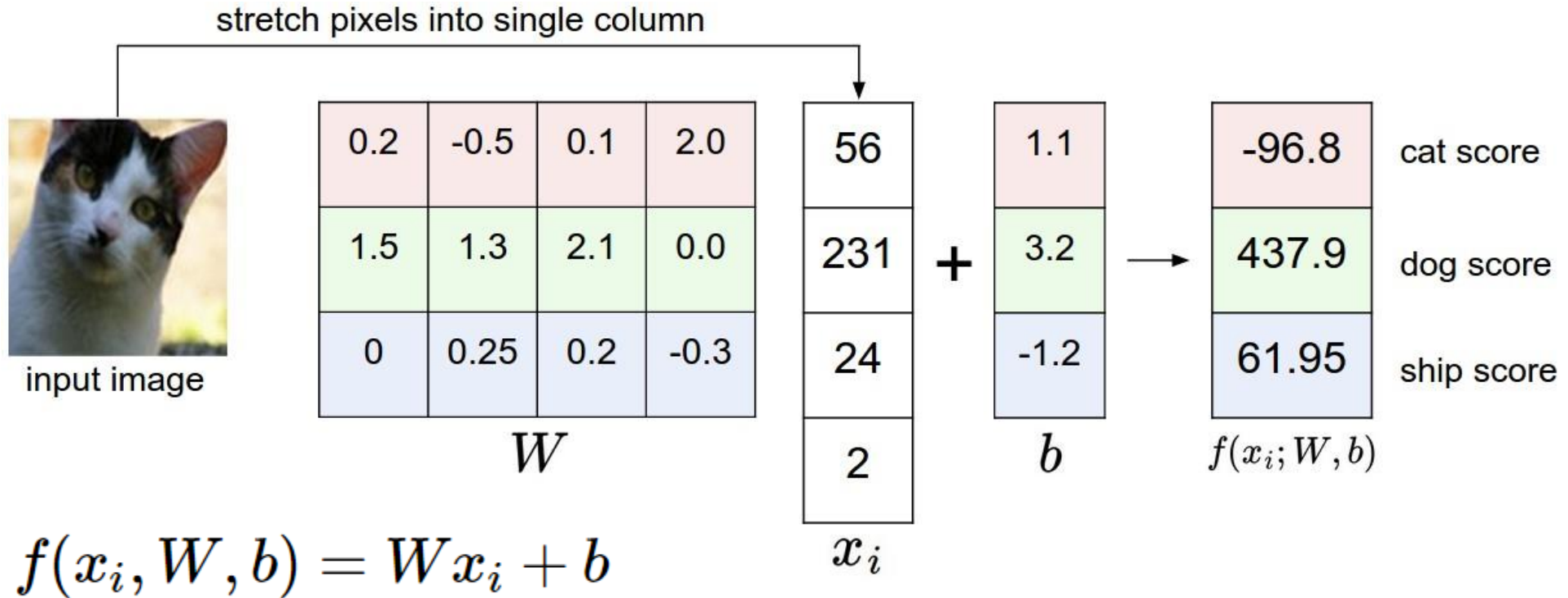
- Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

LINEAR CLASSIFIERS — MORE THAN 2 CLASS



LINEAR CLASSIFIERS – MORE THAN 2 CLASS



LINEAR CLASSIFIERS – MORE THAN 2 CLASS

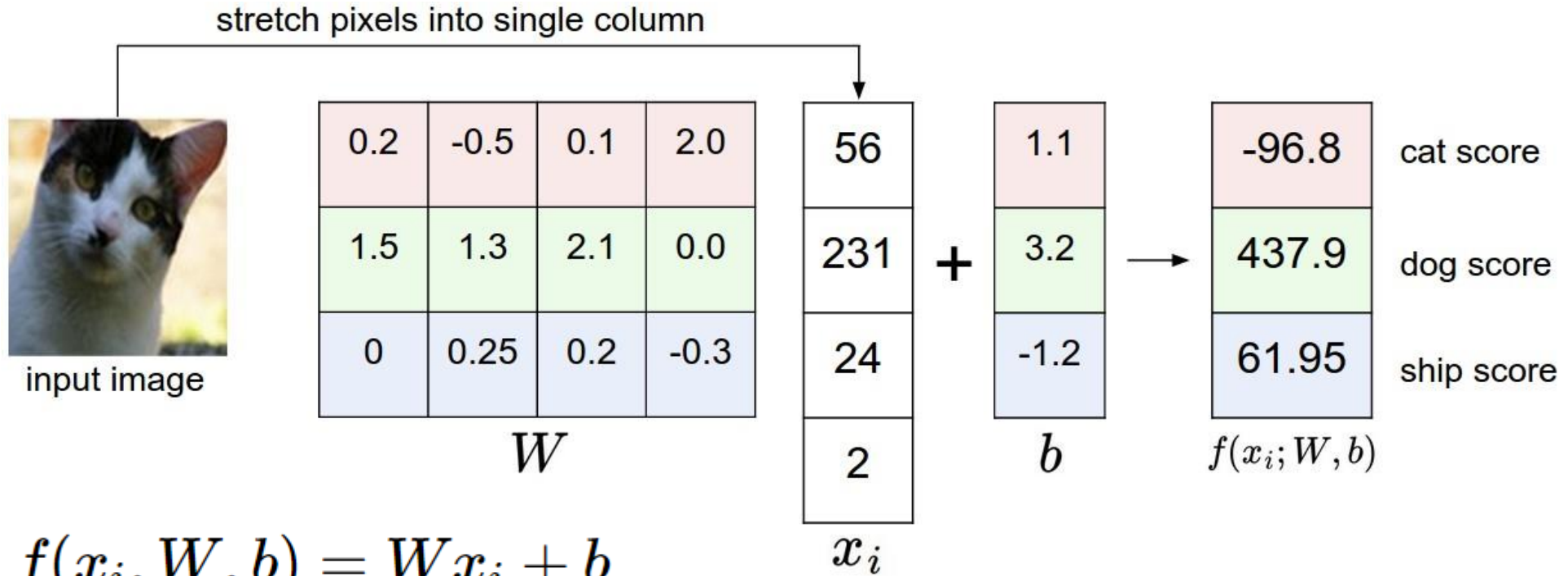
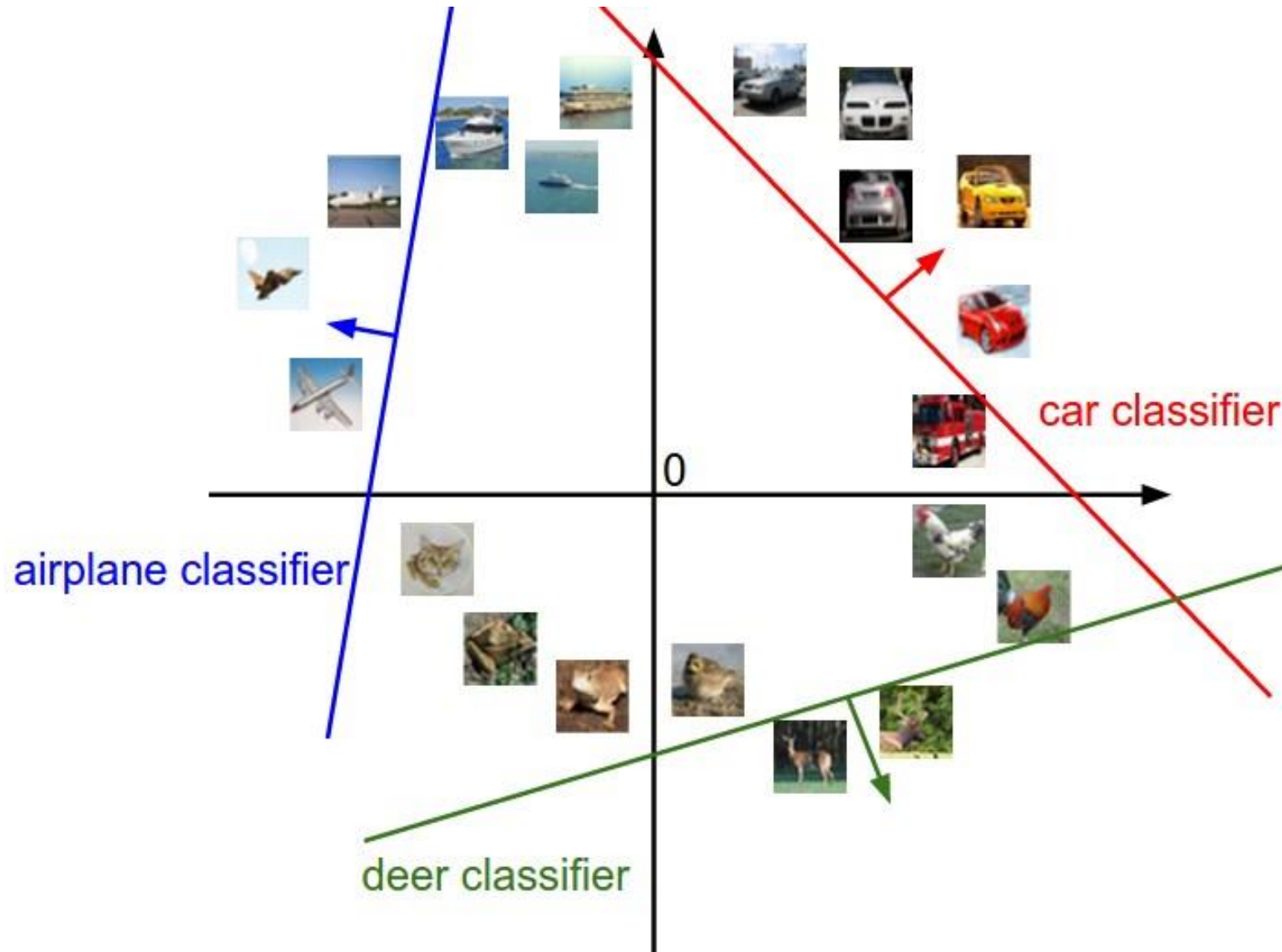


Image x_i has all of its pixels flattened out to a single column vector of shape $[D \times 1]$.

Matrix W (of size $[K \times D]$), and vector b (of size $[K \times 1]$) are the **parameters**.

K is the number of classes.

ANALOGY OF IMAGES AS HIGH-DIMENSIONAL POINTS



Source: cs231n, <http://cs231n.github.io/linear-classify/>

NEAREST NEIGHBOR VS. LINEAR CLASSIFIERS

- **Linear pros:**

- Low-dimensional *parametric* representation
- Very fast at test time

- **Linear cons:**

- How to train the linear function?
- What if data is not linearly separable?

SOFTMAX CLASSIFIER

- Interprets the class scores as the unnormalized log probabilities for each class and replace the *hinge loss* with a **cross-entropy loss** that has the form:

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) = -s_{y_i} + \log \sum_j e^{s_j}$$

SOFTMAX CLASSIFIER

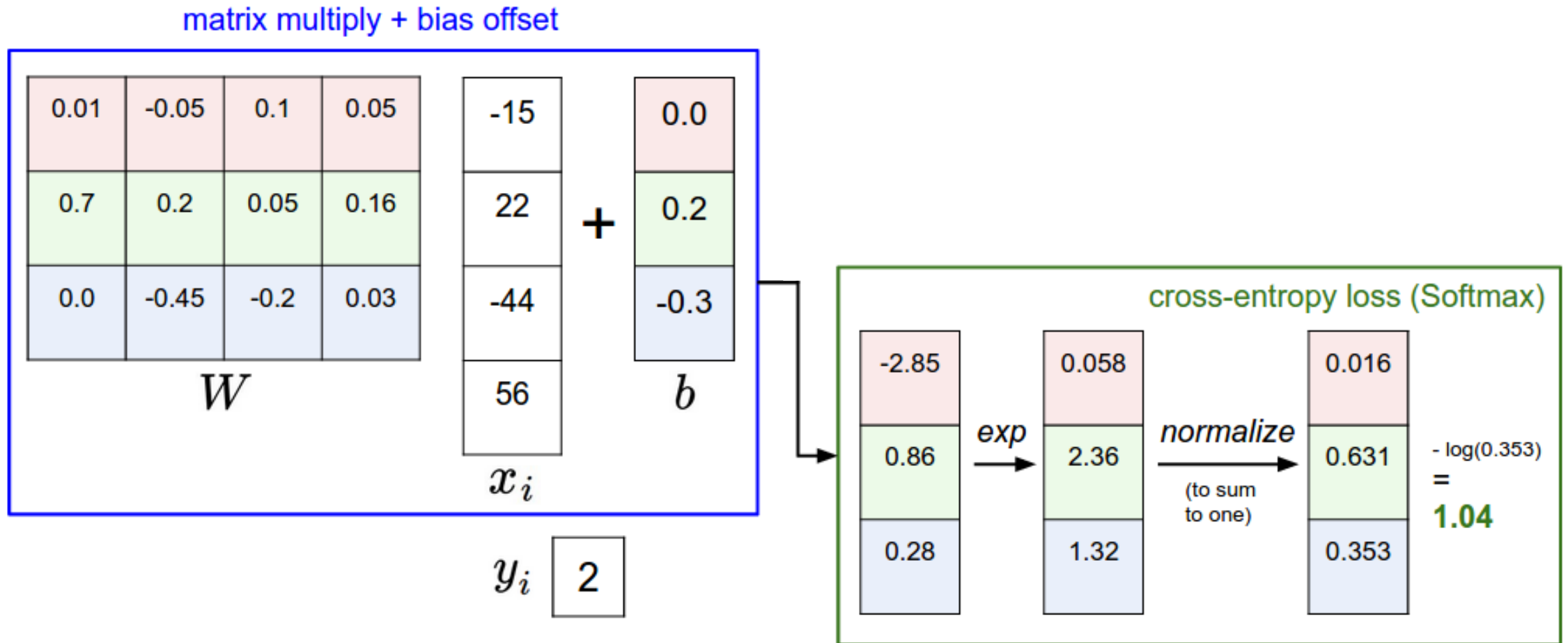
- Interprets the class scores as the unnormalized log probabilities for each class and replace the *hinge loss* with a **cross-entropy loss** that has the form:

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) = -s_{y_i} + \log \sum_j e^{s_j}$$

- Softmax Loss:

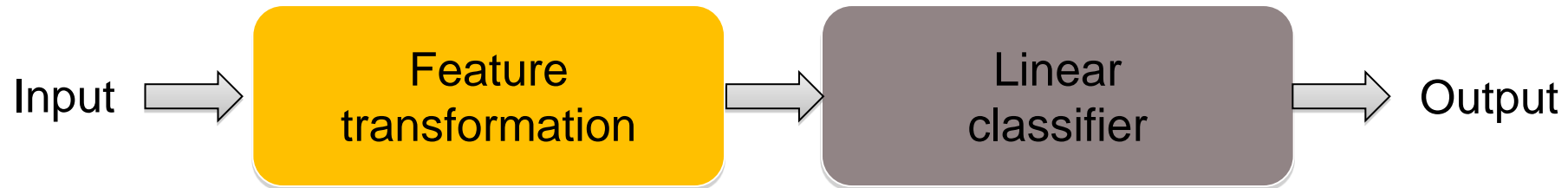
$$L = \underbrace{\frac{1}{N} \sum_i L_i}_{\text{data loss}} + \underbrace{\lambda R(W)}_{\text{regularization loss}}$$

HINGE VS CROSS-ENTROPY LOSS



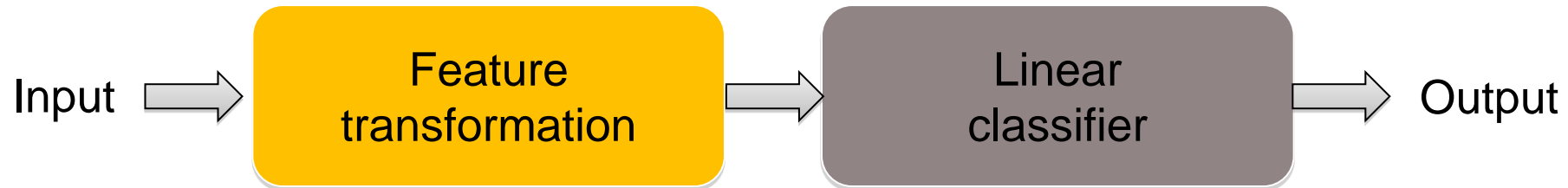
FROM LINEAR TO NONLINEAR CLASSIFIERS

- To achieve good accuracy on challenging problems, we need to be able to train *nonlinear* models
- Two strategies for making nonlinear predictors out of linear ones:
 - **“Shallow” approach:** nonlinear feature transformation followed by linear classifier

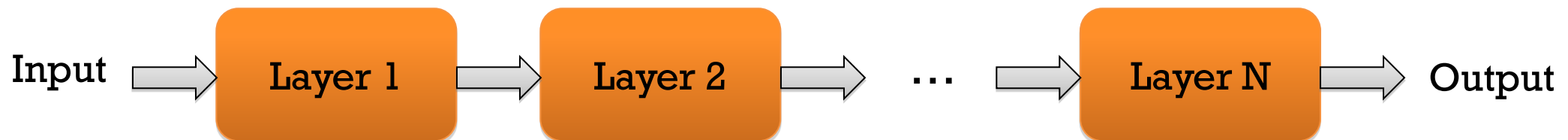


FROM LINEAR TO NONLINEAR CLASSIFIERS

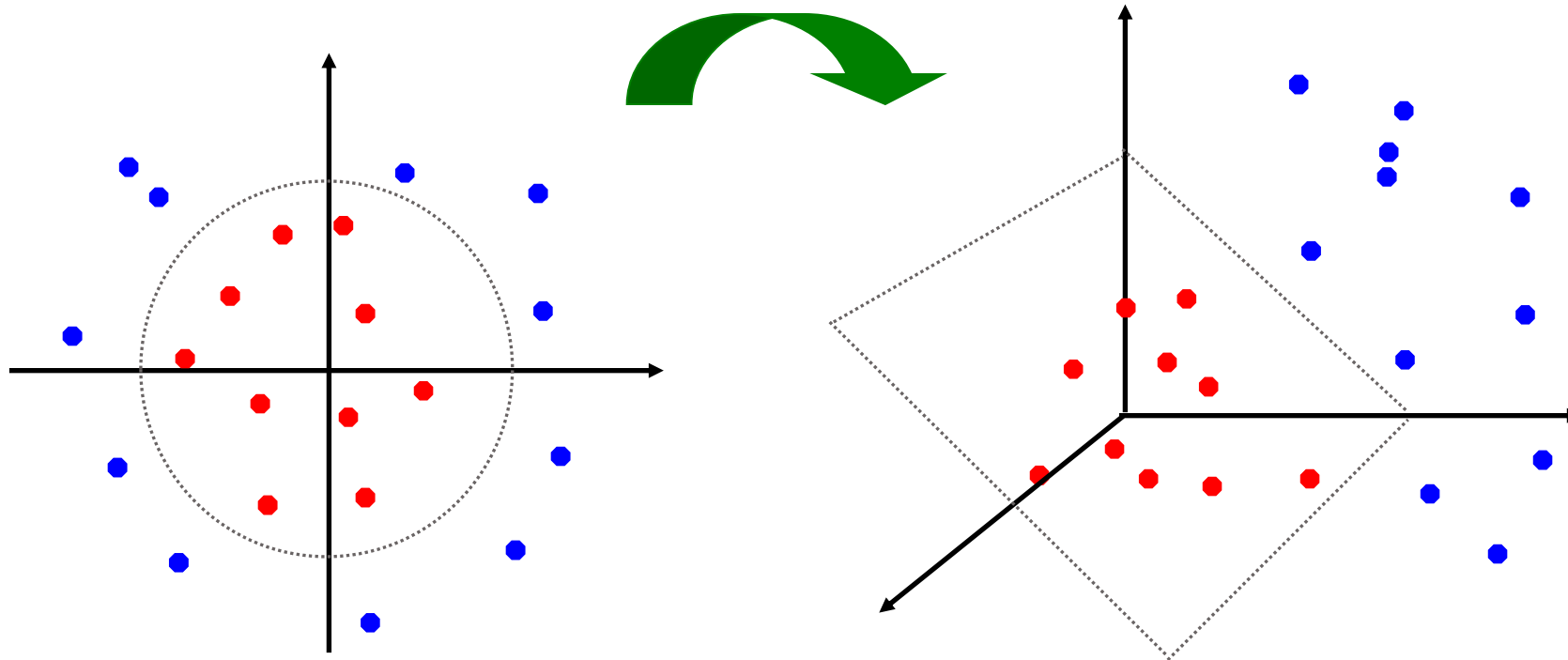
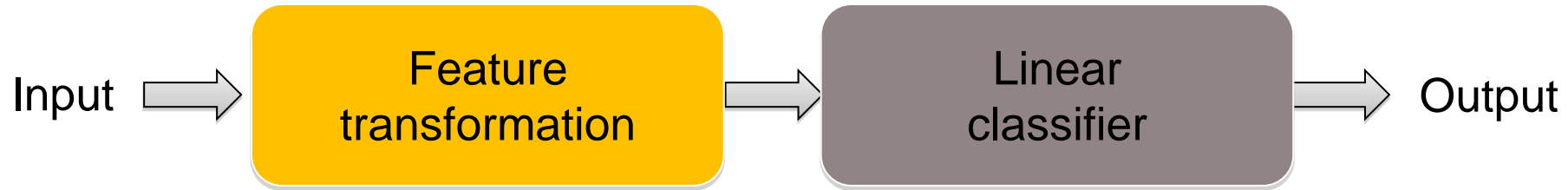
- To achieve good accuracy on challenging problems, we need to be able to train *nonlinear* models
- Two strategies for making nonlinear predictors out of linear ones:
 - **“Shallow” approach:** nonlinear feature transformation followed by linear classifier



- **“Deep” approach:** stack multiple layers of linear predictors (interspersed with nonlinearities)

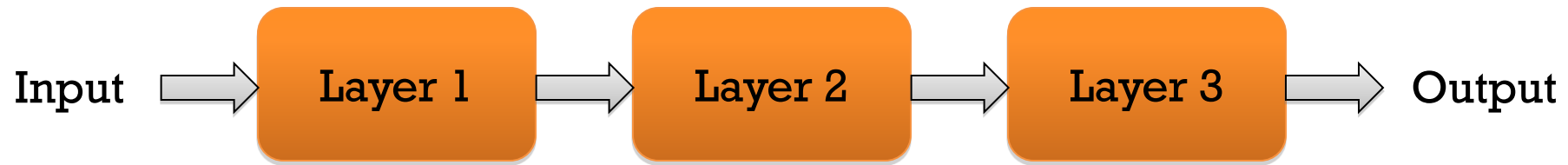


SHALLOW APPROACH



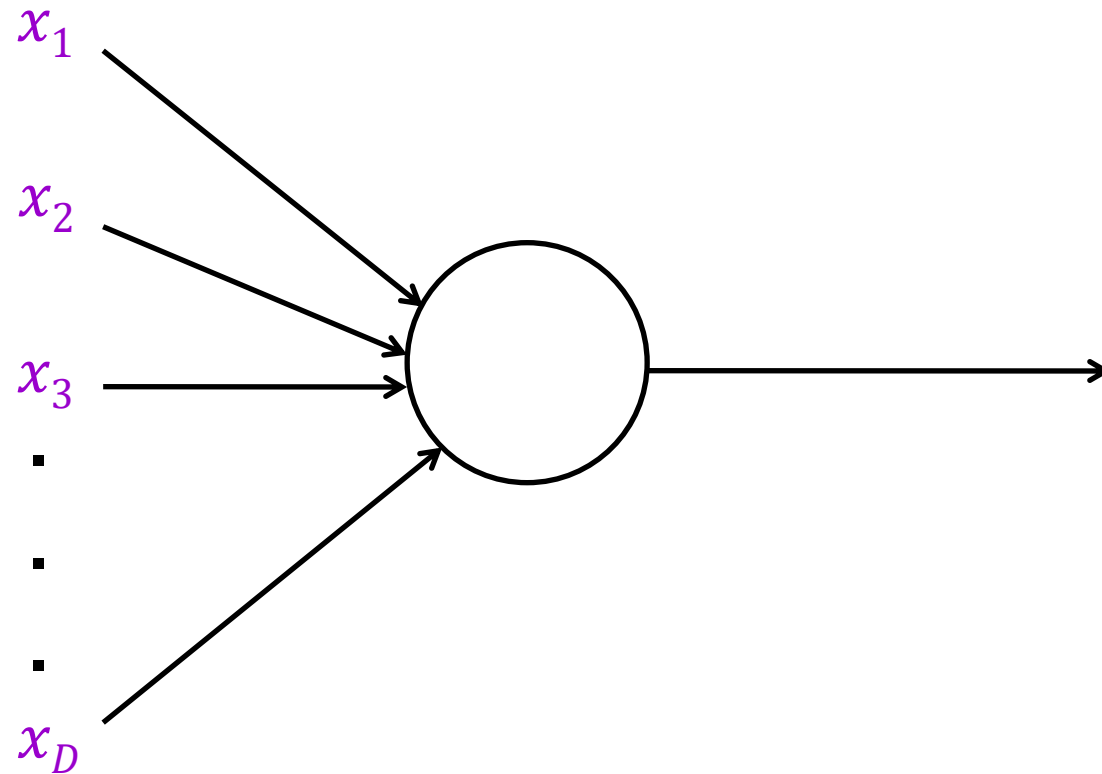
MULTI-LAYER NEURAL NETWORKS

- “Deep” approach: stack multiple layers of linear predictors (perceptrons) interspersed with nonlinearities

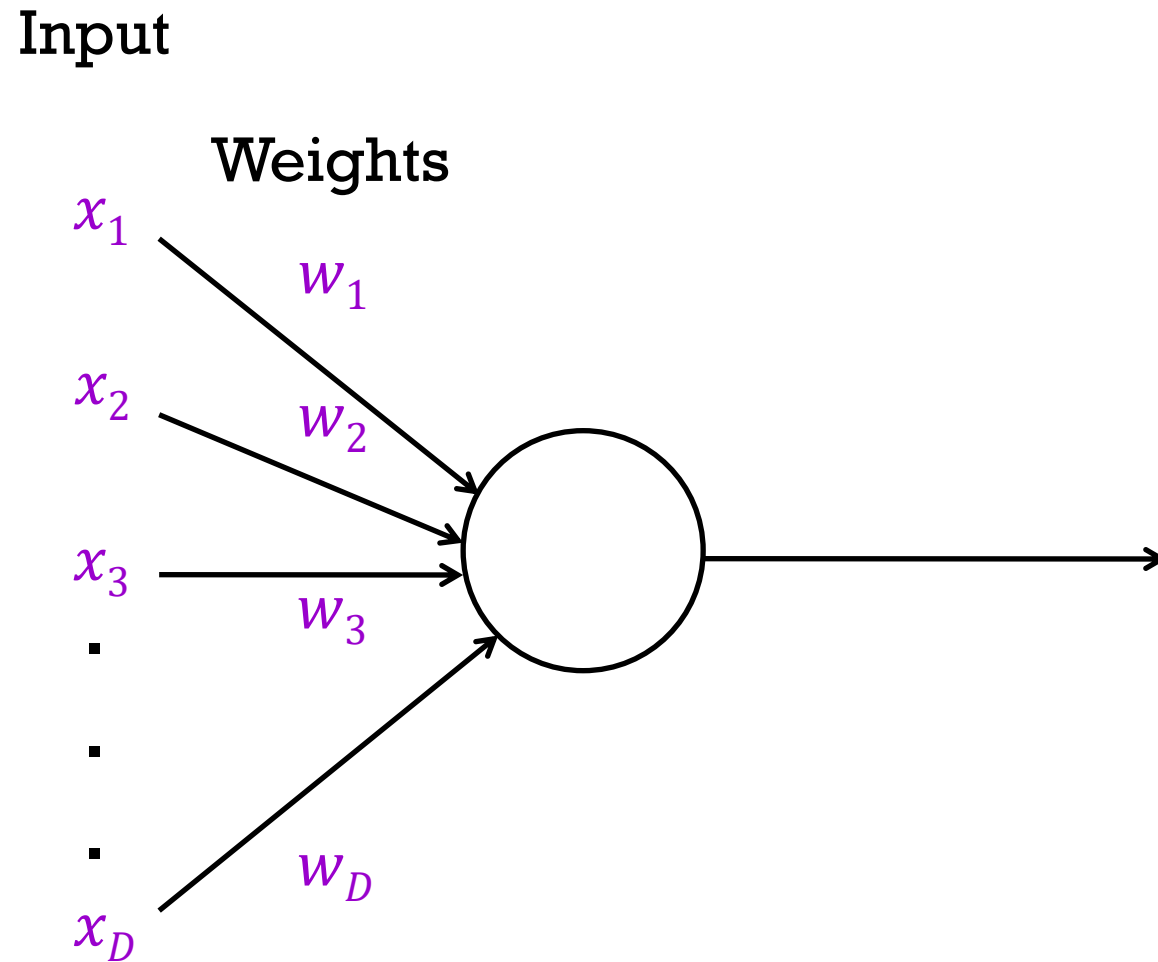


PERCEPTRON

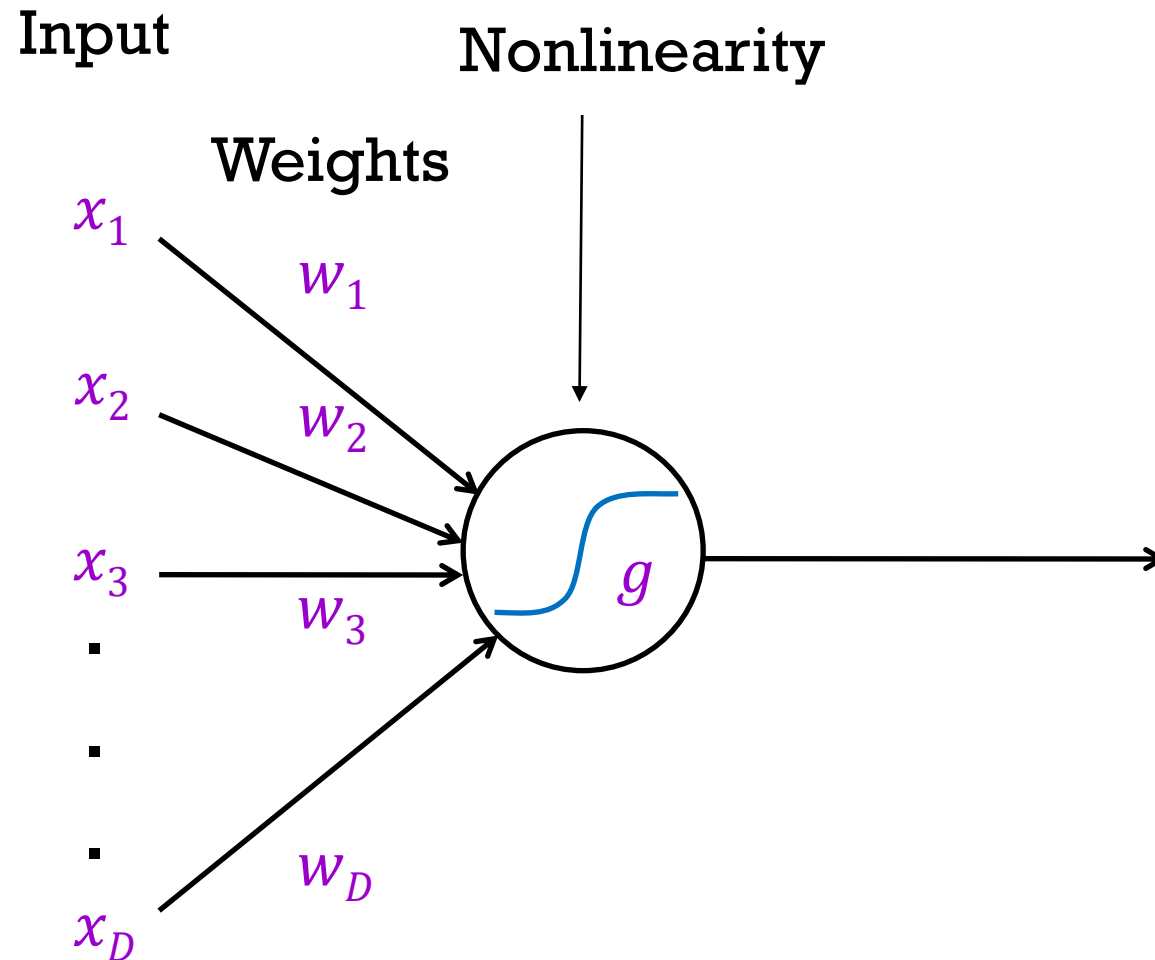
Input



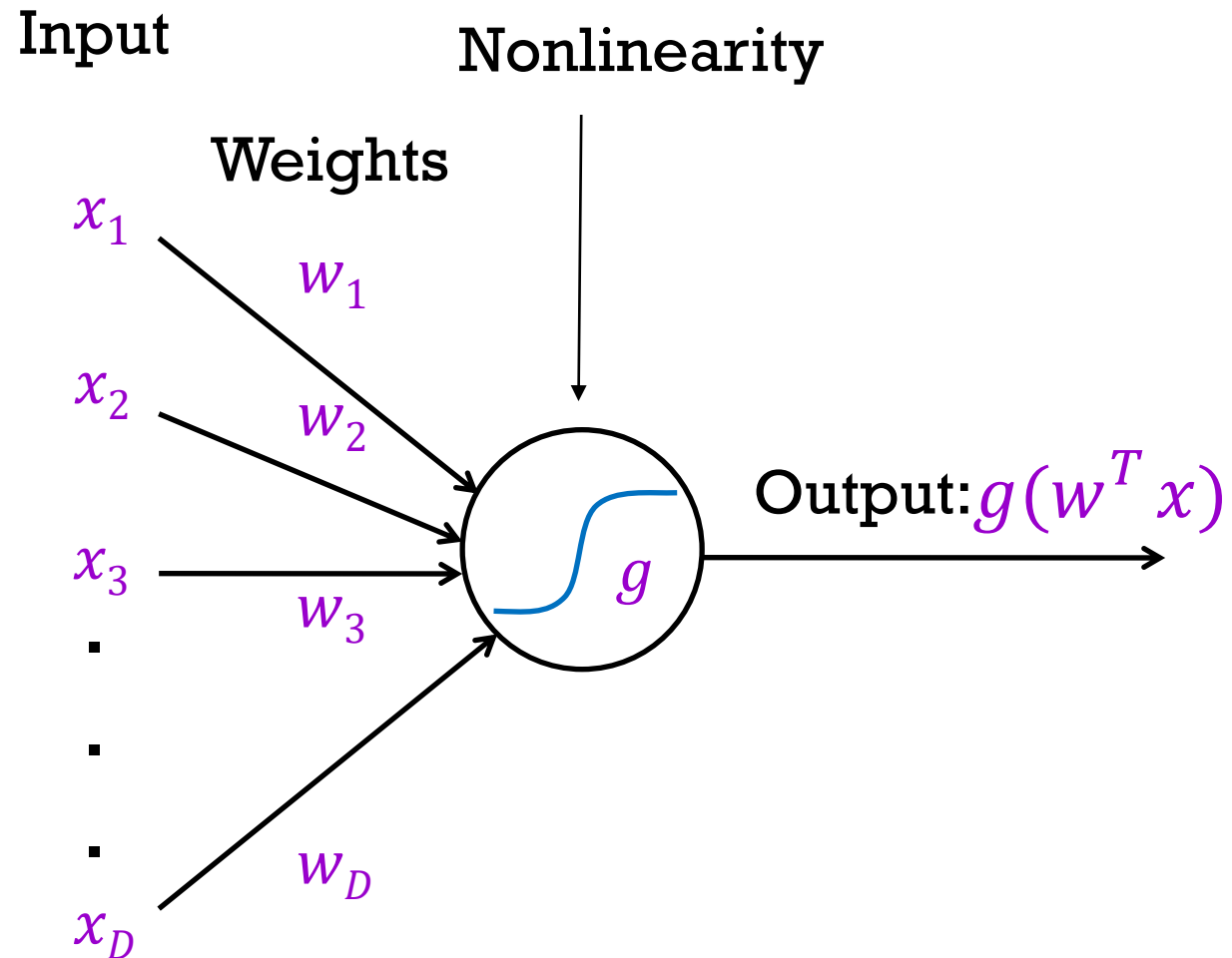
PERCEPTRON



PERCEPTRON



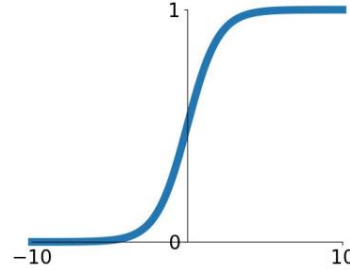
PERCEPTRON



COMMON NONLINEARITIES (OR *ACTIVATION FUNCTIONS*)

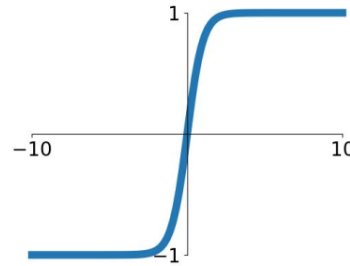
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



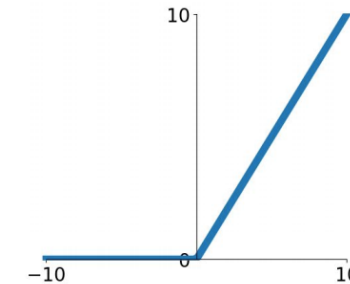
tanh

$$\tanh(x)$$



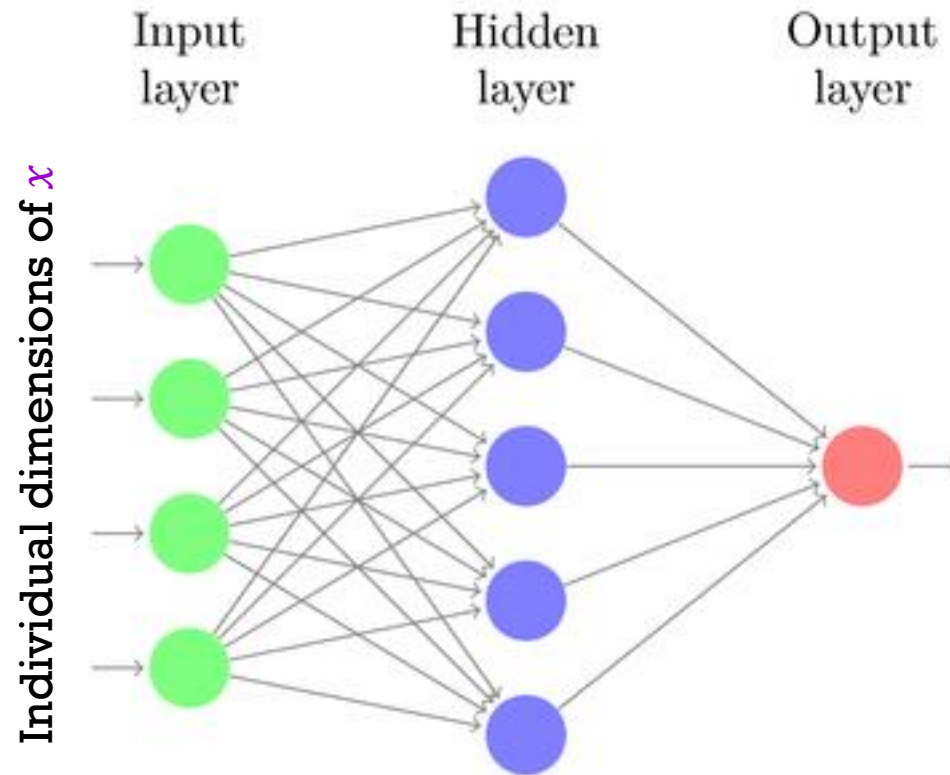
ReLU

$$\max(0, x)$$



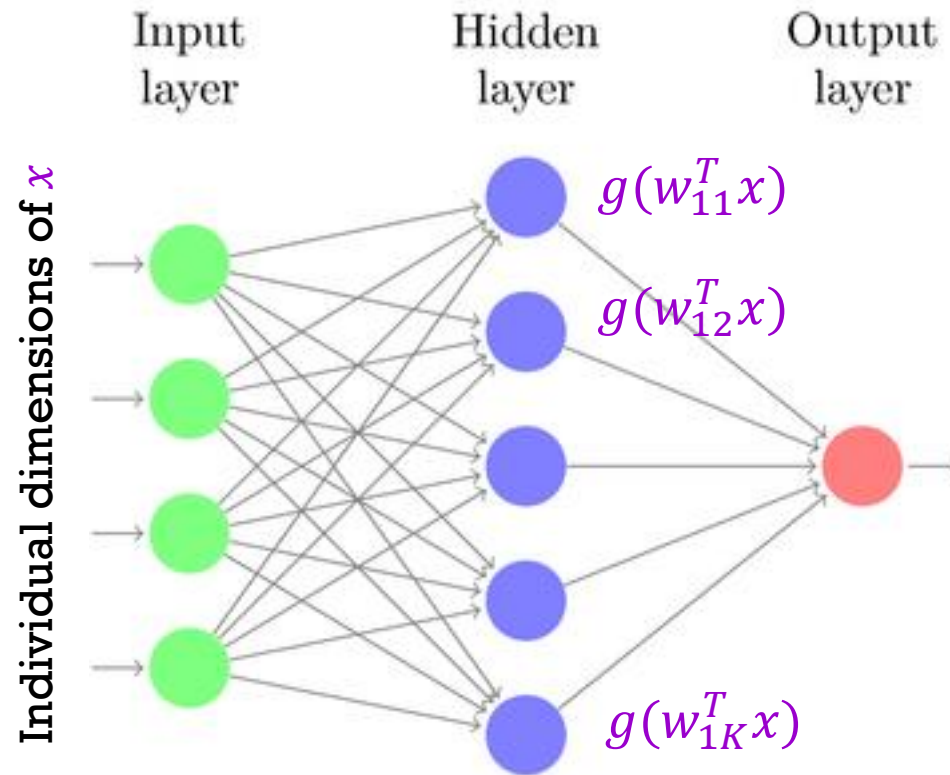
TWO-LAYER NEURAL NETWORK

- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities



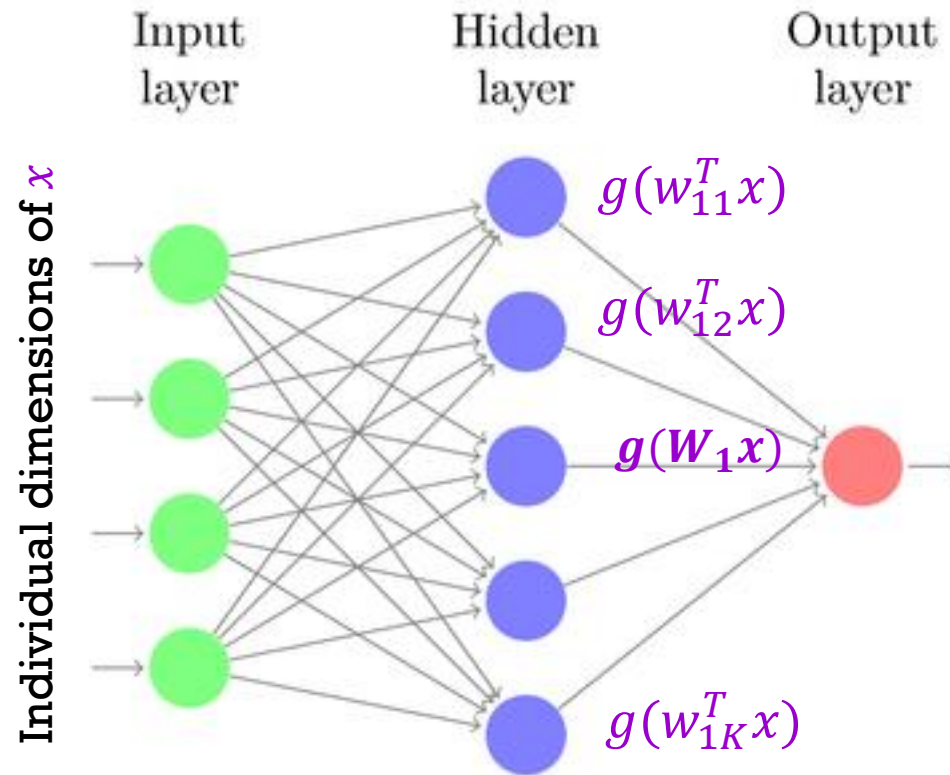
TWO-LAYER NEURAL NETWORK

- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities



TWO-LAYER NEURAL NETWORK

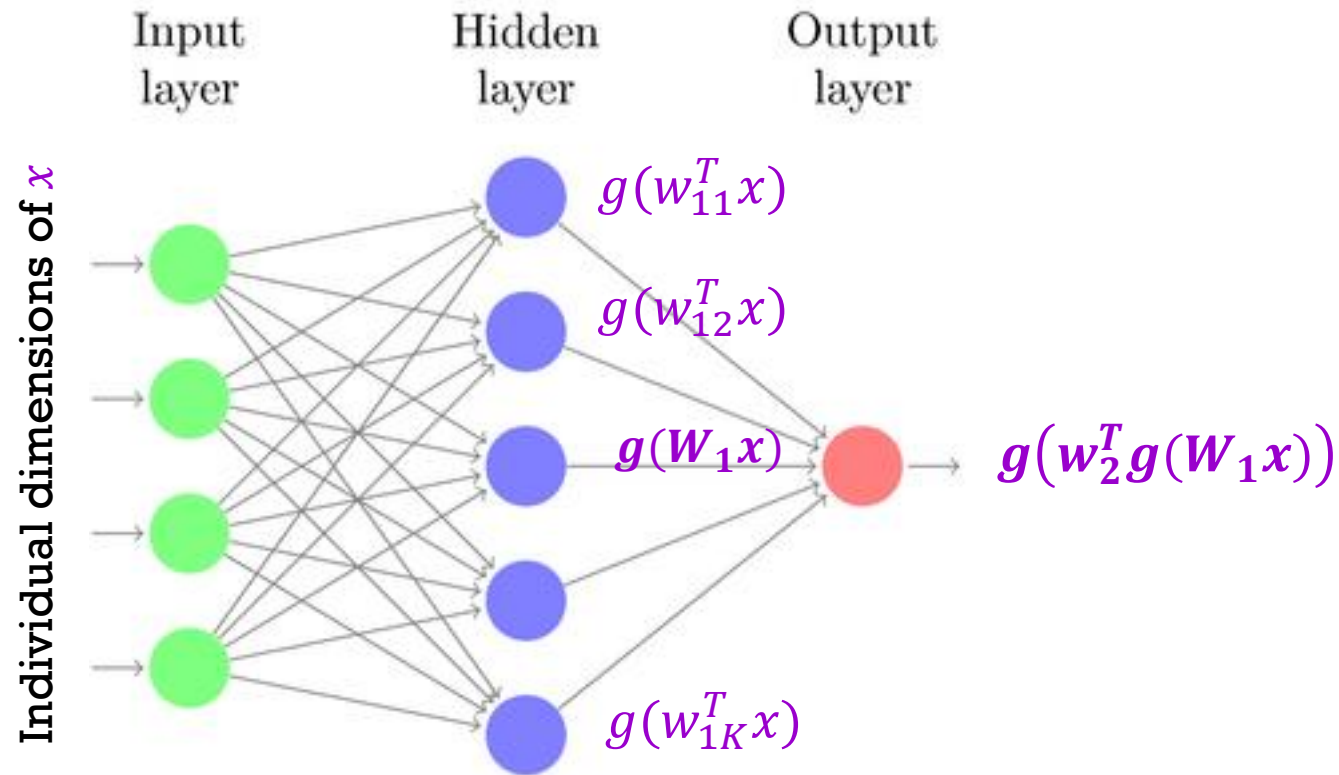
- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities



W_1 – matrix with rows w_{1j}^T

TWO-LAYER NEURAL NETWORK

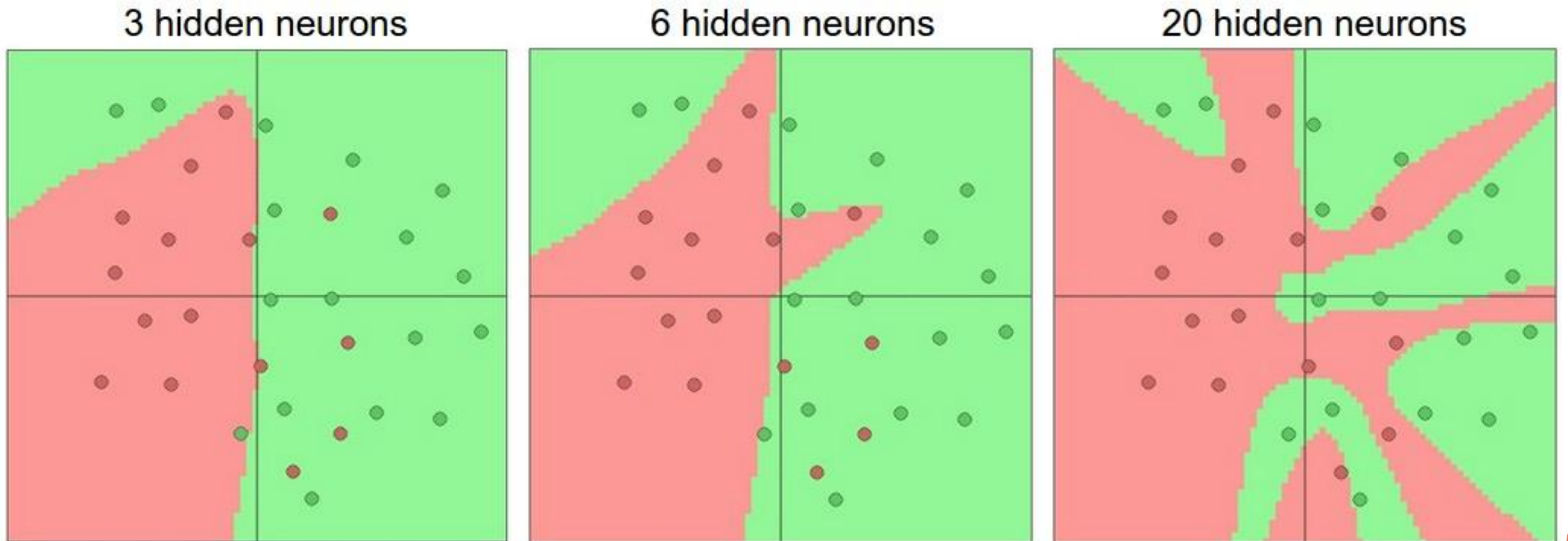
- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities



W_1 – matrix with rows w_{1j}^T

TWO-LAYER NEURAL NETWORK

- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities
- The bigger the hidden layer, the more expressive the model



TWO-LAYER NEURAL NETWORK

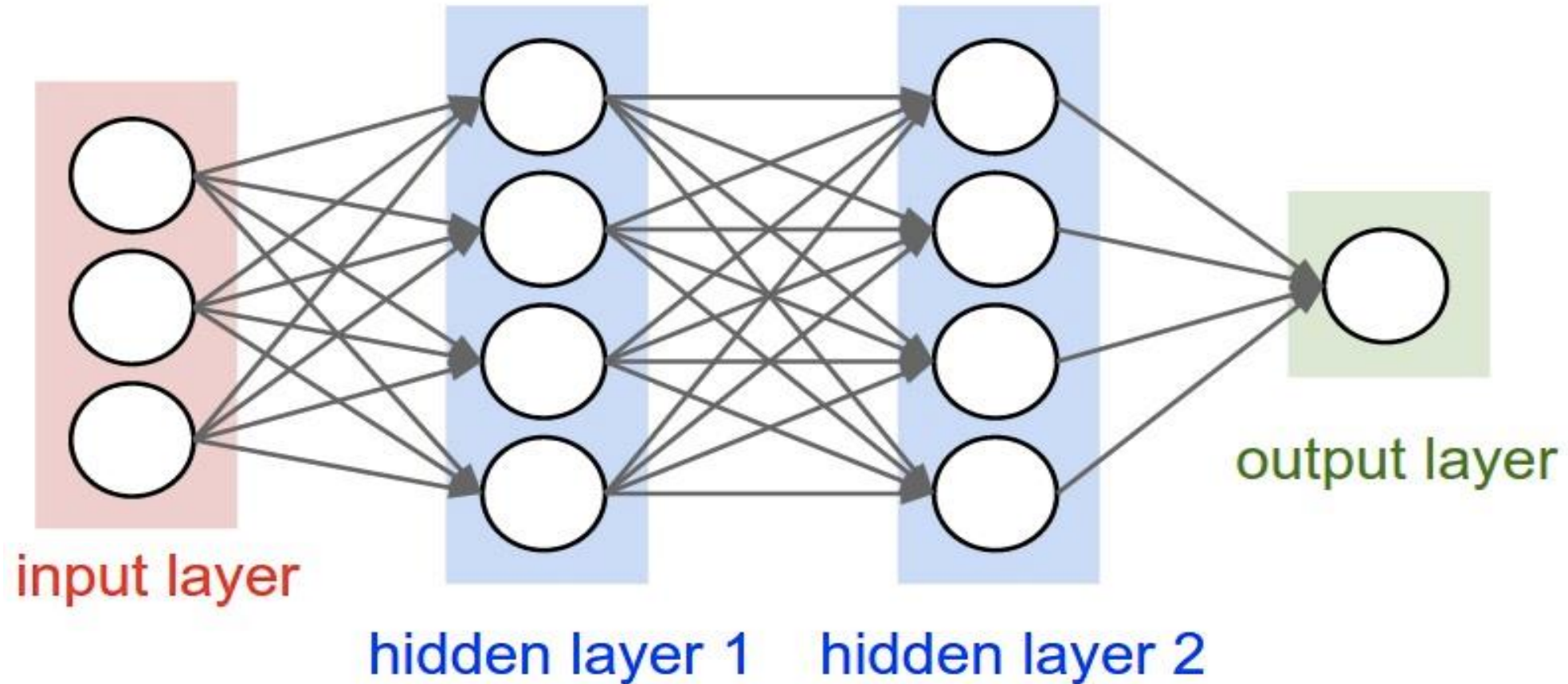
- Introduce a *hidden layer* of perceptrons computing linear combinations of inputs followed by nonlinearities
- The bigger the hidden layer, the more expressive the model



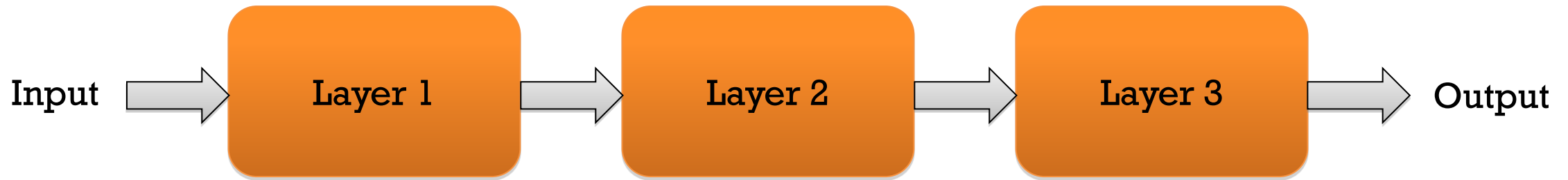
- A two-layer network is a universal function approximator
- But the hidden layer may need to be huge



BEYOND TWO LAYERS



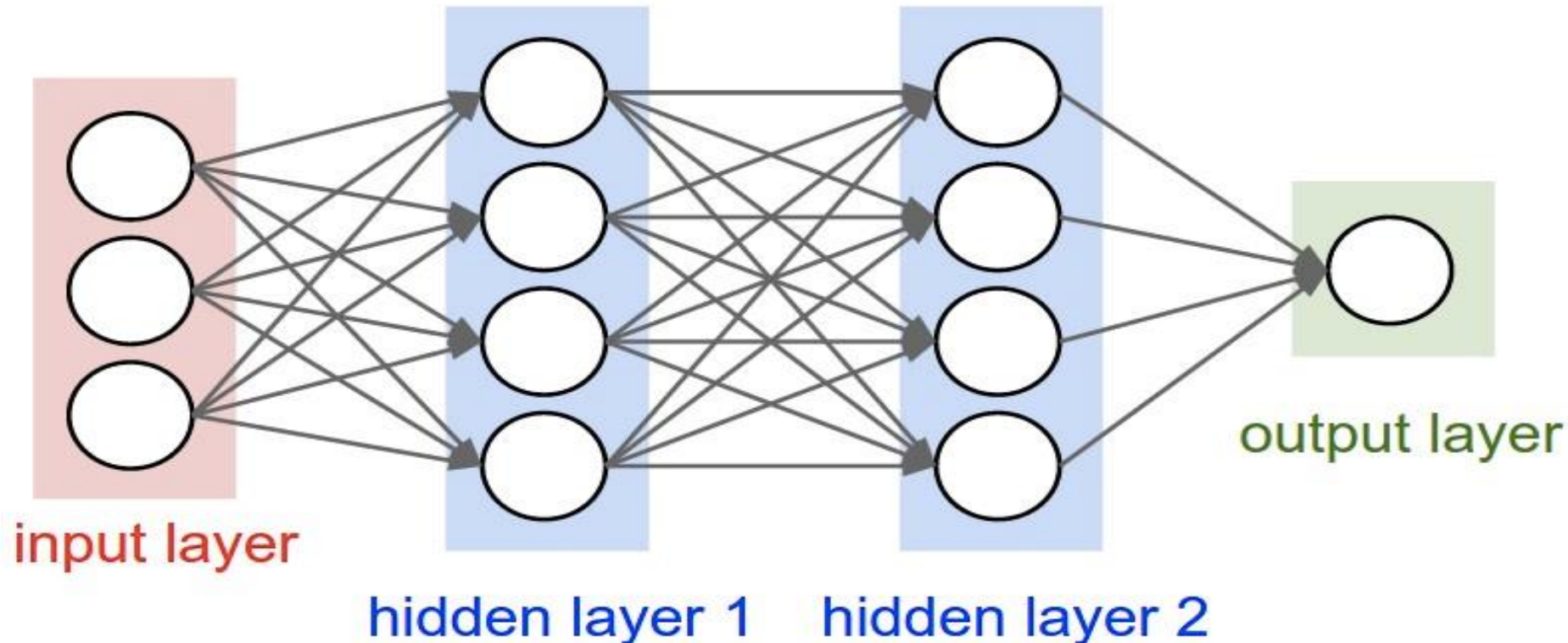
“DEEP” PIPELINE



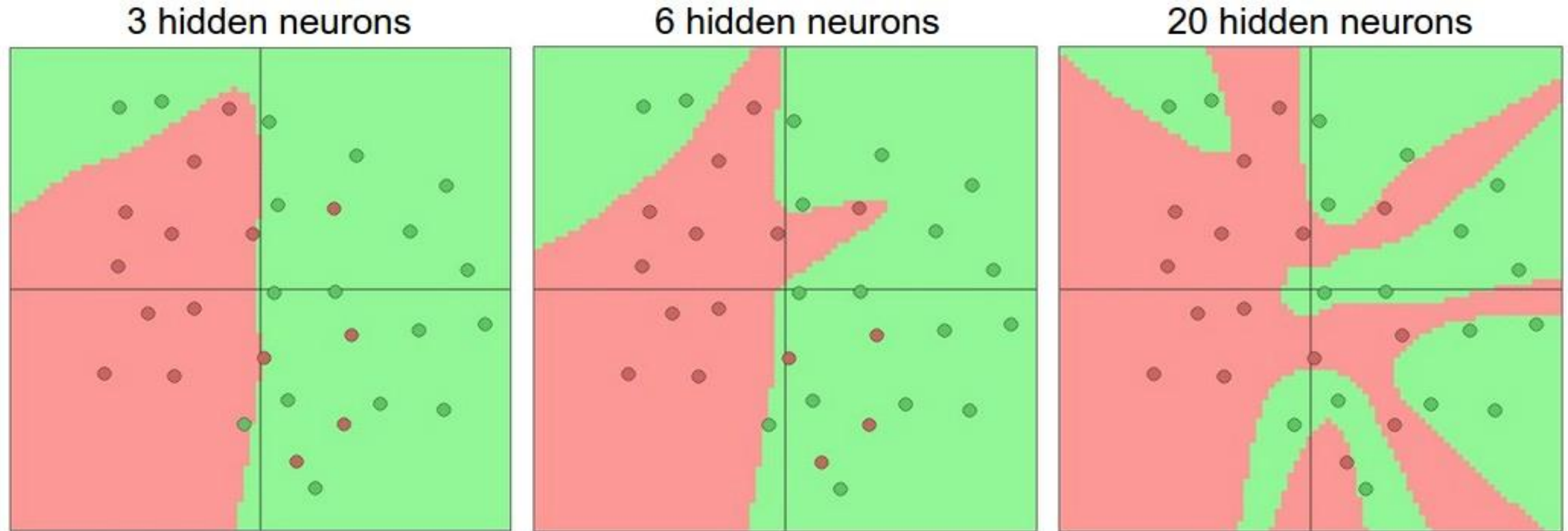
- Learn a *feature hierarchy*
- Each layer extracts features from the output of previous layer
- All layers are trained jointly

HYPERPARAMETERS IN MULTI-LAYER NETWORKS

- Number of layers, number of units per layer, activation function, loss function, regularization, regularization constant, optimizer, optimizer settings: learning rate schedule, number of epochs, batch size, etc.

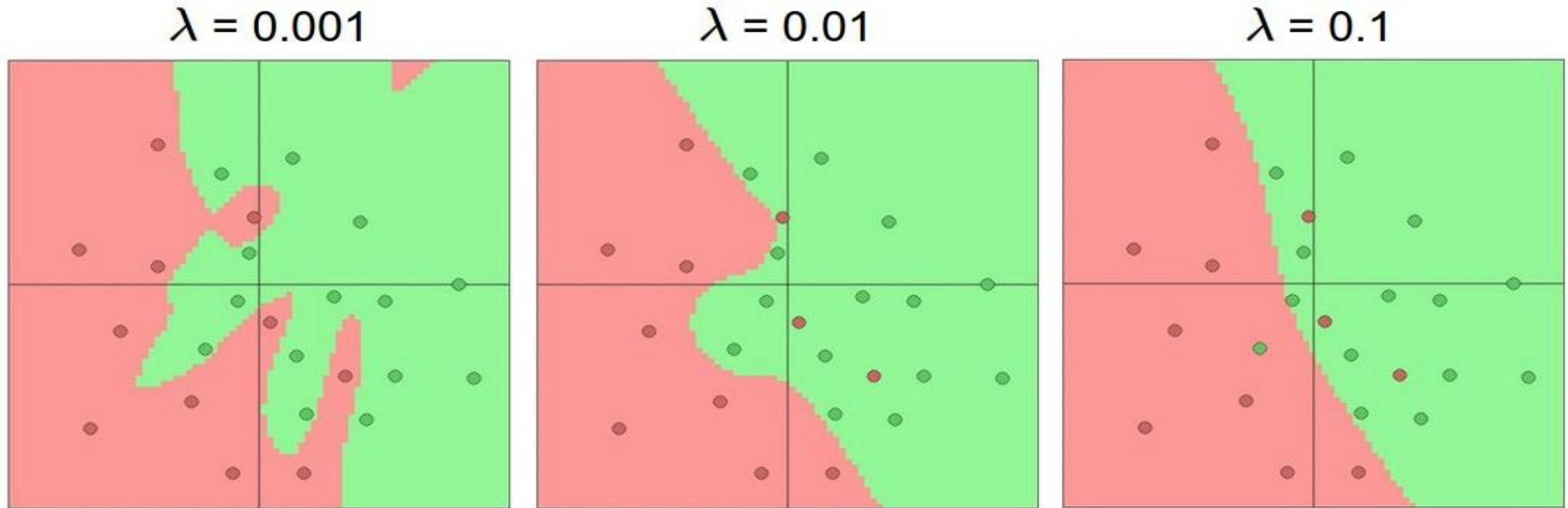


HYPERPARAMETERS IN MULTI-LAYER NETWORKS



Number of hidden units in a two-layer network

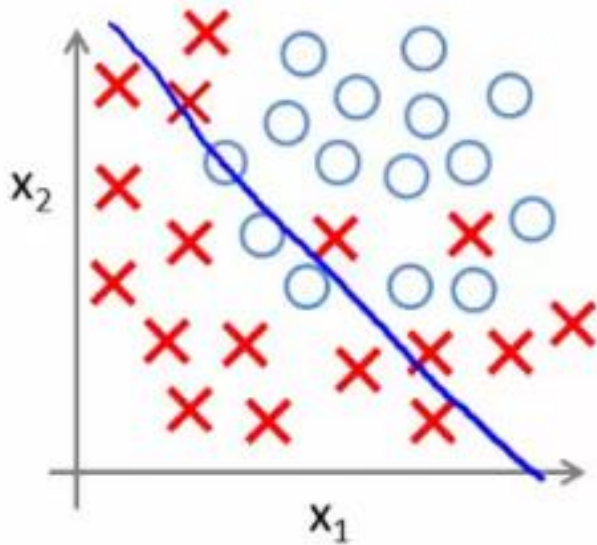
HYPERPARAMETERS IN MULTI-LAYER NETWORKS



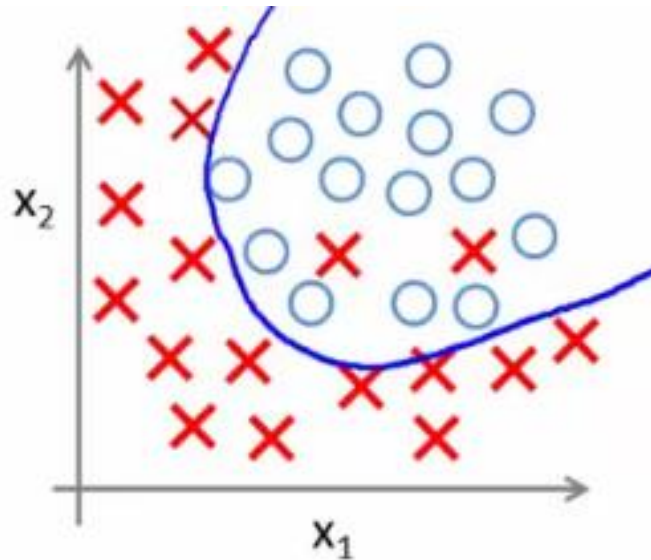
Regularization constant

MODEL COMPLEXITY AND GENERALIZATION

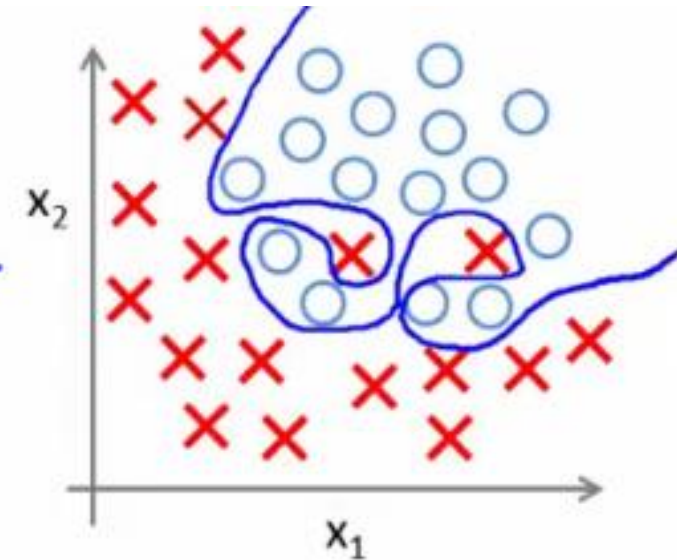
“Simple” model



“Intermediate” model



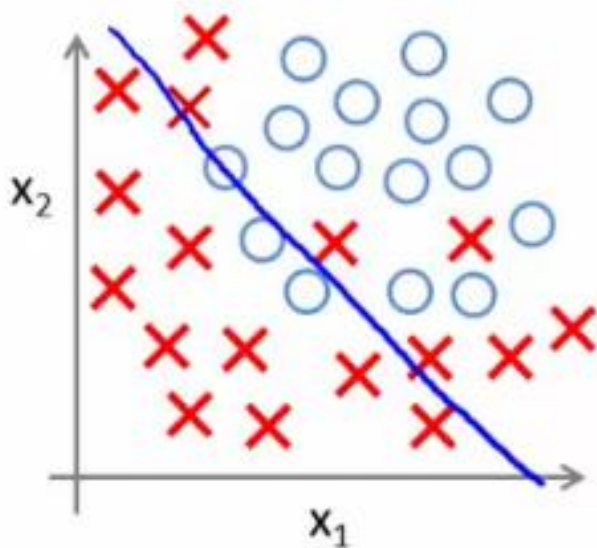
“Complex” model



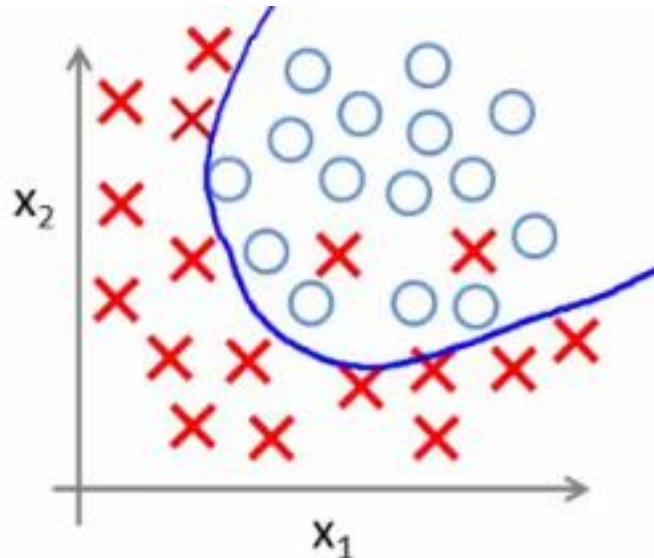
MODEL COMPLEXITY AND GENERALIZATION

- Generalization (test) error of learning algorithms has two main components:
 - **Bias:** error due to simplifying model assumptions
 - **Variance:** error due to randomness of training set

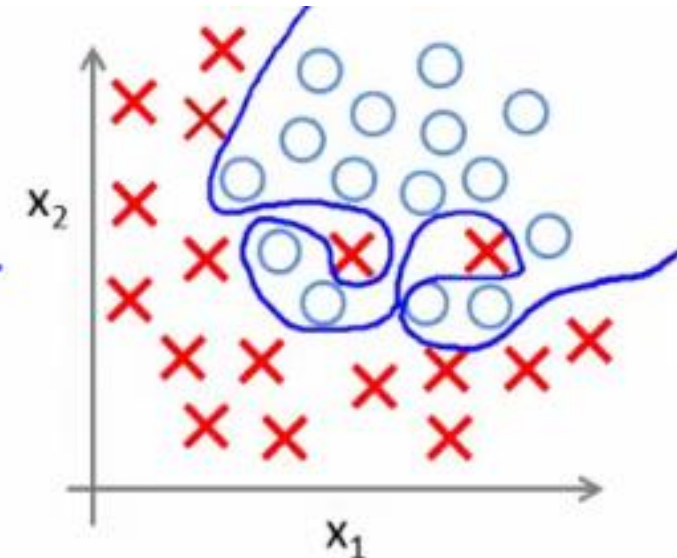
“Simple” model



“Intermediate” model

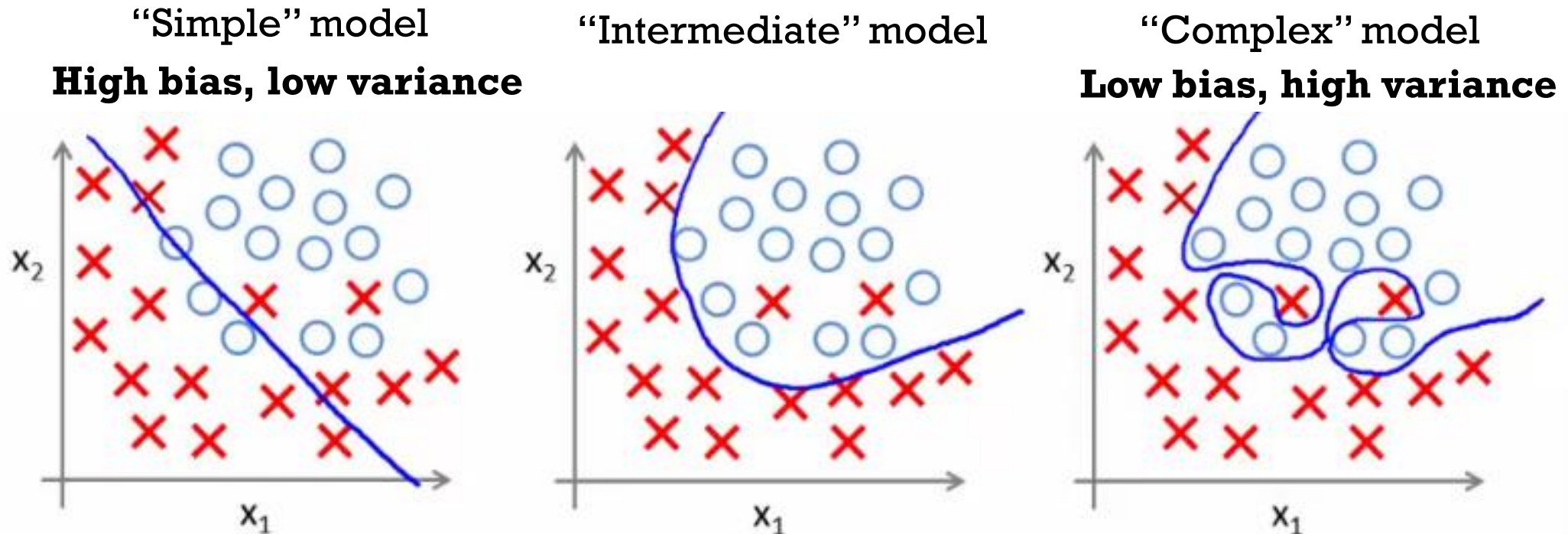


“Complex” model



MODEL COMPLEXITY AND GENERALIZATION

- Generalization (test) error of learning algorithms has two main components:
 - **Bias:** error due to simplifying model assumptions
 - **Variance:** error due to randomness of training set



BIAS-VARIANCE TRADEOFF

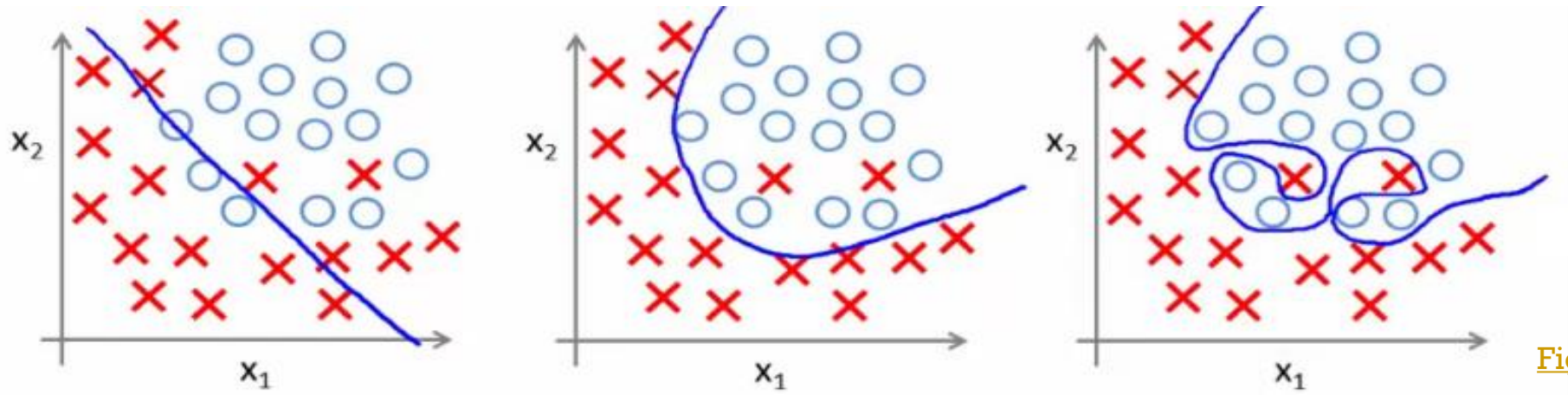
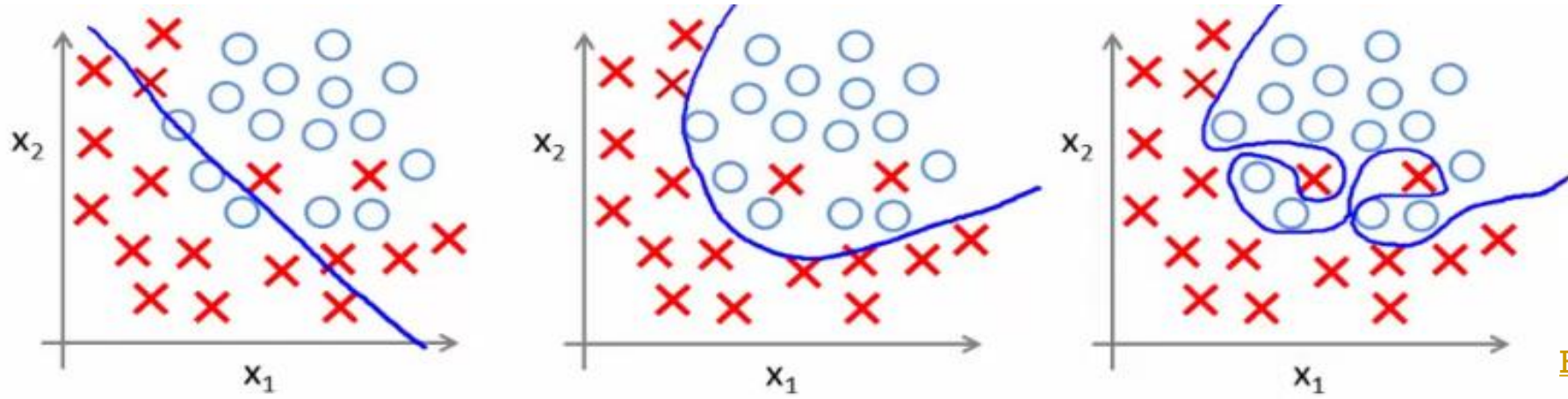


Figure source

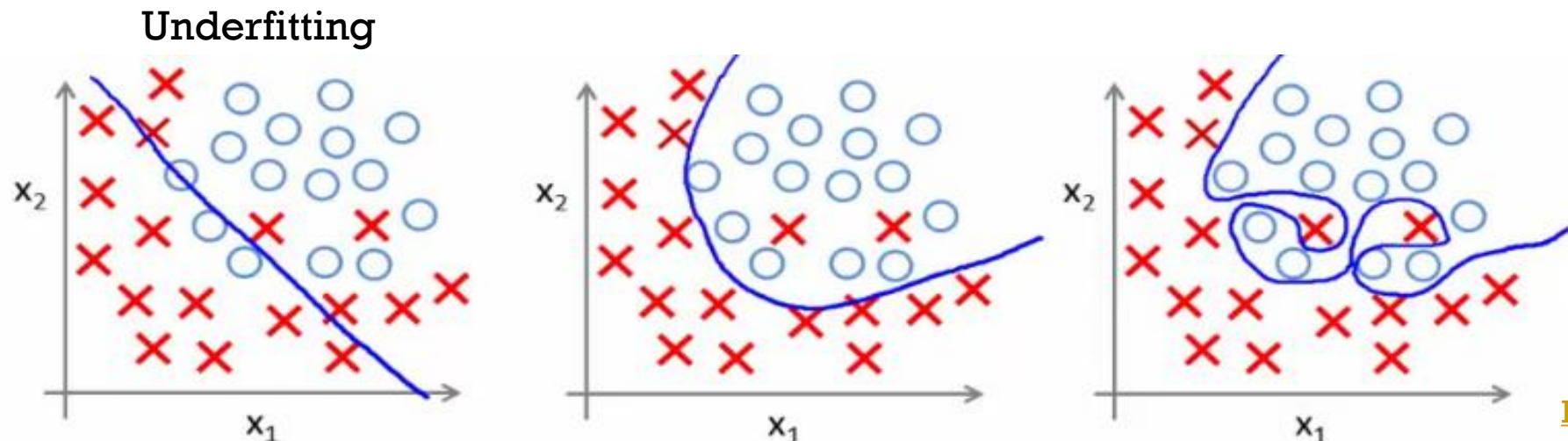
BIAS-VARIANCE TRADEOFF

- What if your model **bias** is too high?



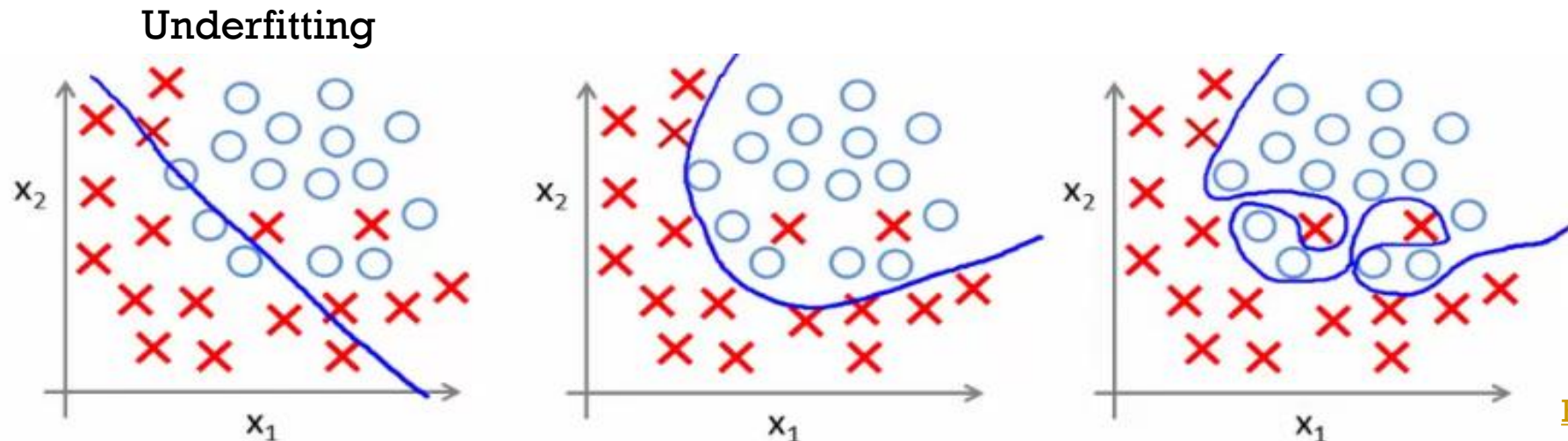
BIAS-VARIANCE TRADEOFF

- What if your model **bias** is too high?
 - Your model is **underfitting** – it is incapable of capturing the important characteristics of the training data



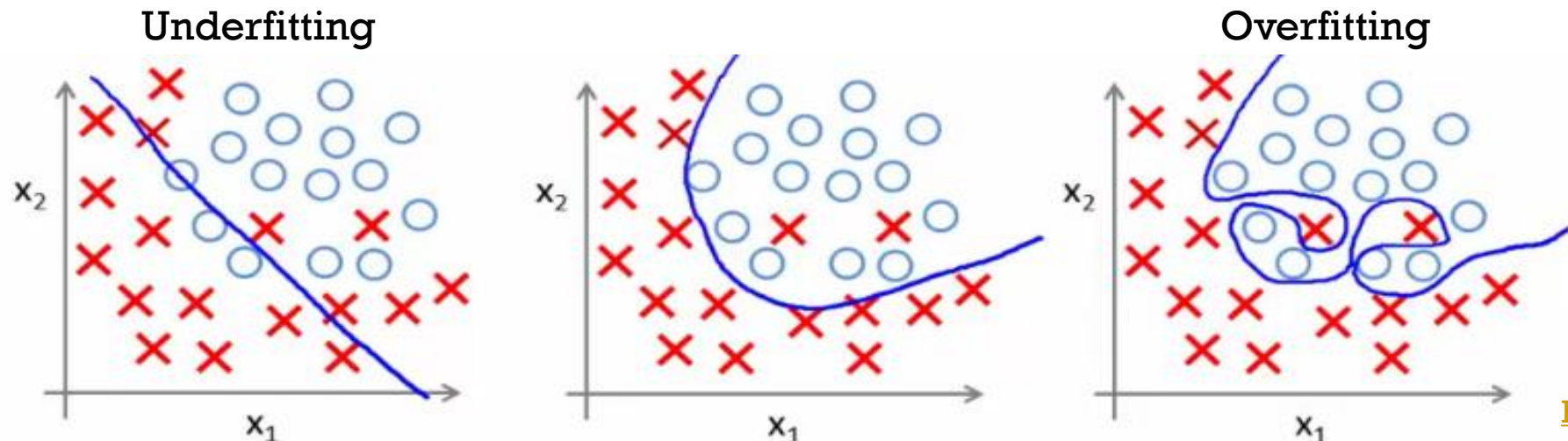
BIAS-VARIANCE TRADEOFF

- What if your model **bias** is too high?
 - Your model is **underfitting** – it is incapable of capturing the important characteristics of the training data
- What if your model **variance** is too high?



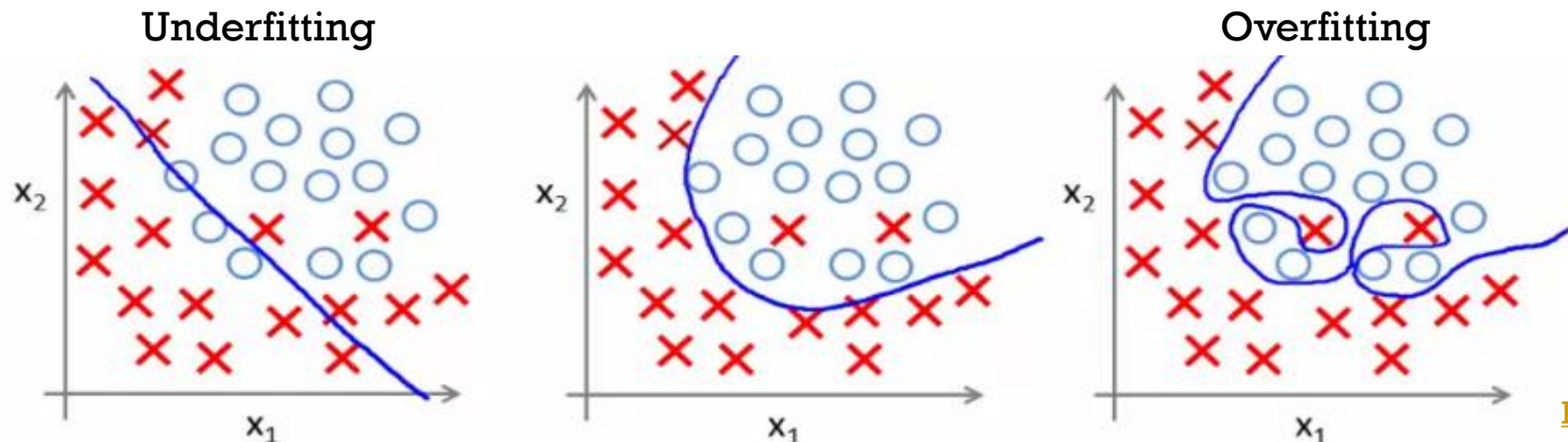
BIAS-VARIANCE TRADEOFF

- What if your model **bias** is too high?
 - Your model is **underfitting** – it is incapable of capturing the important characteristics of the training data
- What if your model **variance** is too high?
 - Your model is **overfitting** – it is fitting noise and unimportant characteristics of the data



BIAS-VARIANCE TRADEOFF

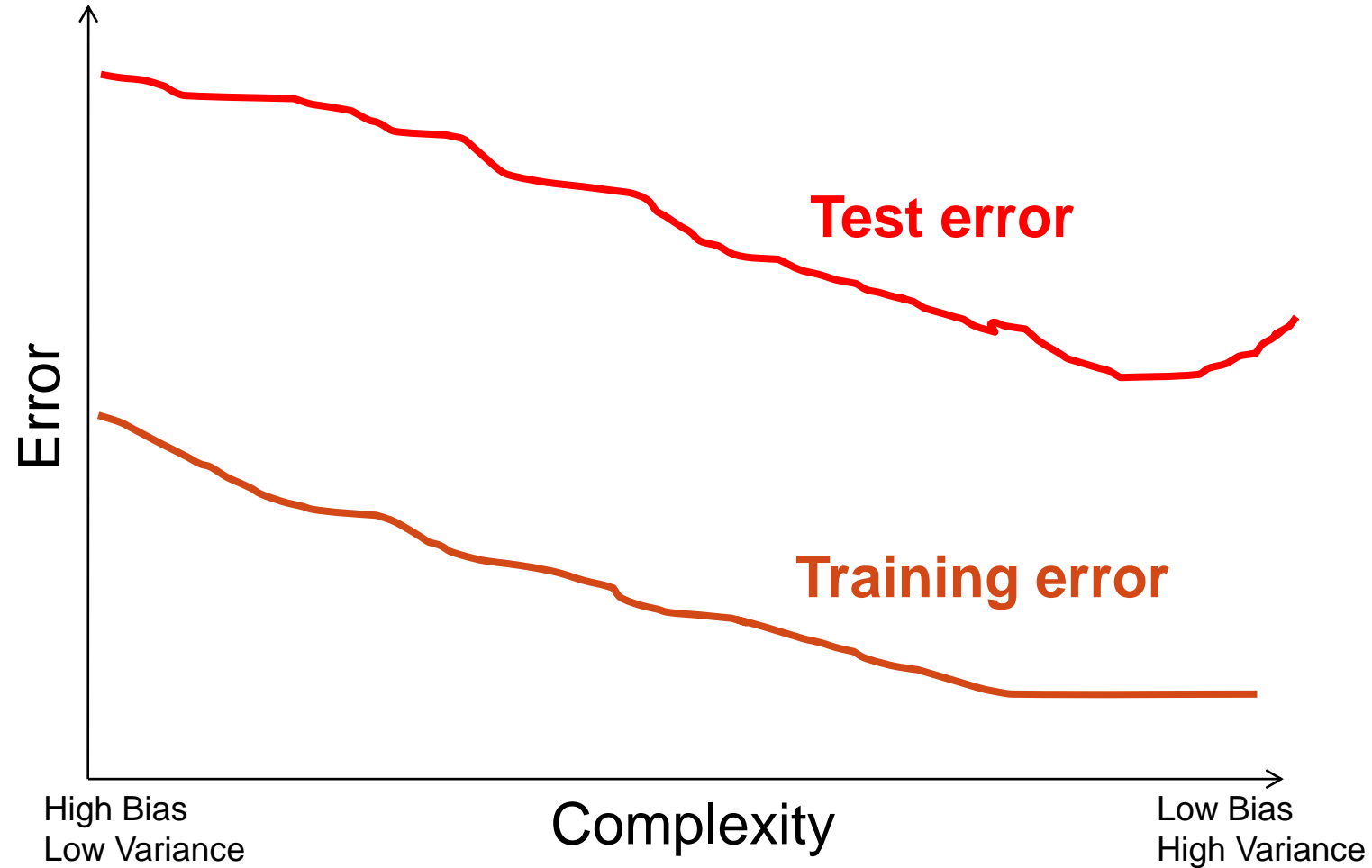
- What if your model **bias** is too high?
 - Your model is **underfitting** – it is incapable of capturing the important characteristics of the training data
- What if your model **variance** is too high?
 - Your model is **overfitting** – it is fitting noise and unimportant characteristics of the data
- How to recognize underfitting or overfitting?



BIAS-VARIANCE TRADEOFF

- What if your model **bias** is too high?
 - Your model is **underfitting** – it is incapable of capturing the important characteristics of the training data
- What if your model **variance** is too high?
 - Your model is **overfitting** – it is fitting noise and unimportant characteristics of the data
- How to recognize underfitting or overfitting?
 - Need to look at both training and test error
 - **Underfitting:** training and test error are both *high*
 - **Overfitting:** training error is *low*, test error is *high*

LOOKING AT TRAINING AND TEST ERROR



OPTIMIZATION

Optimization is the process of finding the set of parameters W that minimize the loss function.

Strategy #1: First very bad idea solution: Random search:

Simply try out many different random weights and keep track of what works best.

OPTIMIZATION

Optimization is the process of finding the set of parameters W that minimize the loss function.

Strategy #1: First very bad idea solution: Random search:

Simply try out many different random weights and keep track of what works best.

Strategy #2: Random local search:

Start out with a random W , generate random changes δW to it and if the loss at the changed $W + \delta W$ is lower, we will perform an update.

OPTIMIZATION

Optimization is the process of finding the set of parameters W that minimize the loss function.

Strategy #1: First very bad idea solution: Random search:

Simply try out many different random weights and keep track of what works best.

Strategy #2: Random local search:

Start out with a random W , generate random changes δW to it and if the loss at the changed $W + \delta W$ is lower, we will perform an update.

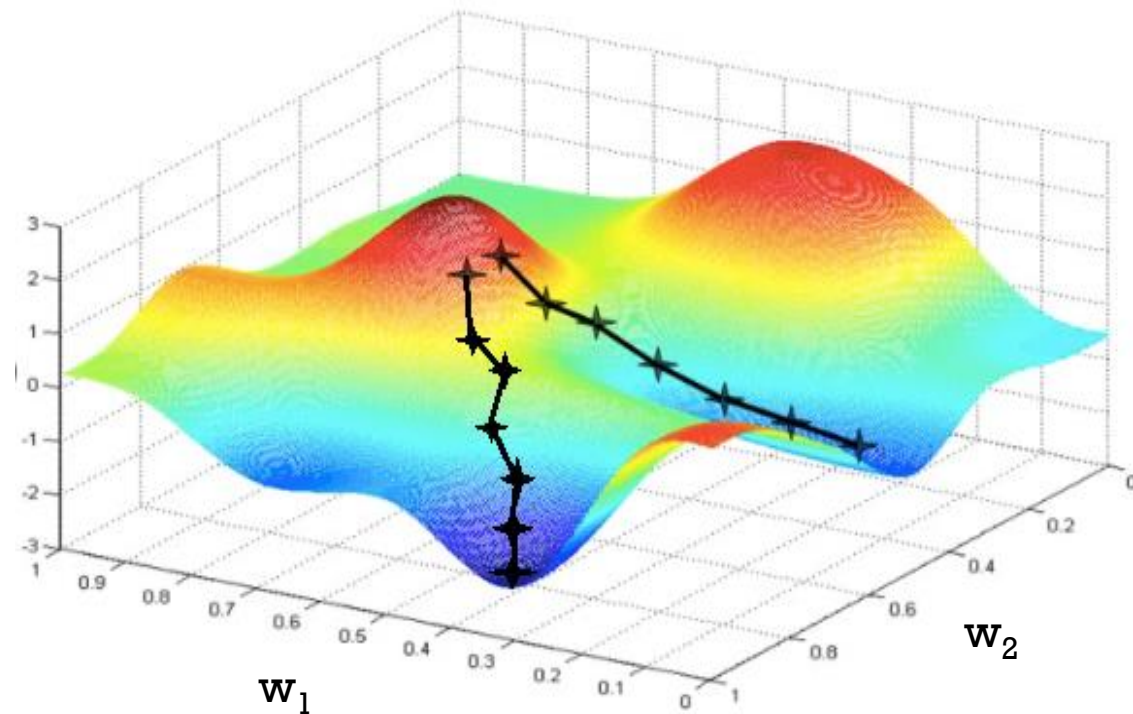
Strategy #3: Following the gradients:

There is no need to randomly search for a good direction: this direction is related to the gradient of the loss function.

GRADIENT DESCENT

- Goal: find w to minimize loss $L(w)$
- Start with some initial estimate of w
- At each step, find $\nabla L(w)$, the *gradient* of the loss w.r.t. w , and take a small step in the *opposite* direction

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



GRADIENT DESCENT

The procedure of repeatedly evaluating the **gradient of loss function** and then performing a **parameter update**.

Vanilla (Original) Gradient Descent:

```
while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

GRADIENT DESCENT

The procedure of repeatedly evaluating the **gradient of loss function** and then performing a **parameter update**.

Vanilla (Original) Gradient Descent:

```
while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

Mini-batch Gradient Descent (MGD):

```
while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```


GRADIENT DESCENT

The procedure of repeatedly evaluating the **gradient of loss function** and then performing a **parameter update**.

Vanilla (Original) Gradient Descent:

```
while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

Mini-batch Gradient Descent (MGD):

```
while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

Stochastic Gradient Descent (SGD):

Special case of MGD when mini-batch contains only a single example

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = \quad \frac{\partial f}{\partial y} =$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = \quad \frac{\partial f}{\partial y} =$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

$$f(x, y) = \max(x, y) \quad \rightarrow \quad \frac{\partial f}{\partial x} = \quad \frac{\partial f}{\partial y} =$$

INTERPRETATION OF THE GRADIENT

Interpretation. Derivatives indicate the rate of change of a function with respect to that variable surrounding an infinitesimally small region near a particular point:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

$$f(x, y) = \max(x, y) \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1(x \geq y) \quad \frac{\partial f}{\partial y} = 1(y \geq x)$$

COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

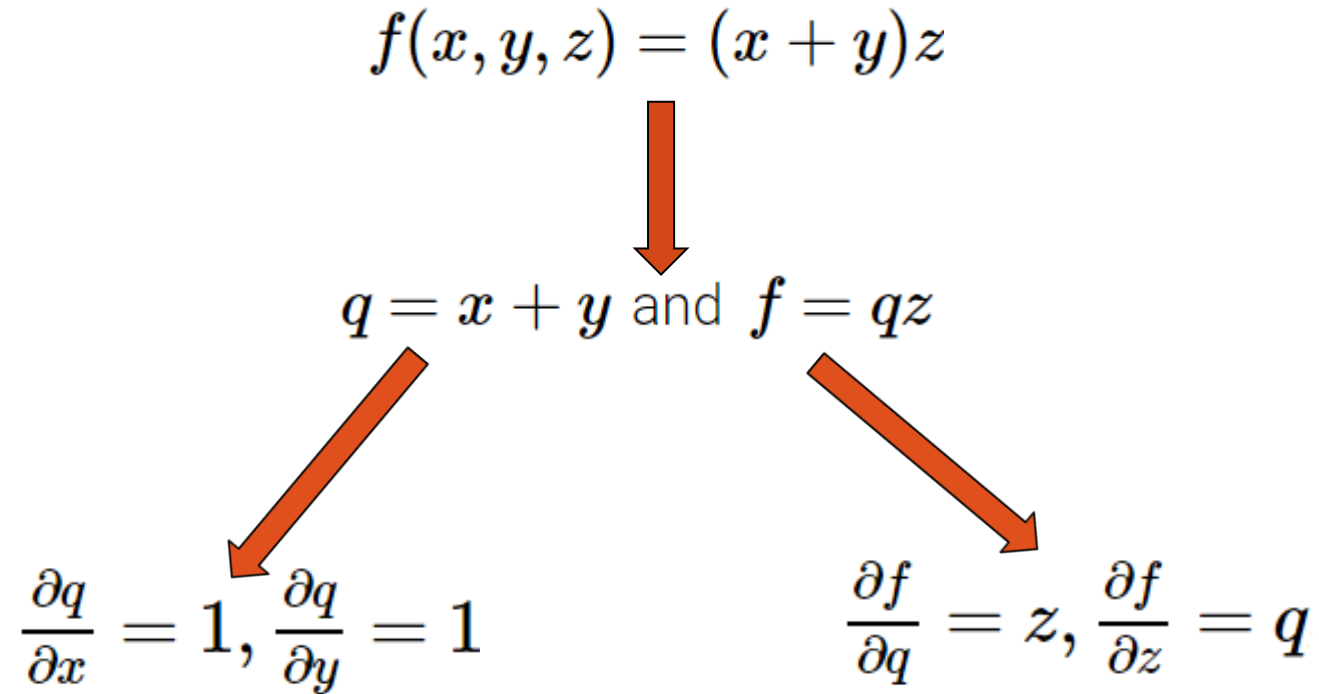
COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

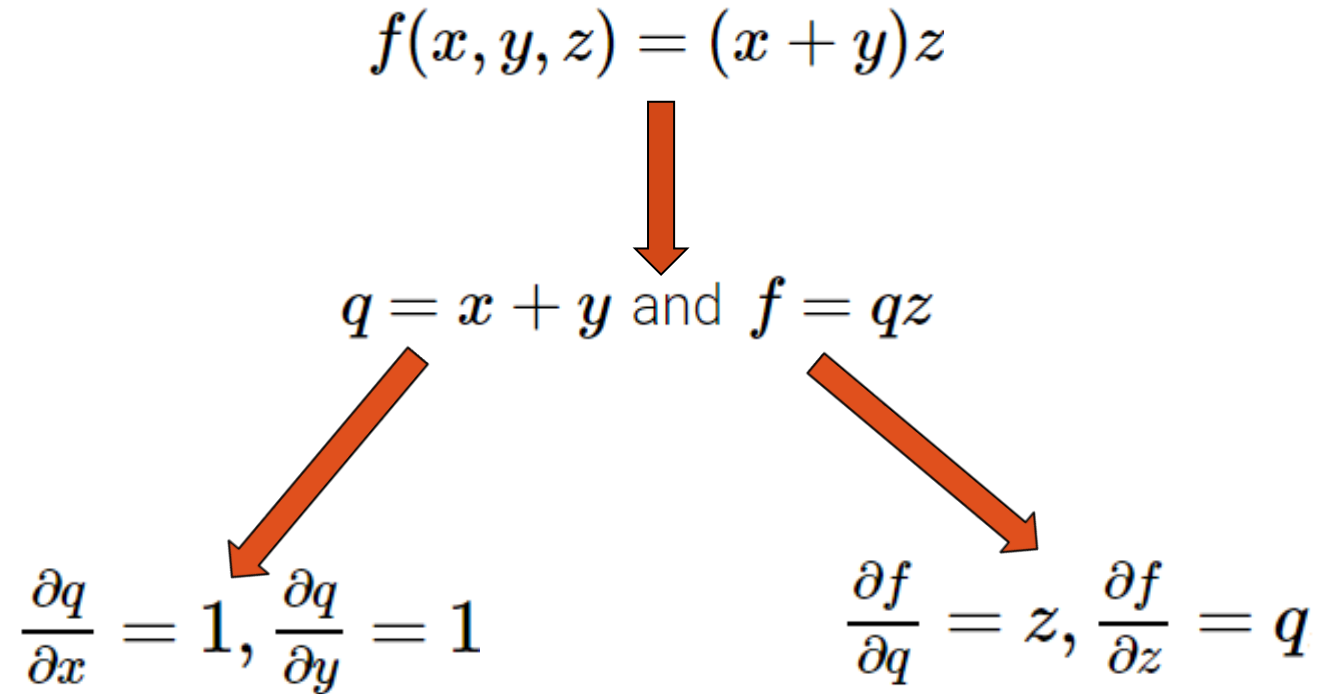


$$q = x + y \text{ and } f = qz$$

COMPOUND EXPRESSIONS WITH CHAIN RULE




COMPOUND EXPRESSIONS WITH CHAIN RULE



Chain rule: $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$

COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

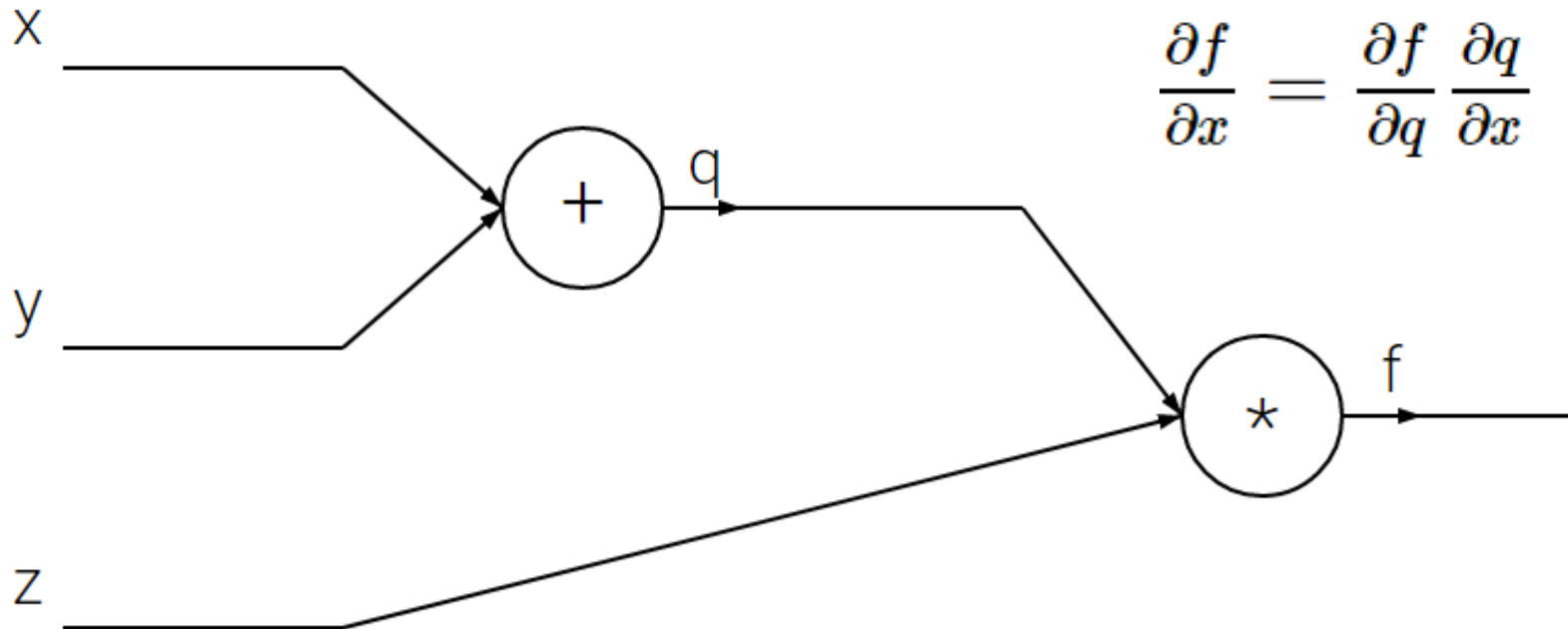


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

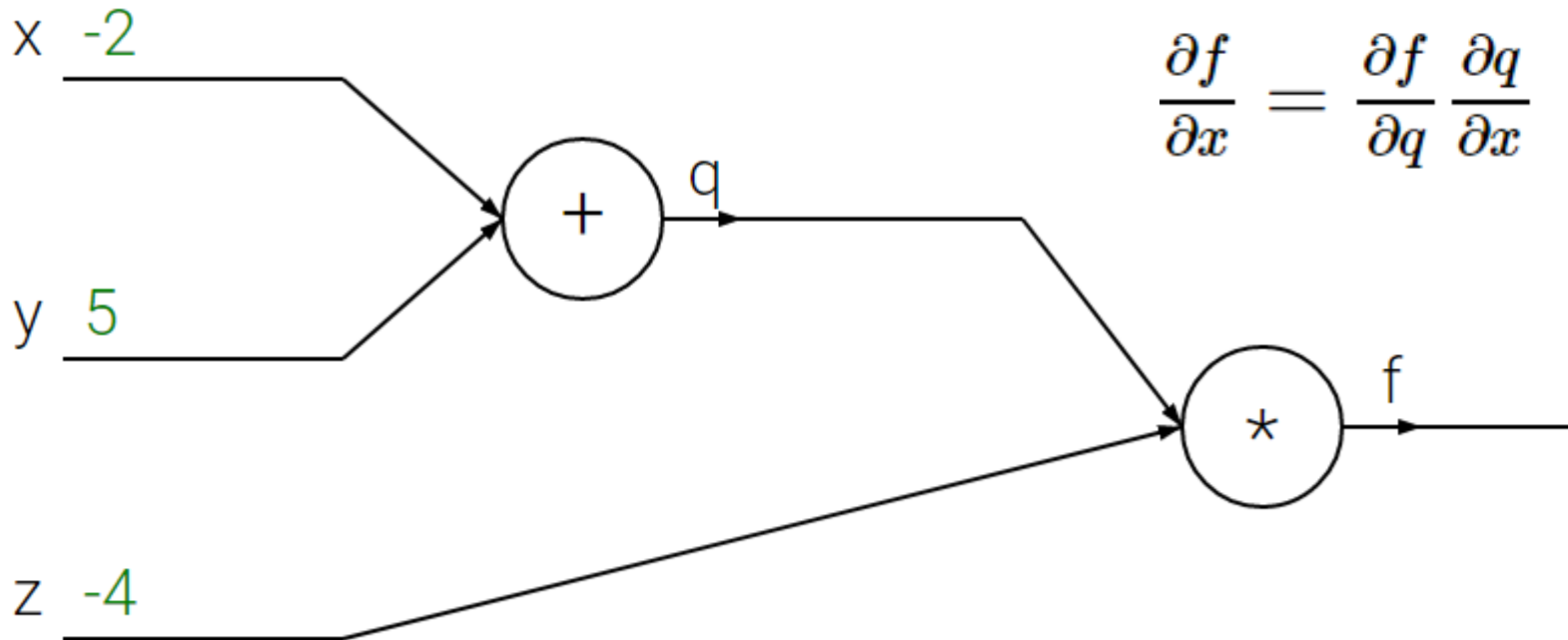


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

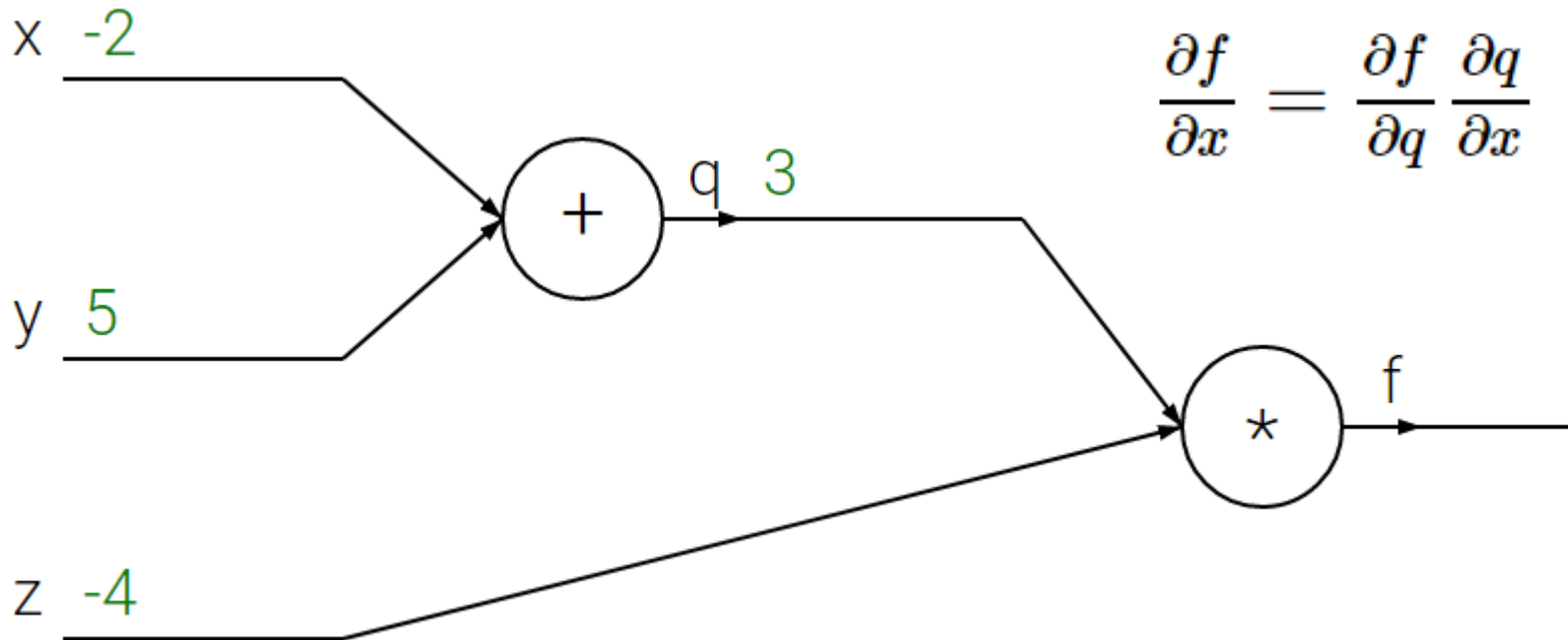


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

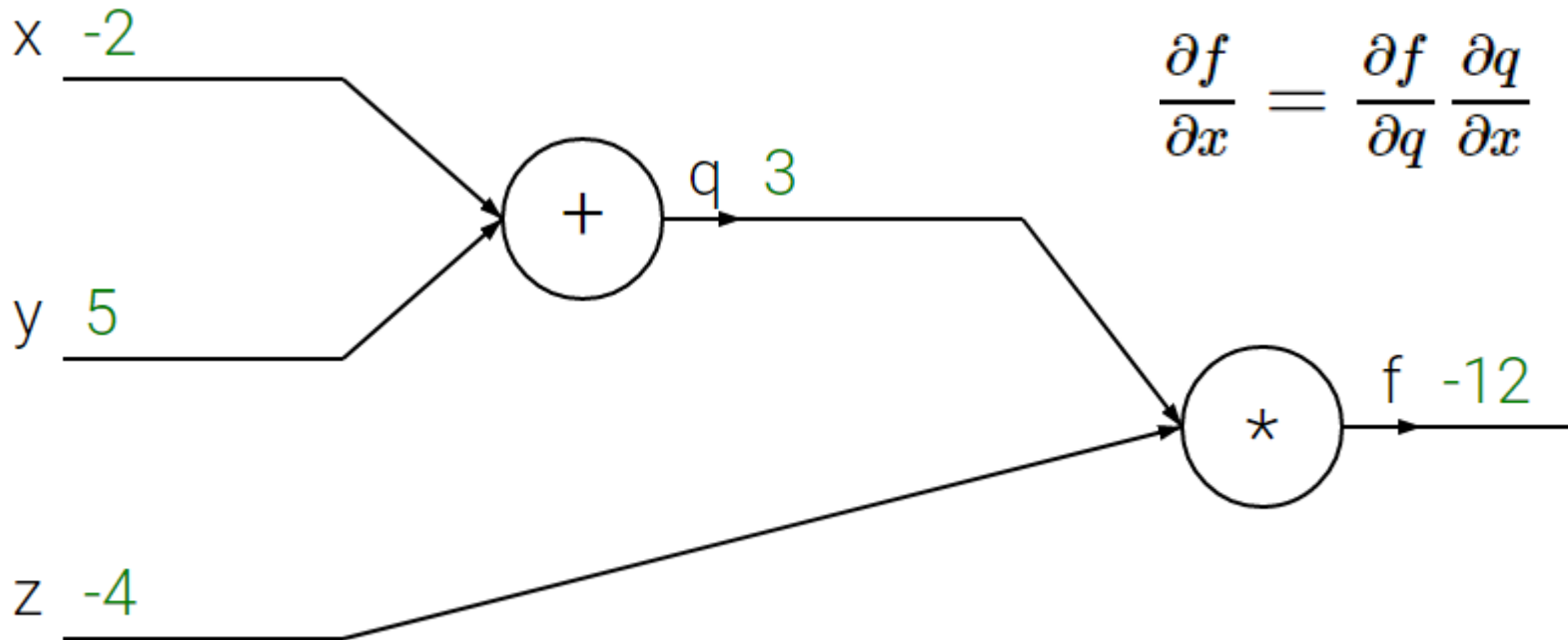


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

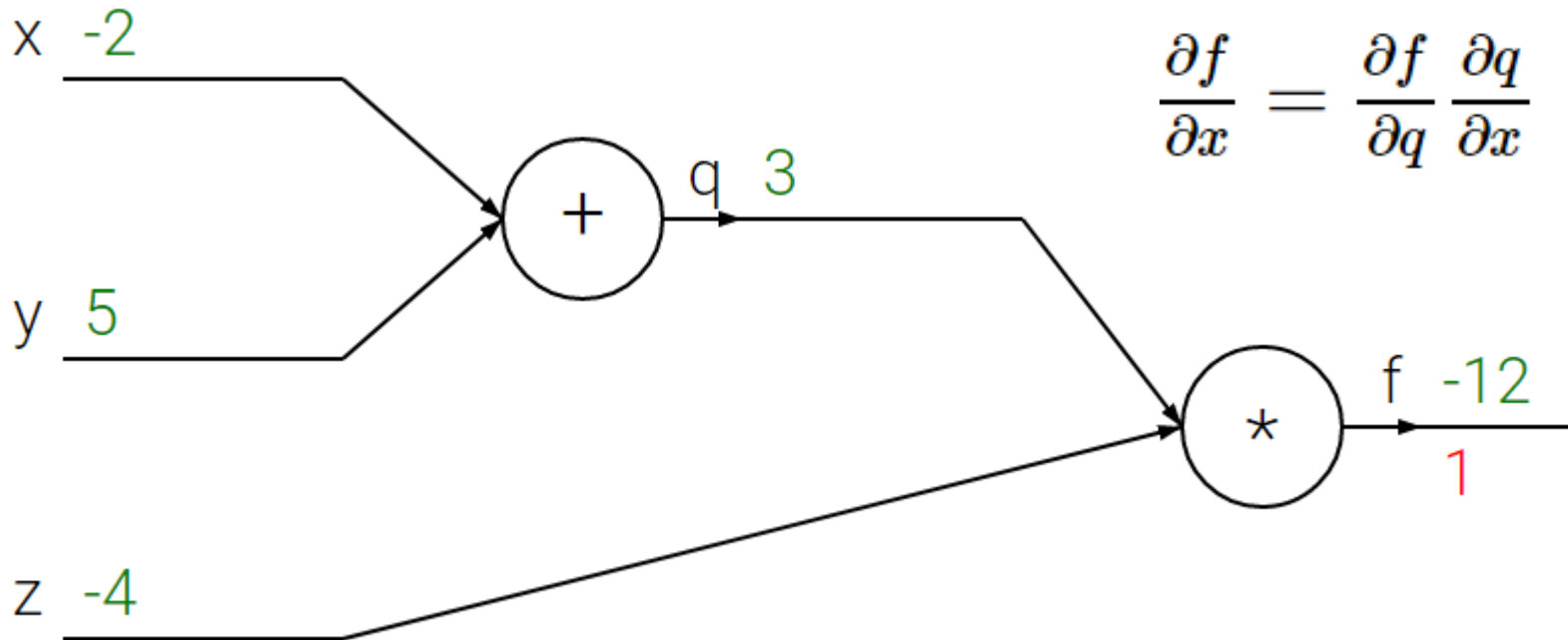


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

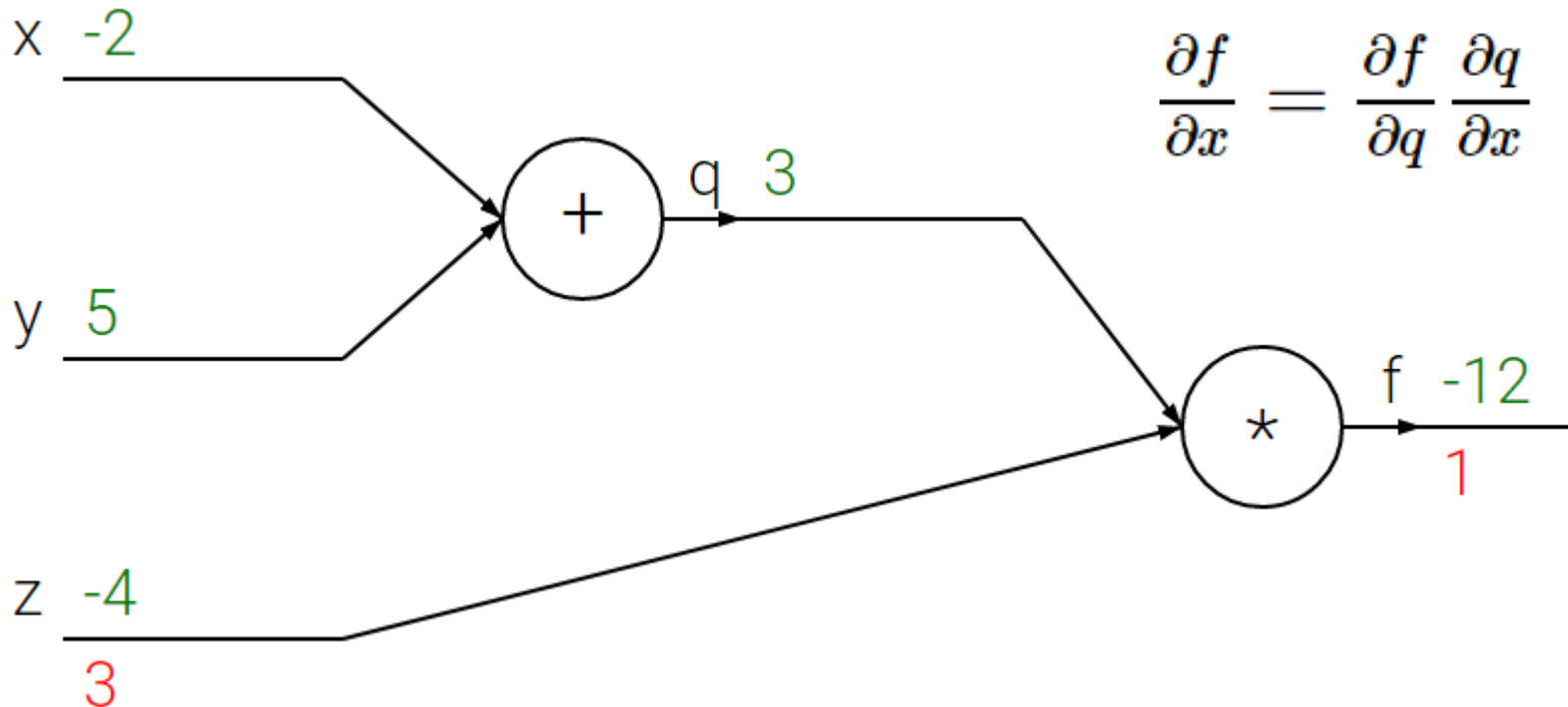


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

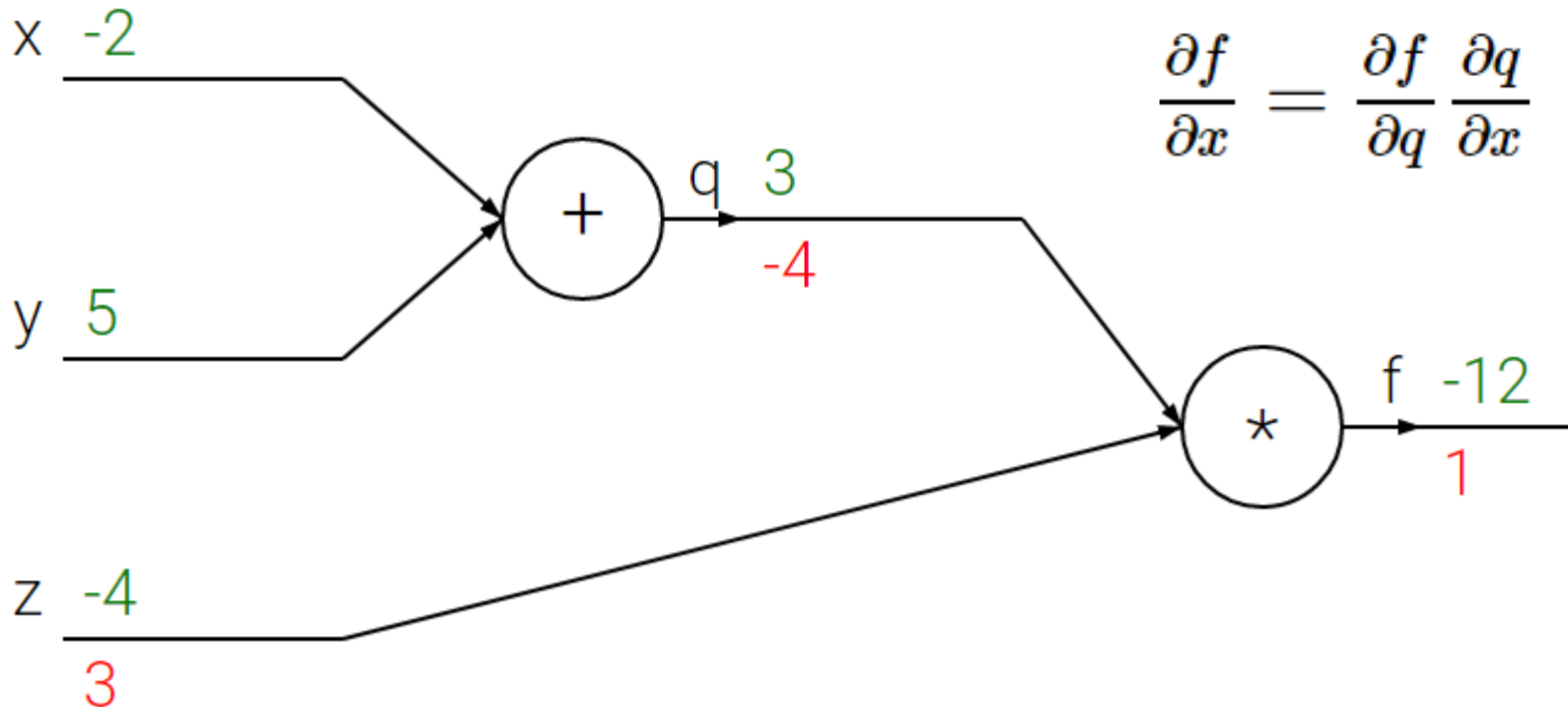


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

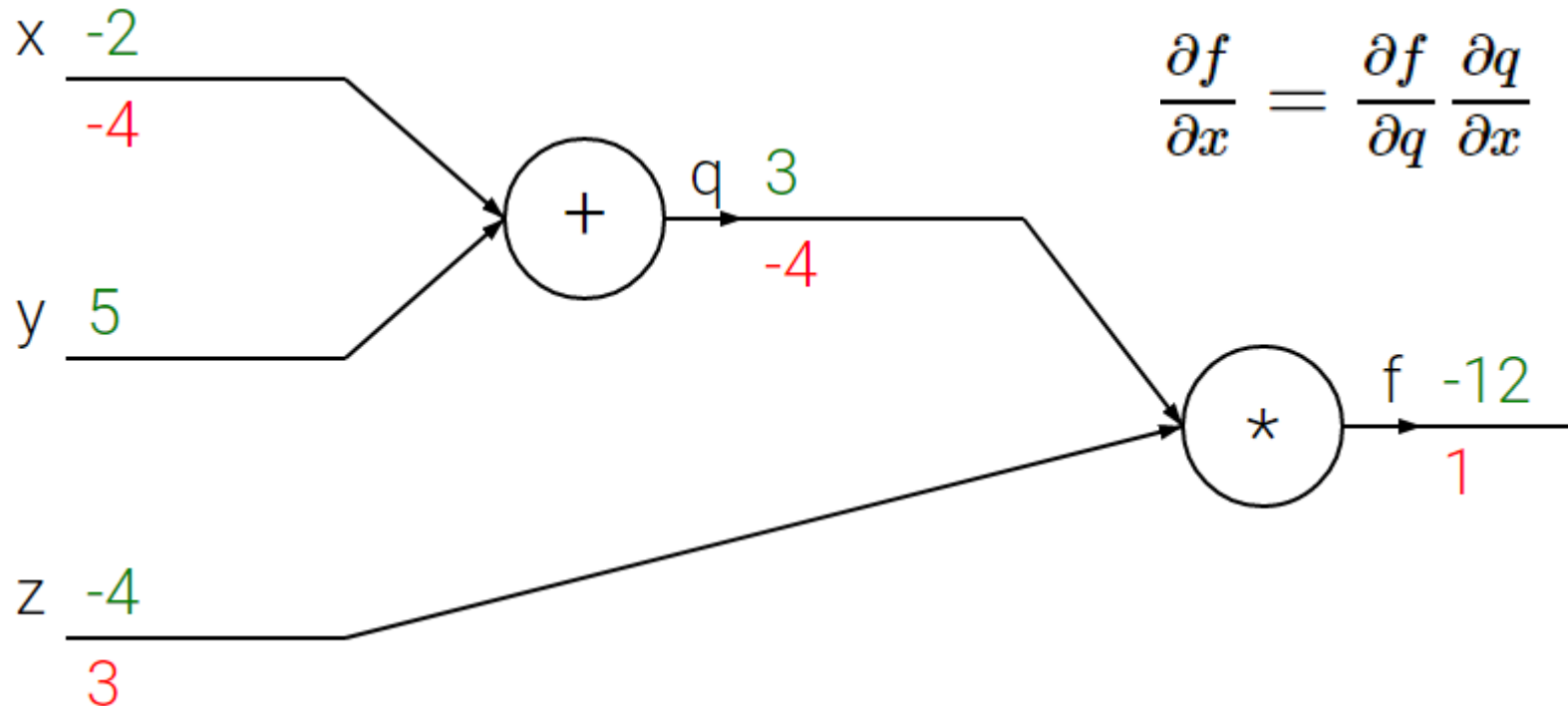


$$q = x + y \text{ and } f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



COMPOUND EXPRESSIONS WITH CHAIN RULE

$$f(x, y, z) = (x + y)z$$

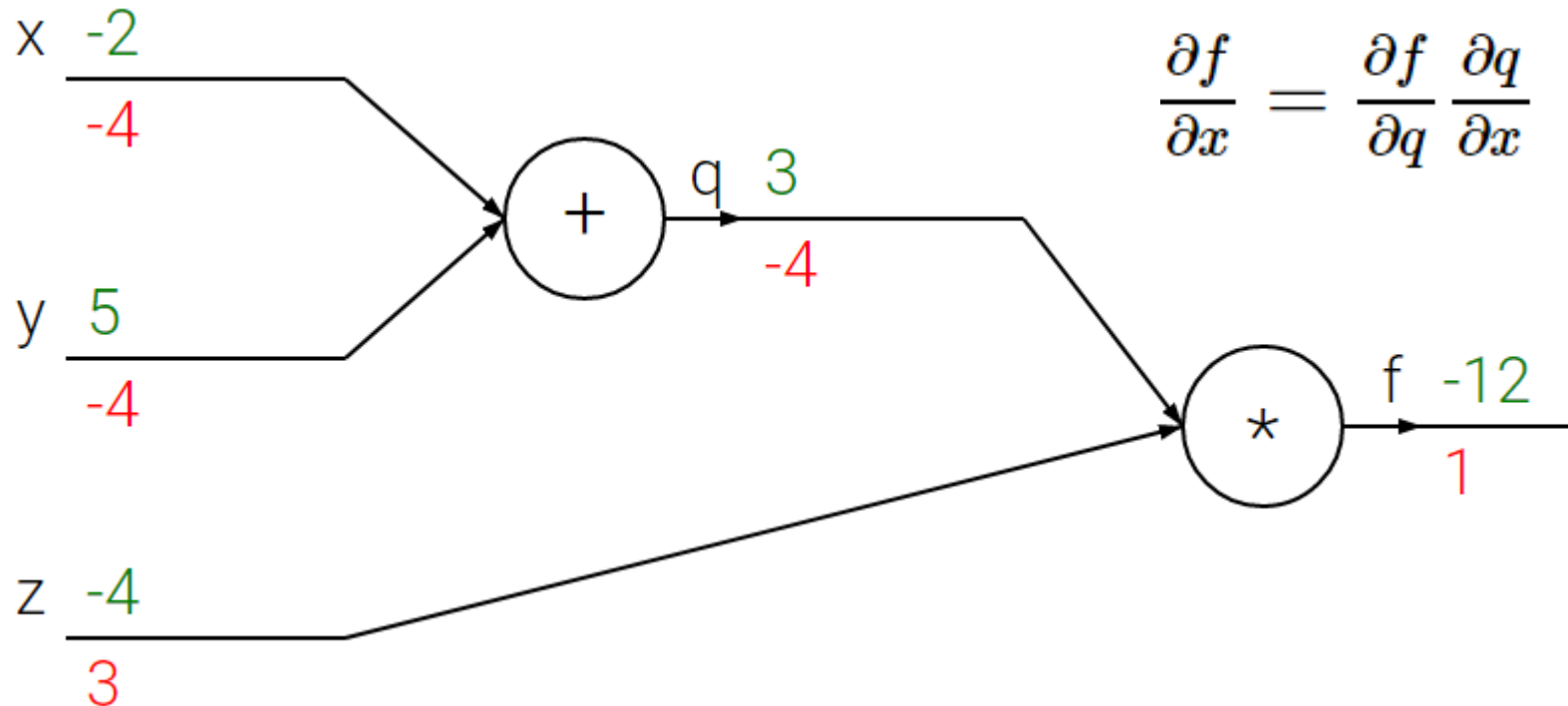


$$q = x + y \text{ and } f = qz$$

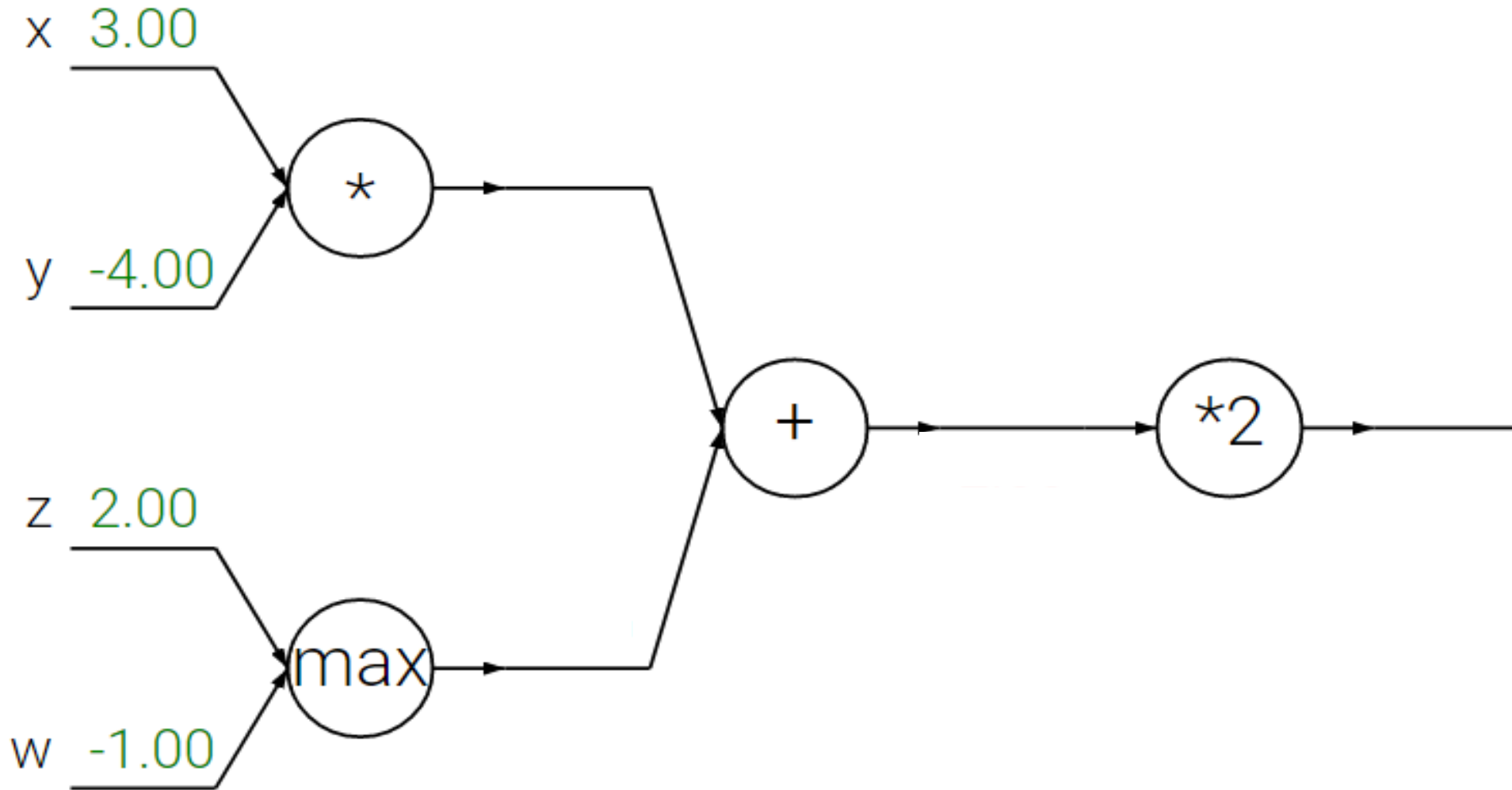
$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

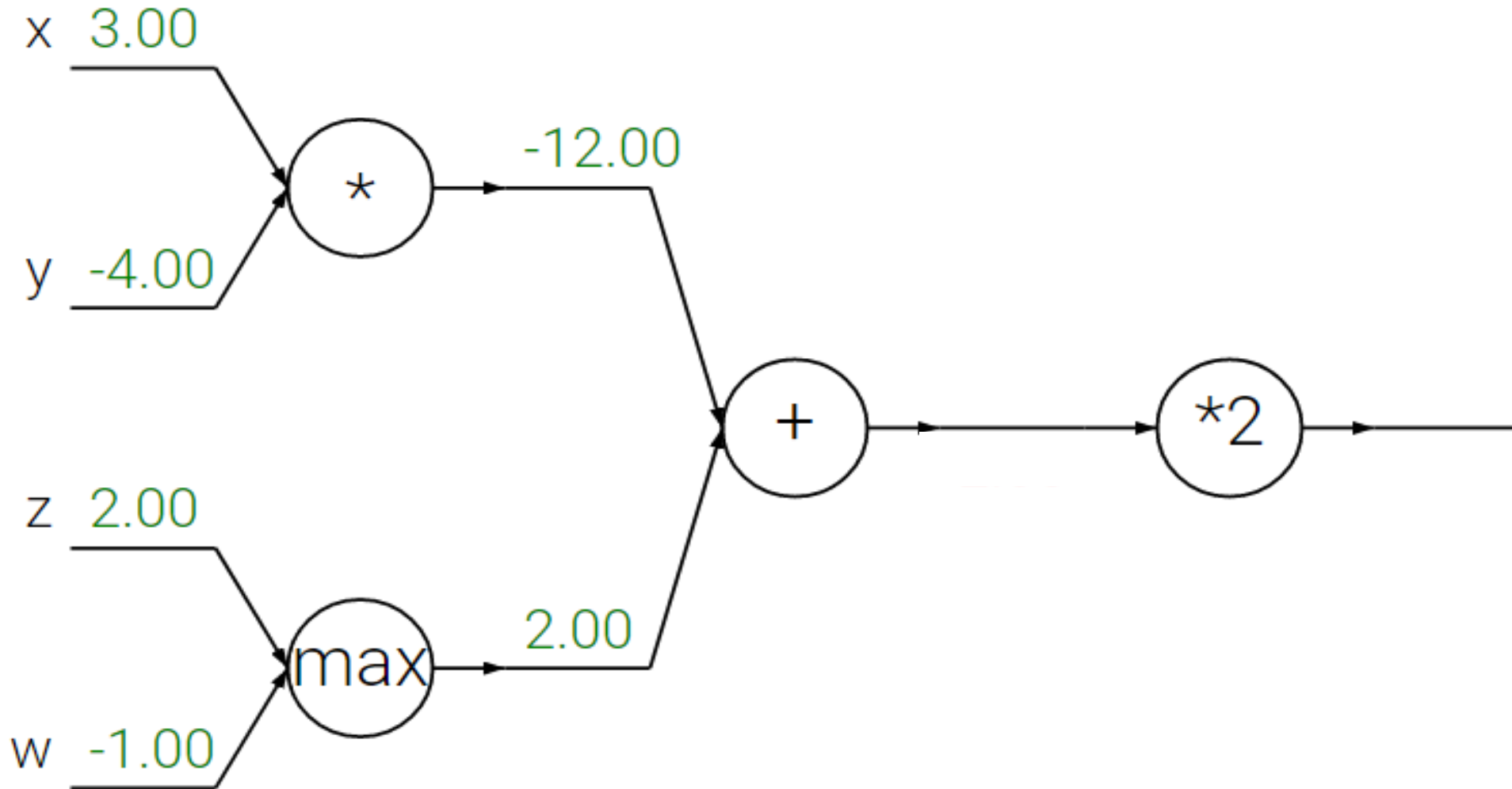
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



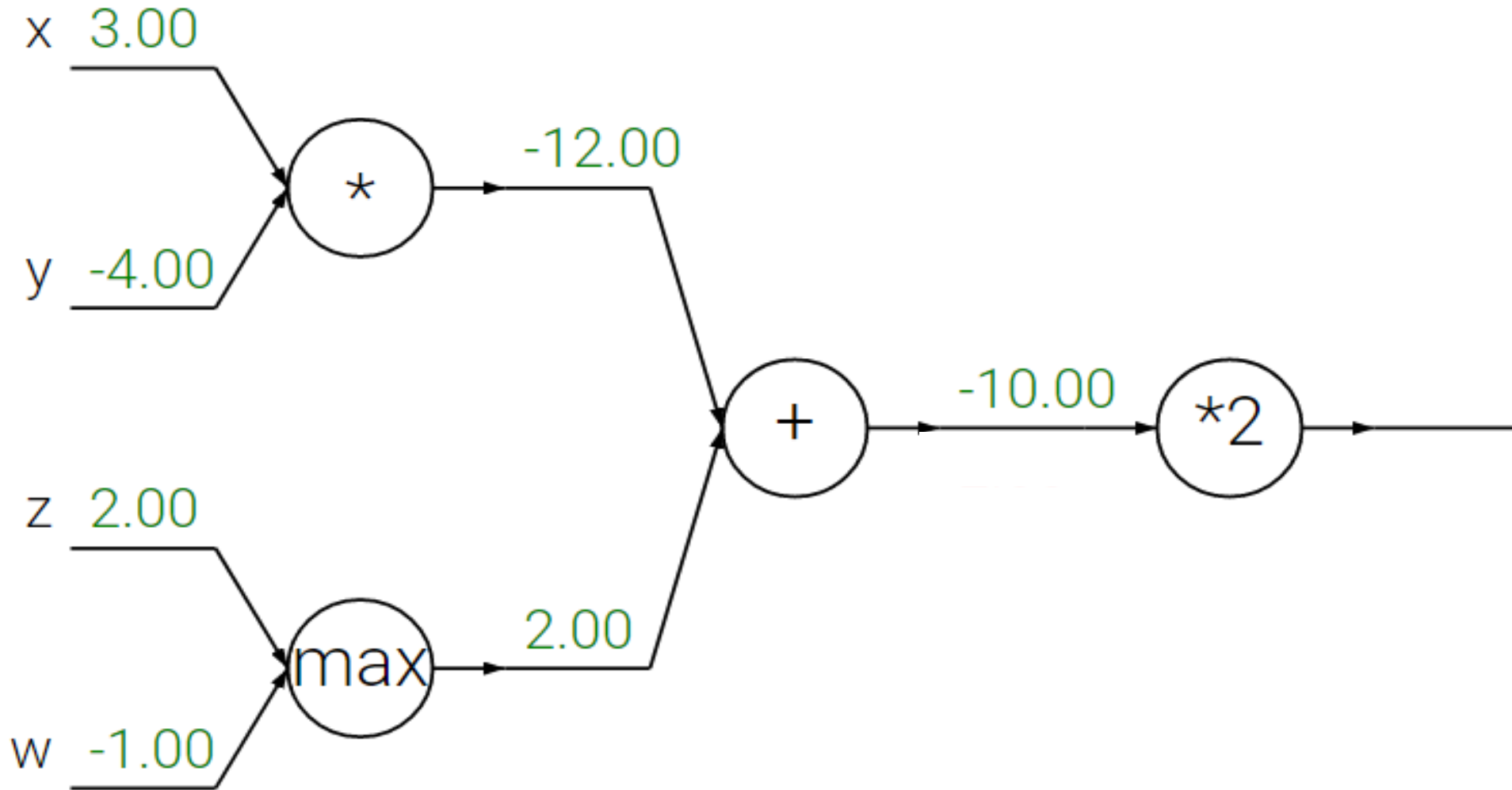
FORWARD AND BACKWARD PASS



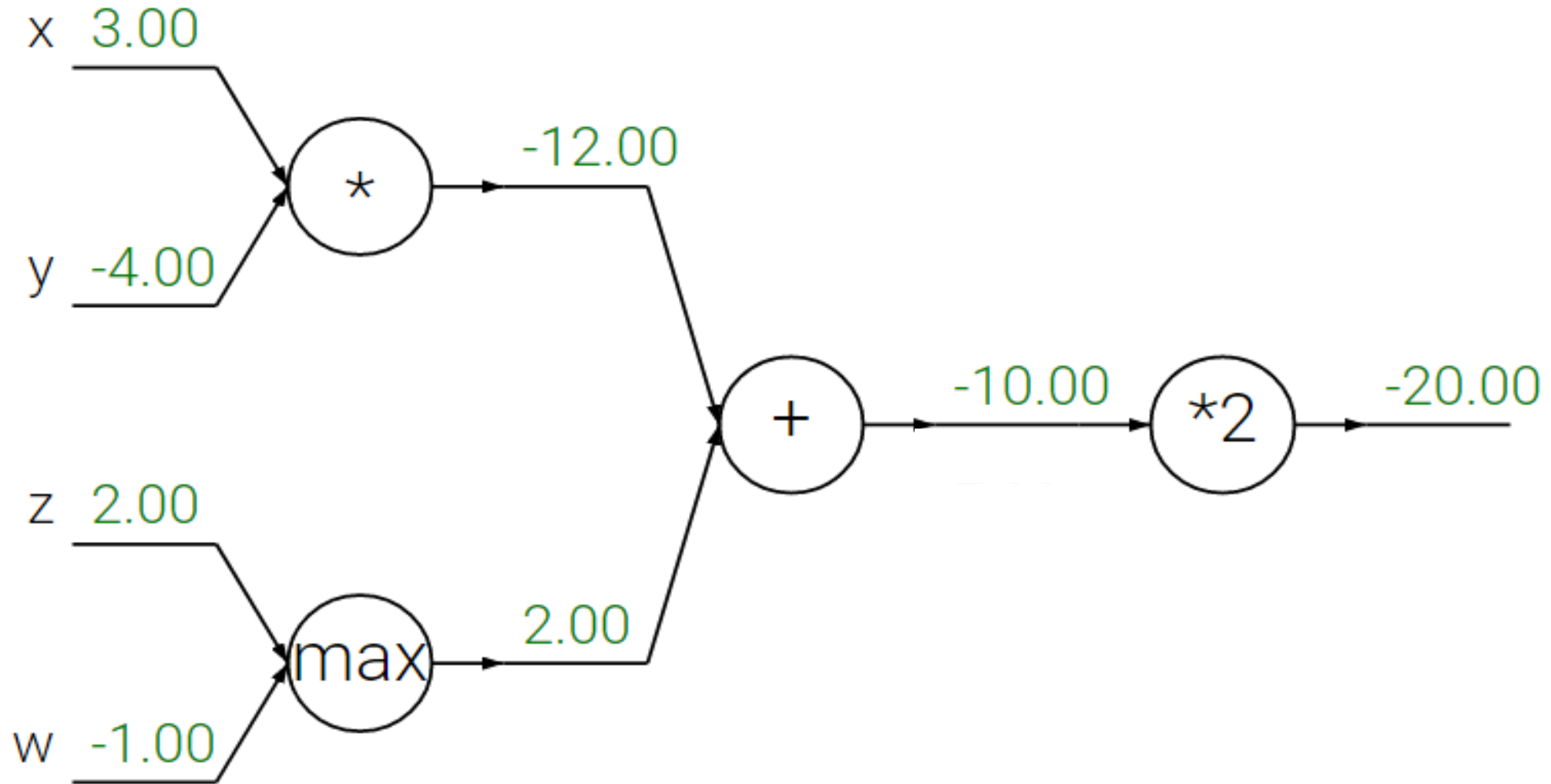
FORWARD AND BACKWARD PASS



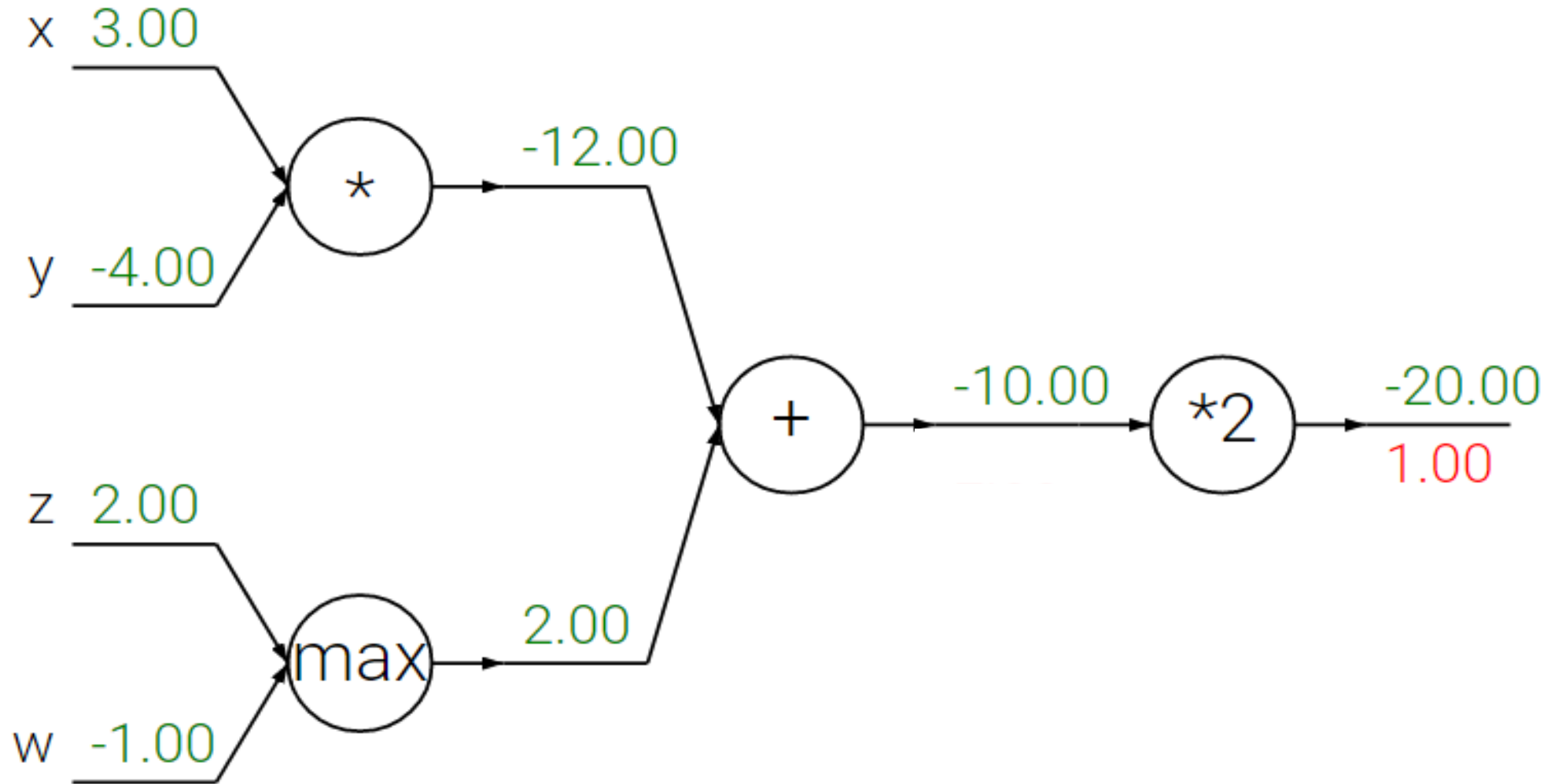
FORWARD AND BACKWARD PASS



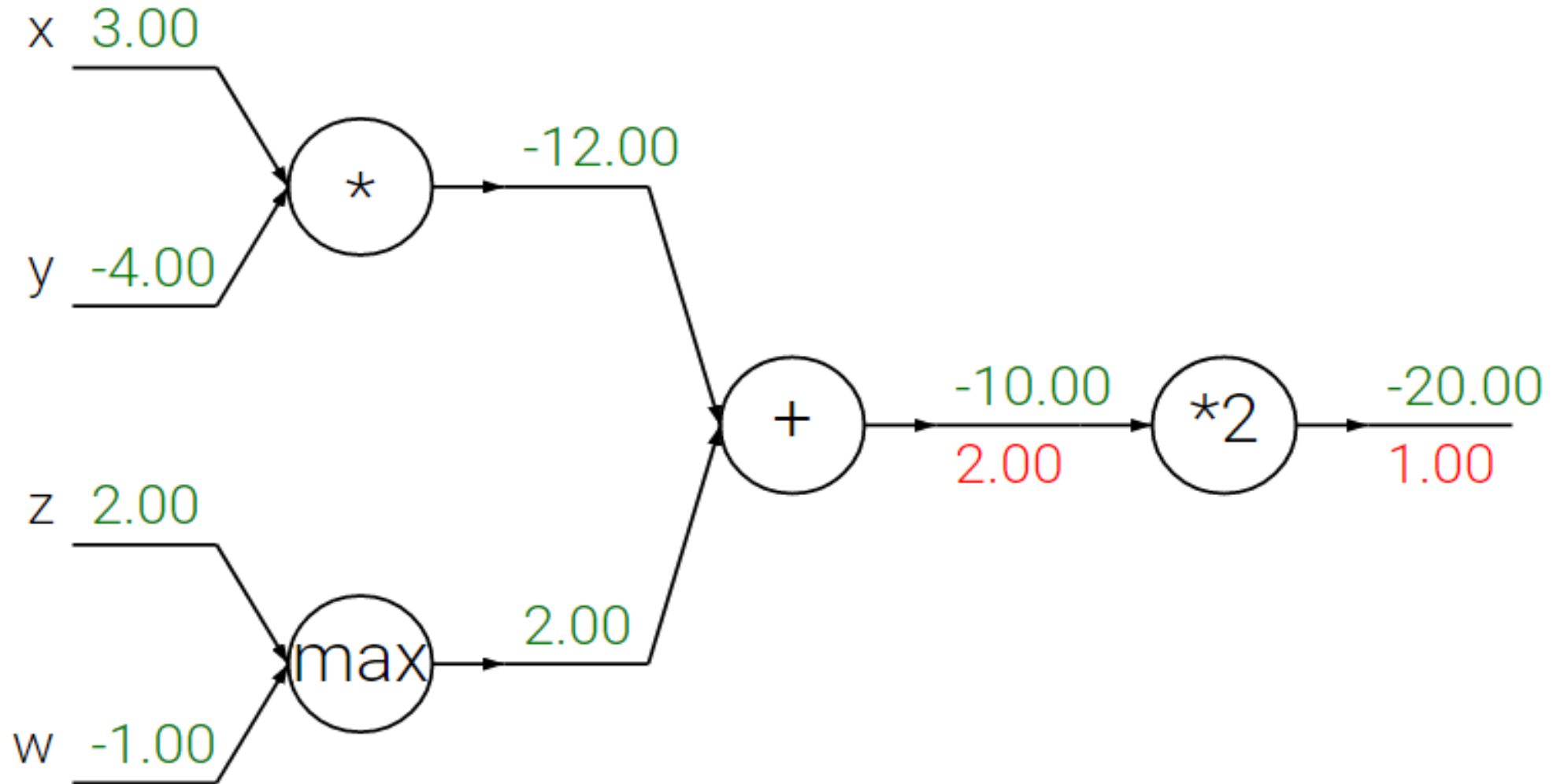
FORWARD AND BACKWARD PASS



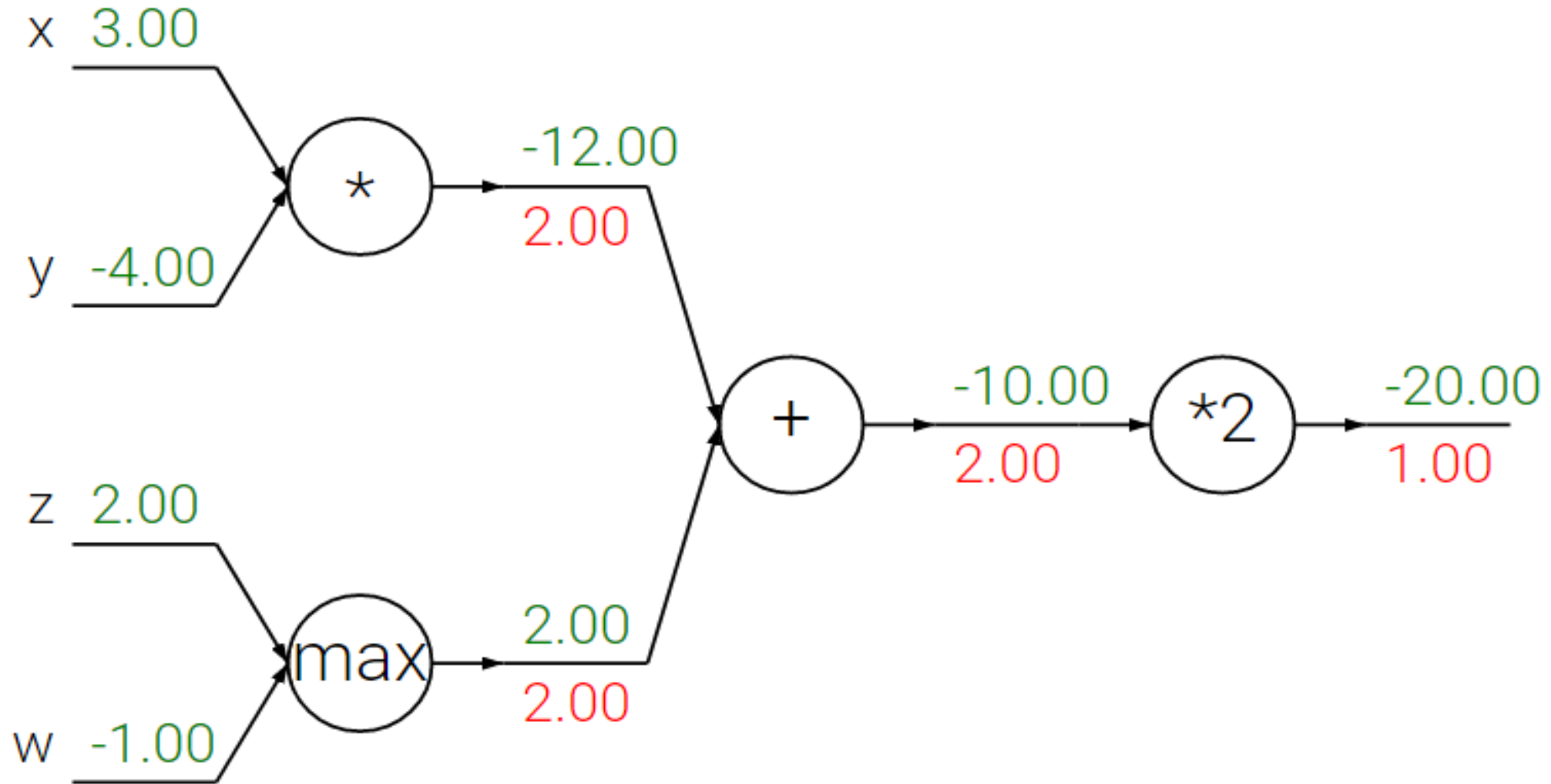
FORWARD AND BACKWARD PASS



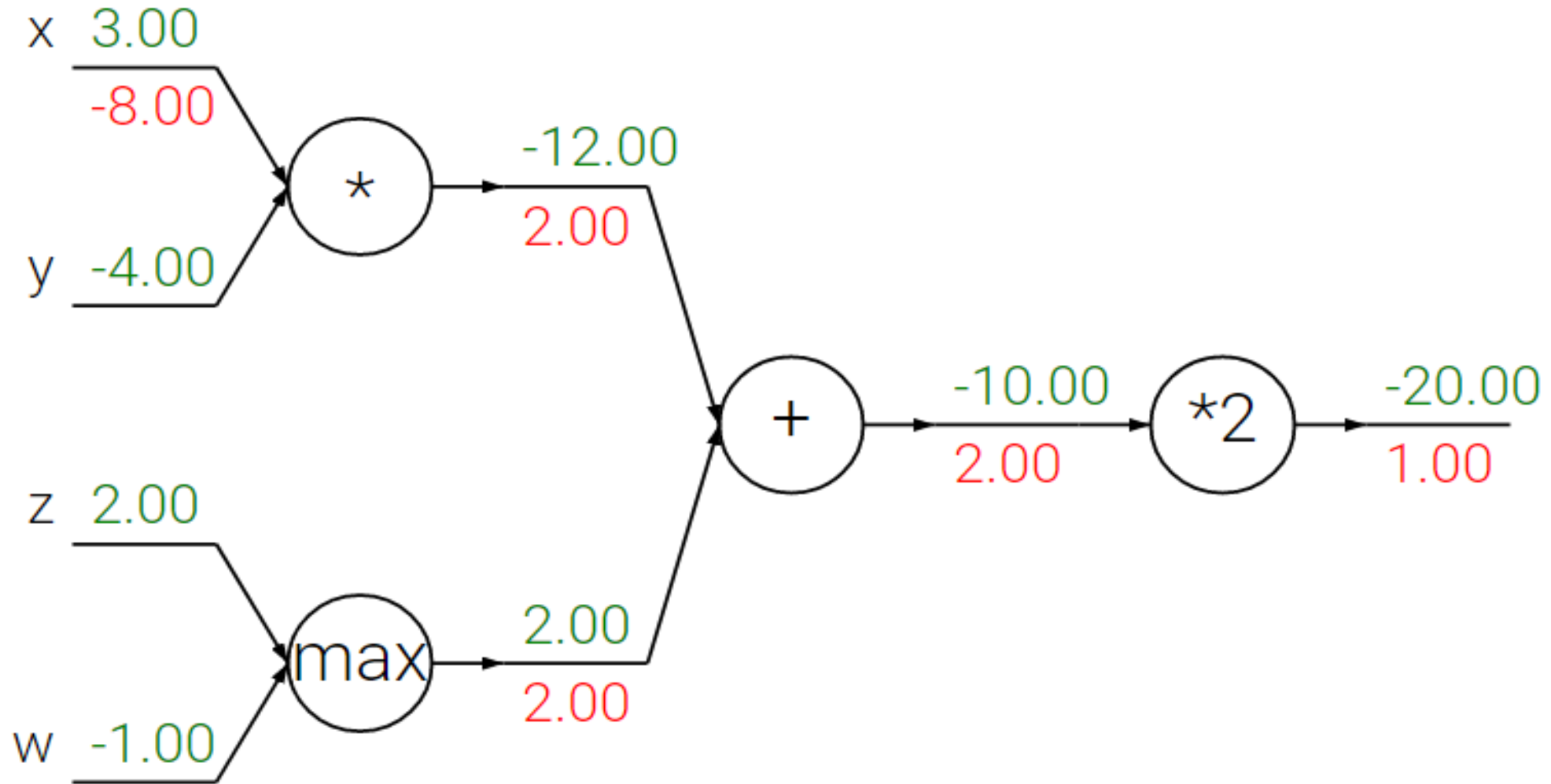
FORWARD AND BACKWARD PASS



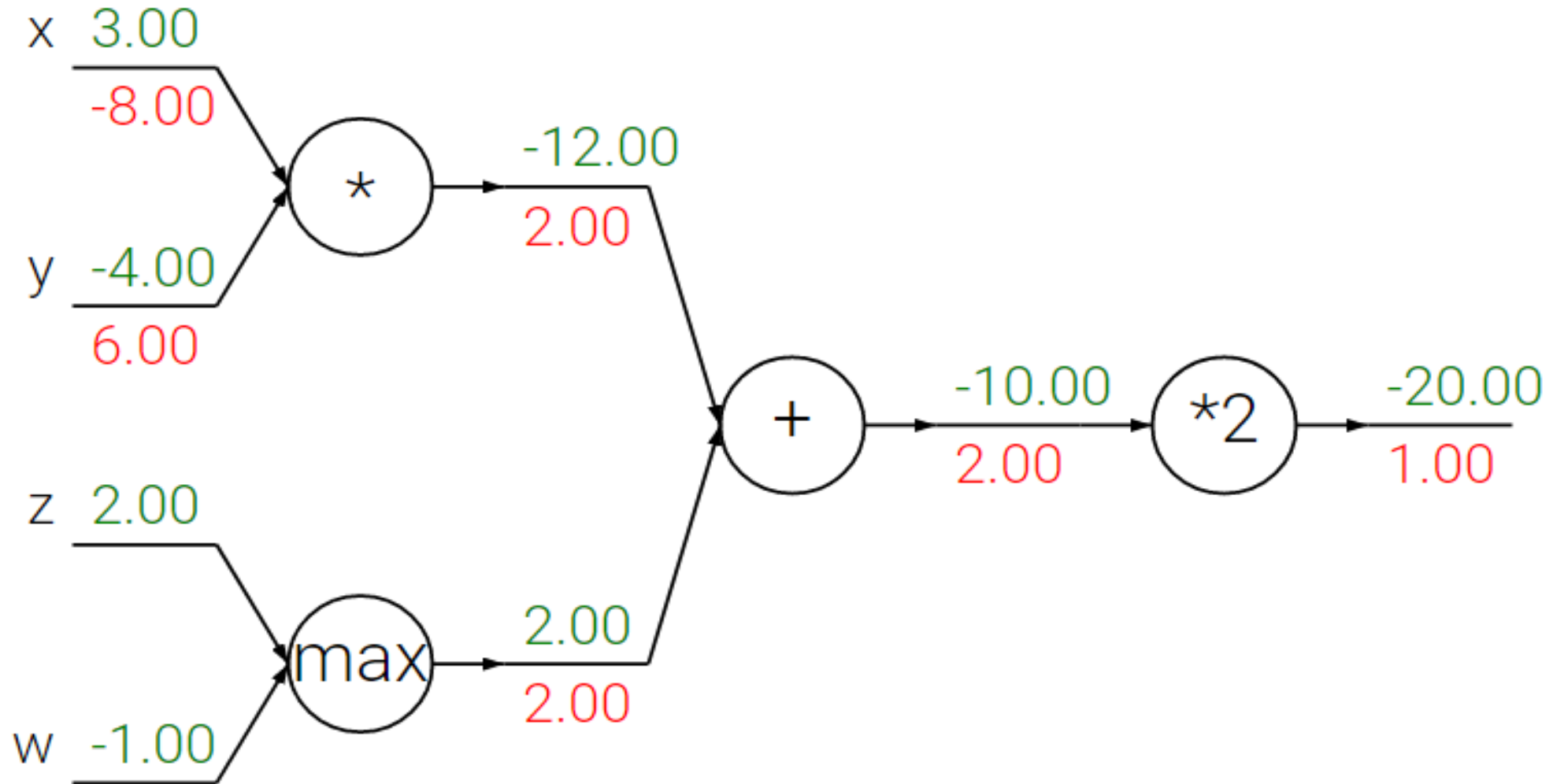
FORWARD AND BACKWARD PASS



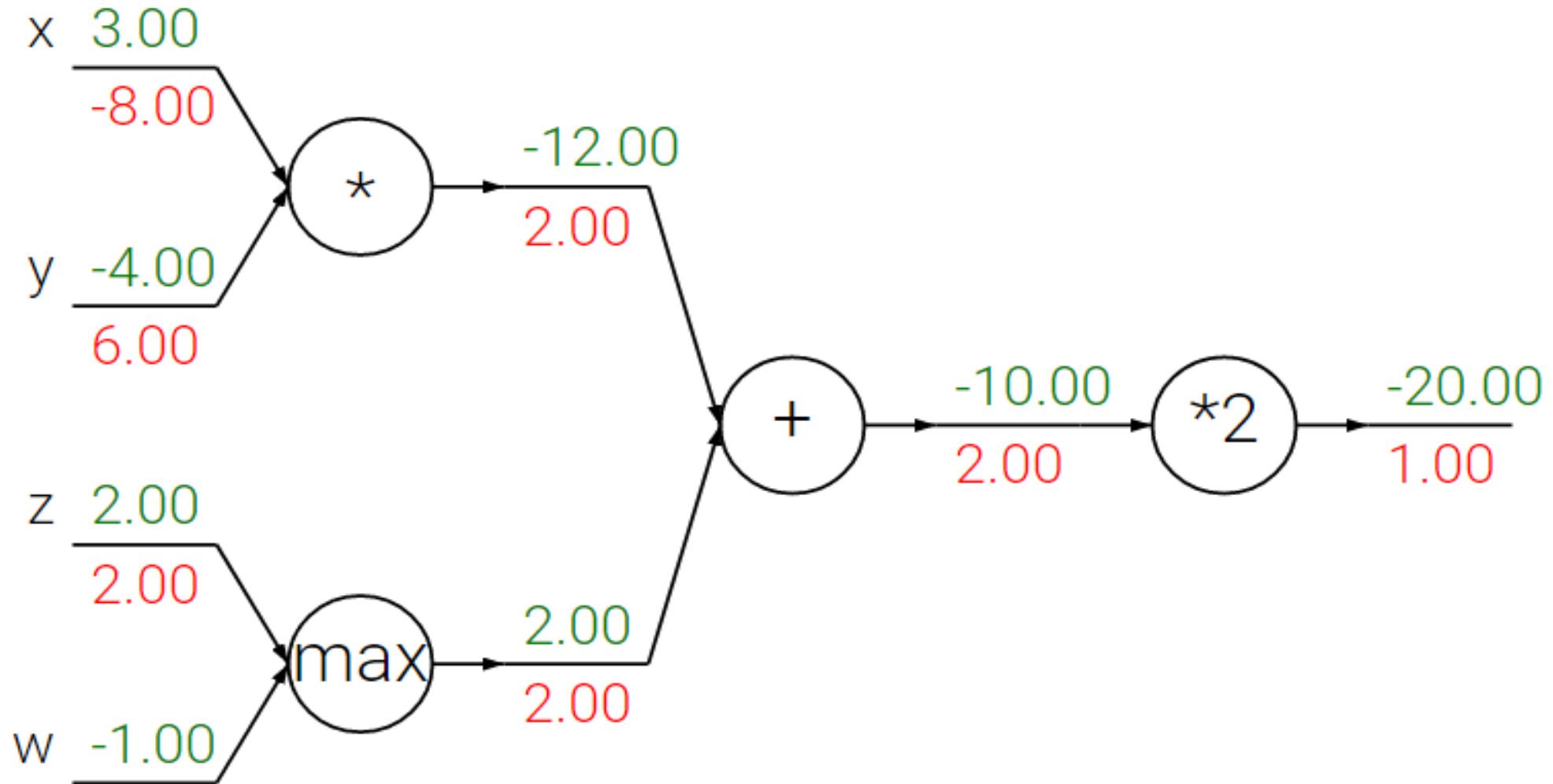
FORWARD AND BACKWARD PASS



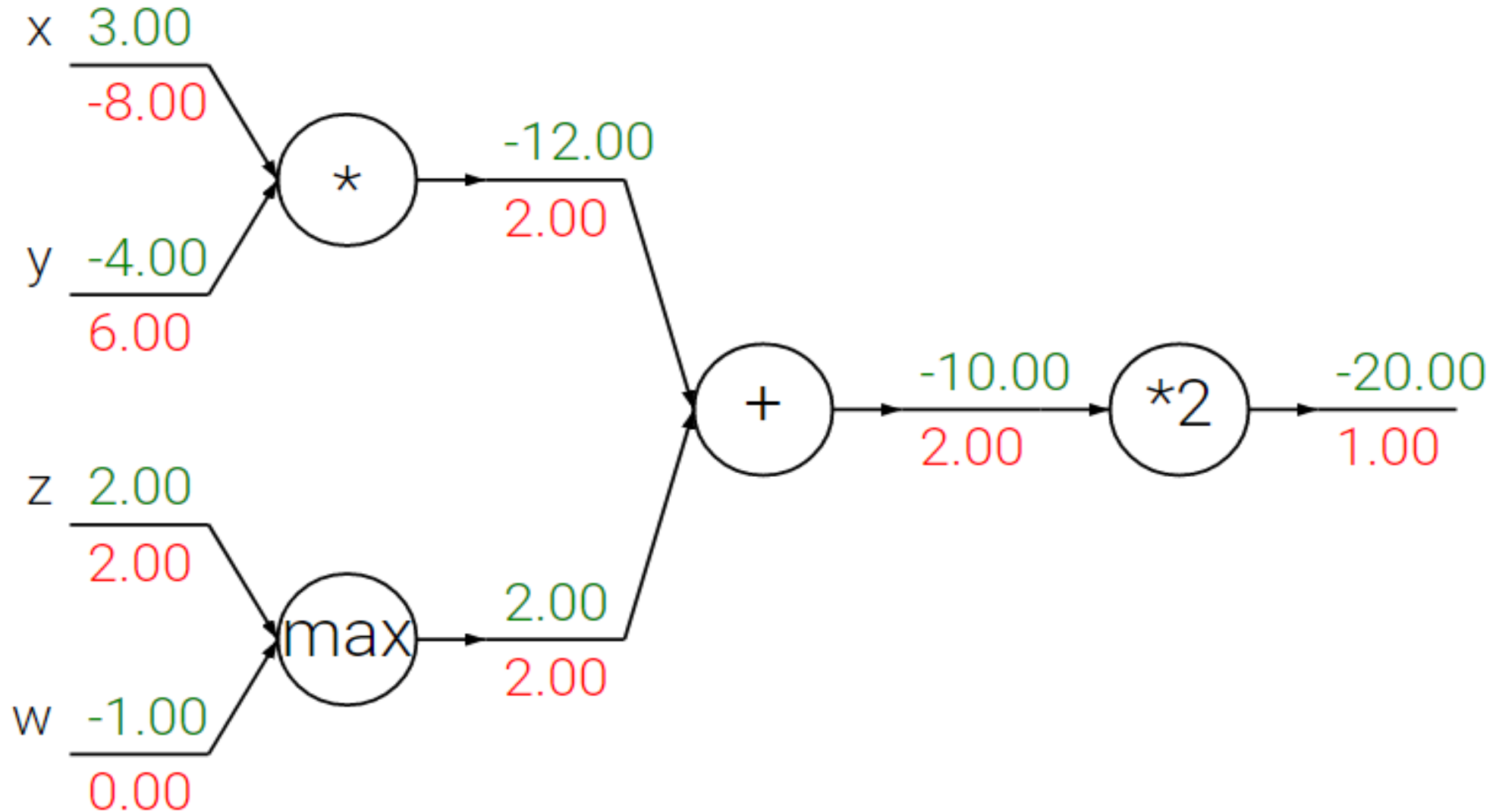
FORWARD AND BACKWARD PASS



FORWARD AND BACKWARD PASS

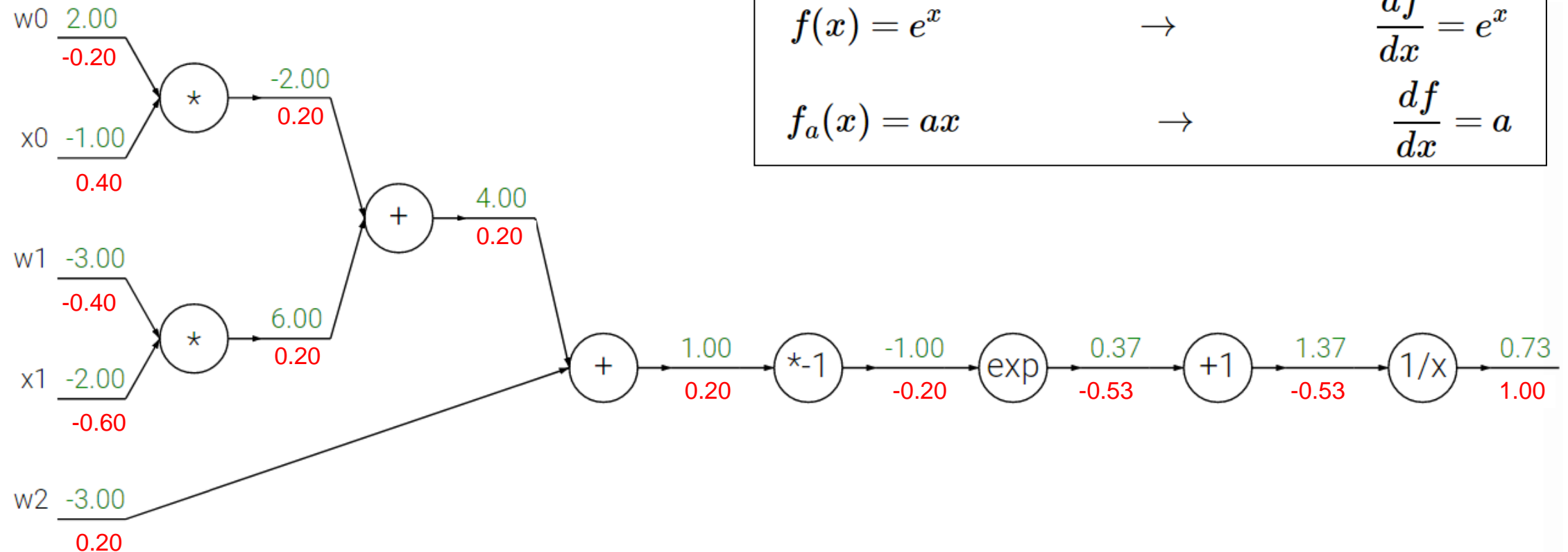


FORWARD AND BACKWARD PASS



SIGMOID EXAMPLE

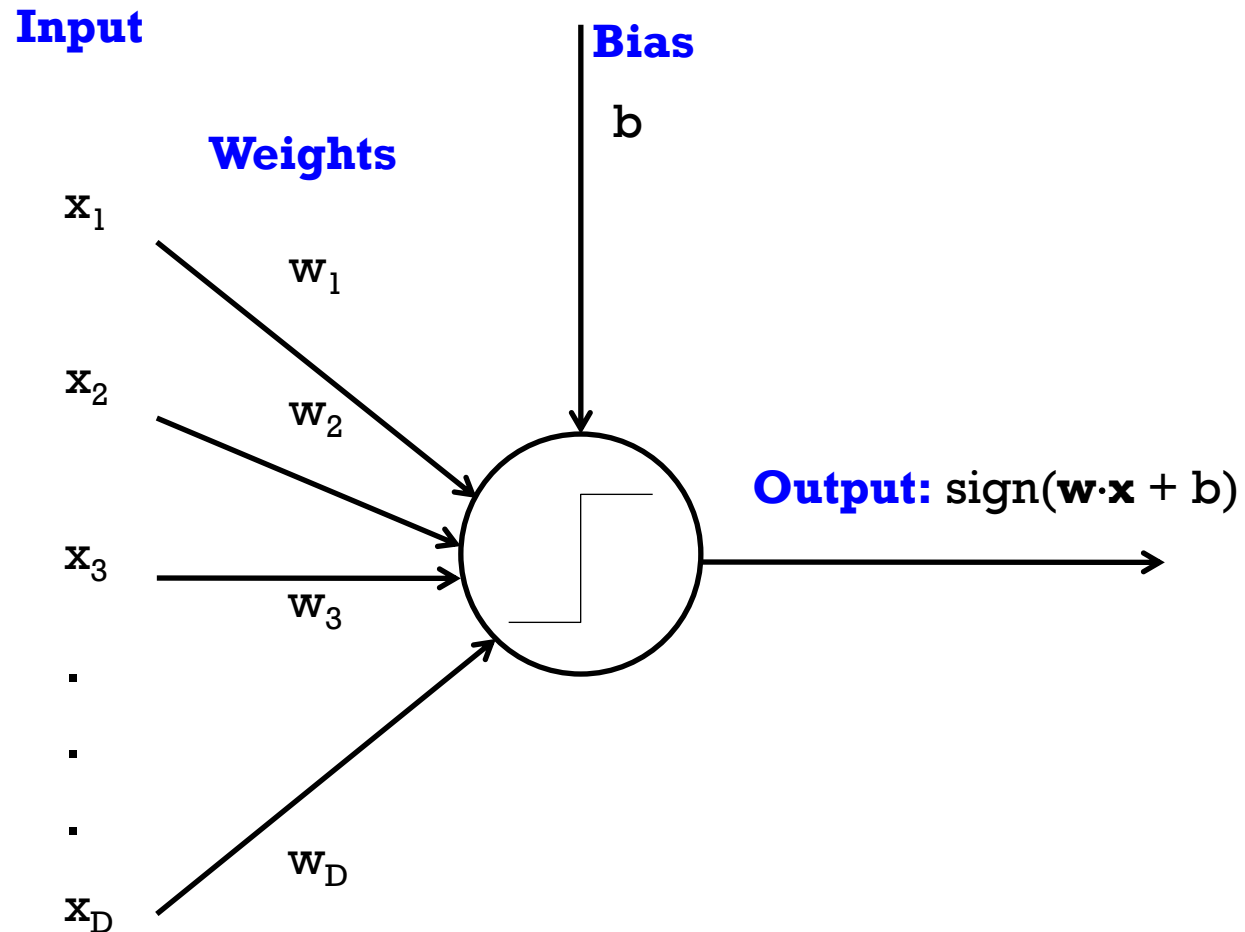
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_c(x) = c + x$	\rightarrow	$\frac{df}{dx} = 1$
$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$

PERCEPTRON

- Supervised learning of binary classifier



SINGLE NEURON AS A LINEAR CLASSIFIER

Binary Softmax classifier (*Logistic Regression*) $\sigma(\sum_i w_i x_i + b)$

SINGLE NEURON AS A LINEAR CLASSIFIER

Binary Softmax classifier (*Logistic Regression*)

$$\sigma(\sum_i w_i x_i + b)$$



Probability of one of the classes:

$$P(y_i = 1 \mid x_i; w)$$

SINGLE NEURON AS A LINEAR CLASSIFIER

Binary Softmax classifier (*Logistic Regression*)

$$\sigma(\sum_i w_i x_i + b)$$



Probability of one of the classes:

$$P(y_i = 1 \mid x_i; w)$$

Probability of the other class would be:

$$P(y_i = 0 \mid x_i; w) = 1 - P(y_i = 1 \mid x_i; w)$$

SINGLE NEURON AS A LINEAR CLASSIFIER

Binary Softmax classifier (*Logistic Regression*)

$$\sigma(\sum_i w_i x_i + b)$$



Probability of one of the classes:

$$P(y_i = 1 \mid x_i; w)$$

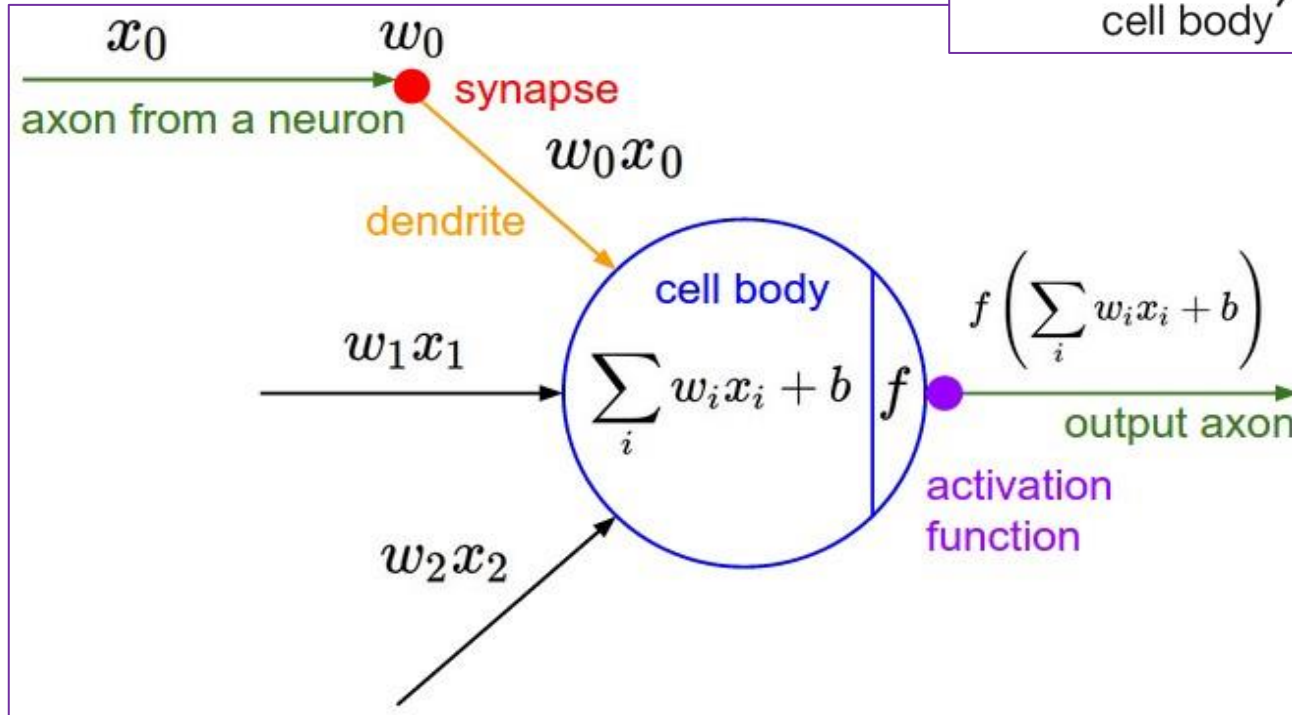
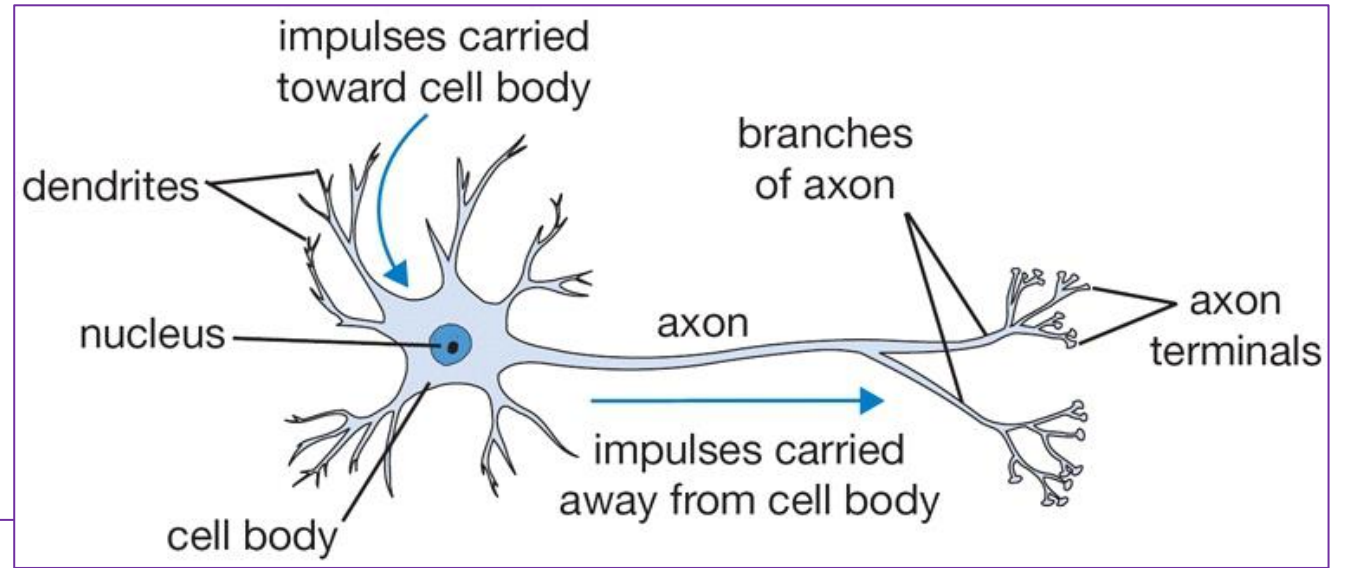
Probability of the other class would be:

$$P(y_i = 0 \mid x_i; w) = 1 - P(y_i = 1 \mid x_i; w)$$

Binary SVM classifier:

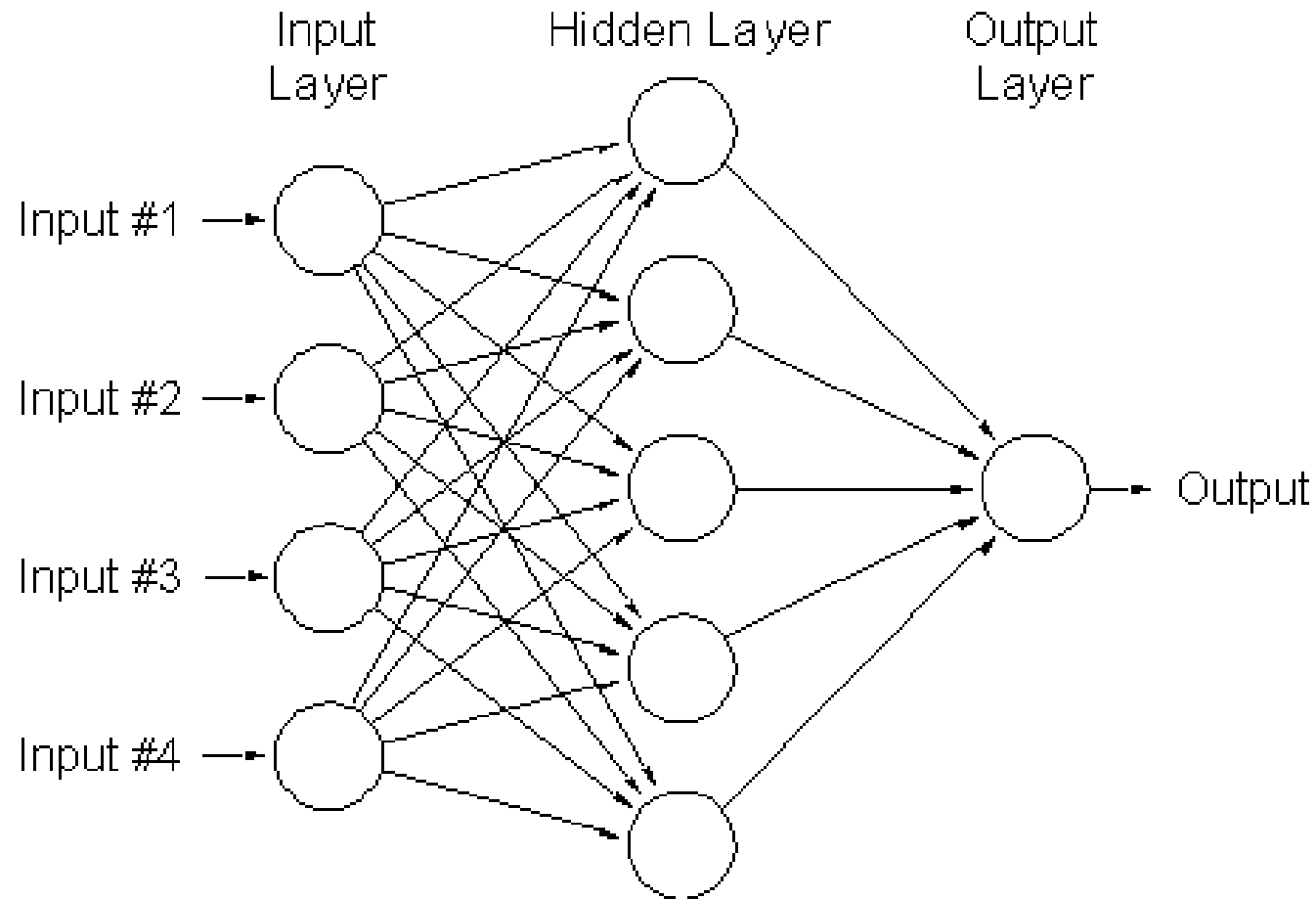
Alternatively, we could attach a max-margin hinge loss to the output of the neuron and train it to become a binary Support Vector Machine.

LOOSE INSPIRATION: HUMAN NEURONS



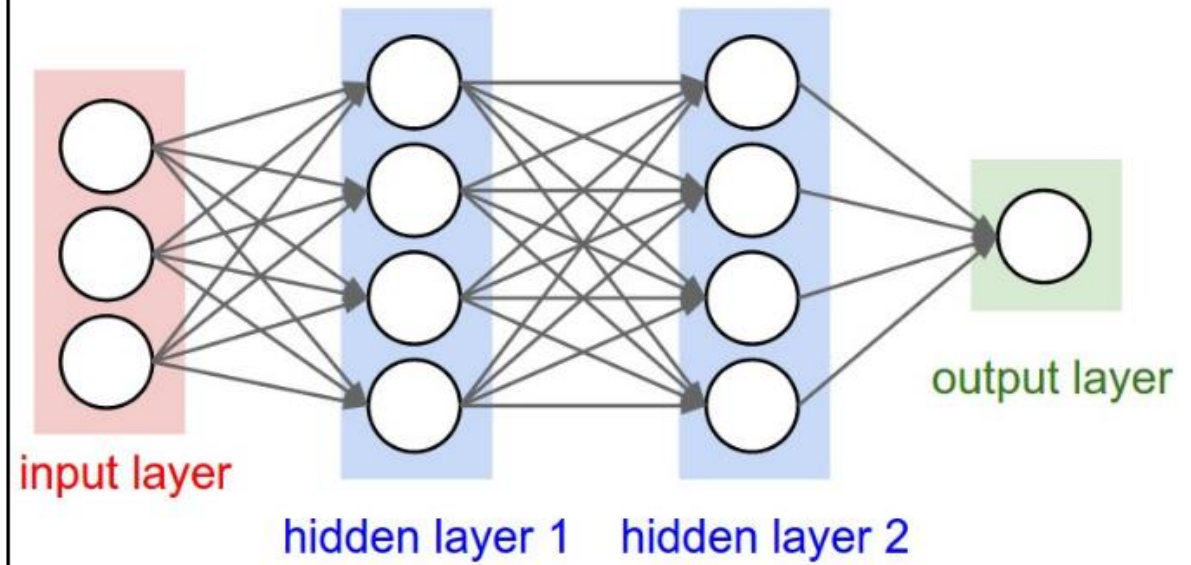
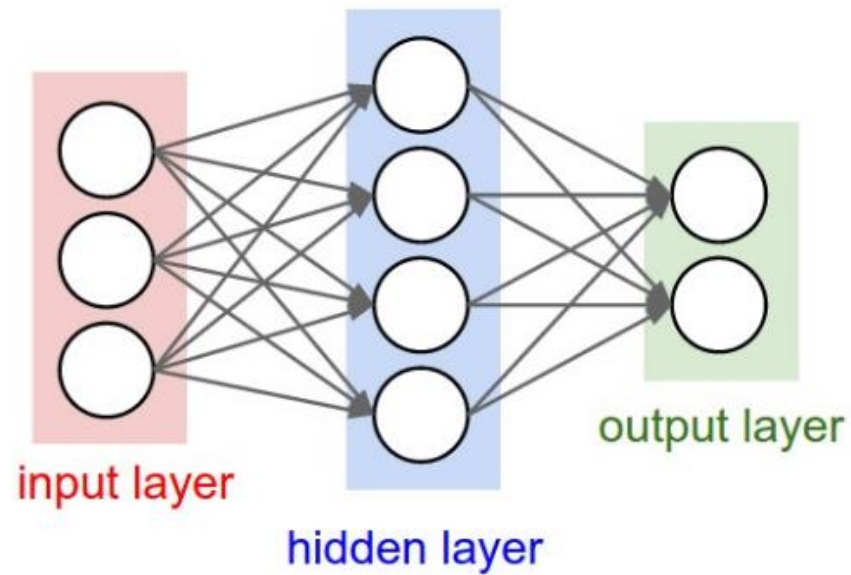
MULTI-LAYER NEURAL NETWORKS

Network with a hidden layer:



Can represent nonlinear functions (provided each perceptron has a nonlinearity)

SIZING NEURAL NETWORKS



First network (left):

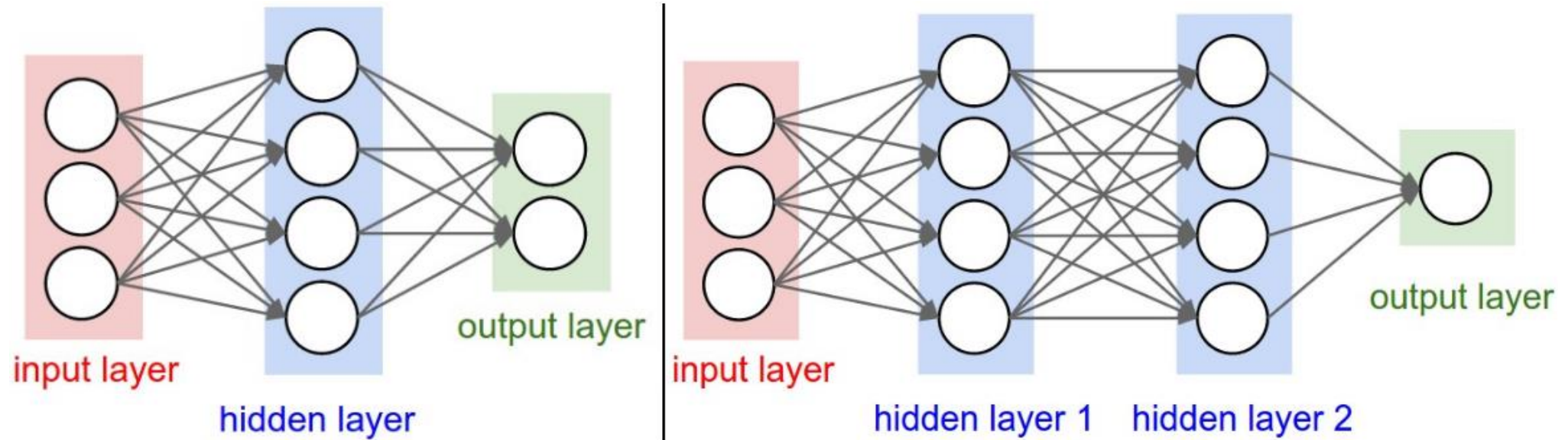
No. of neurons (not counting the inputs):

No. of learnable parameters:

Source:

<http://cs231n.github.io>

SIZING NEURAL NETWORKS

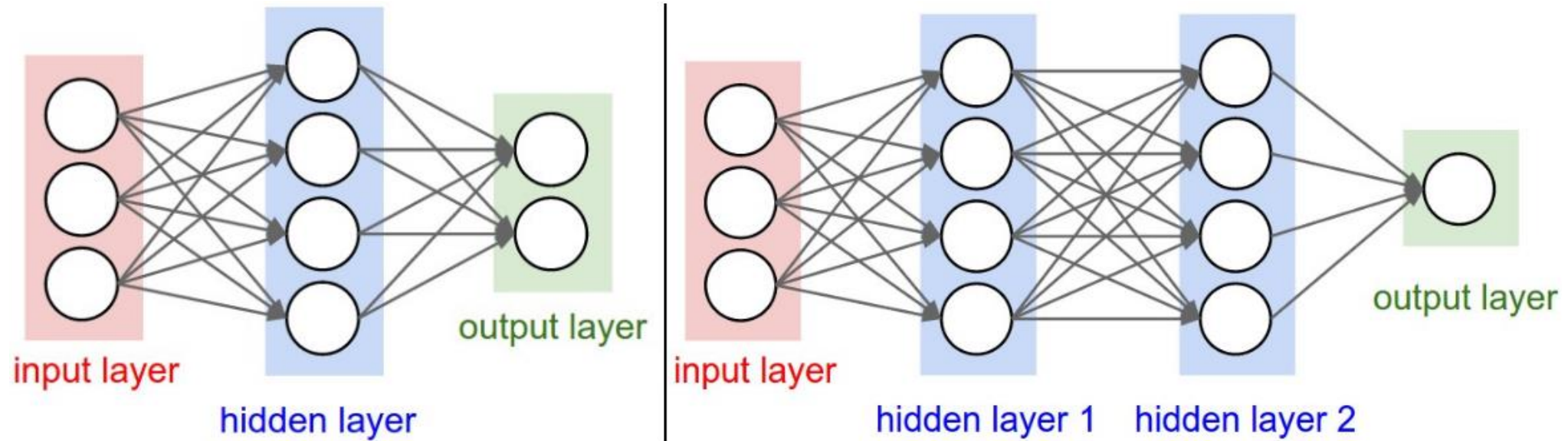


First network (left):

No. of neurons (not counting the inputs): $4 + 2 = 6$

No. of learnable parameters:

SIZING NEURAL NETWORKS

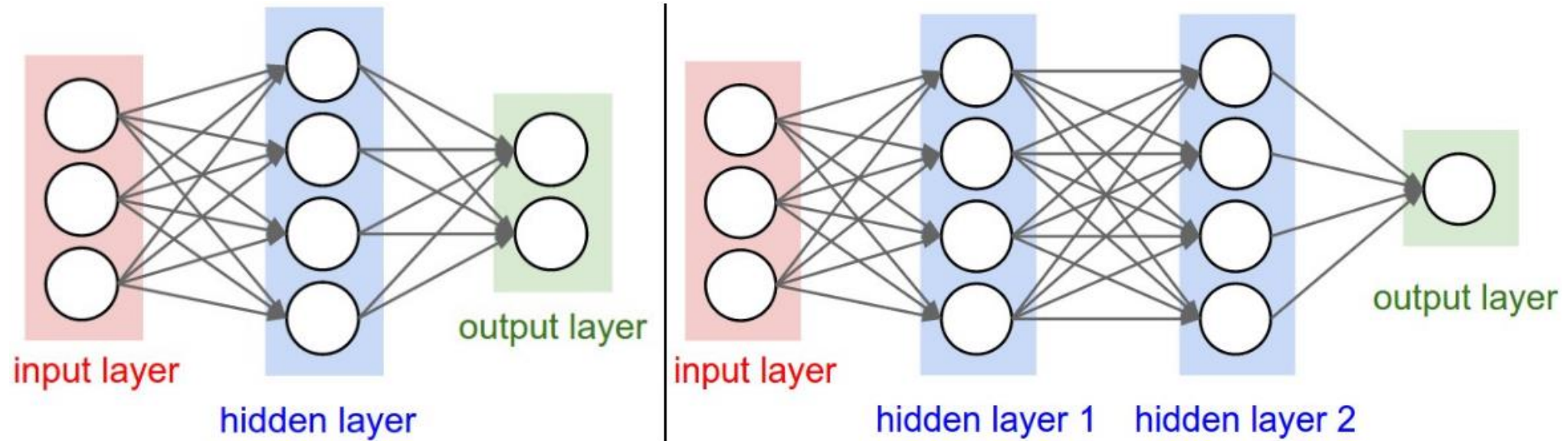


First network (left):

No. of neurons (not counting the inputs): $4 + 2 = 6$

No. of learnable parameters: $[3 \times 4] + [4 \times 2] = 20$ weights +
 $4 + 2 = 6$ biases = 26.

SIZING NEURAL NETWORKS



First network (left):

No. of neurons (not counting the inputs): $4 + 2 = 6$

No. of learnable parameters: $[3 \times 4] + [4 \times 2] = 20$ weights +
 $4 + 2 = 6$ biases = 26.

Second network (right):

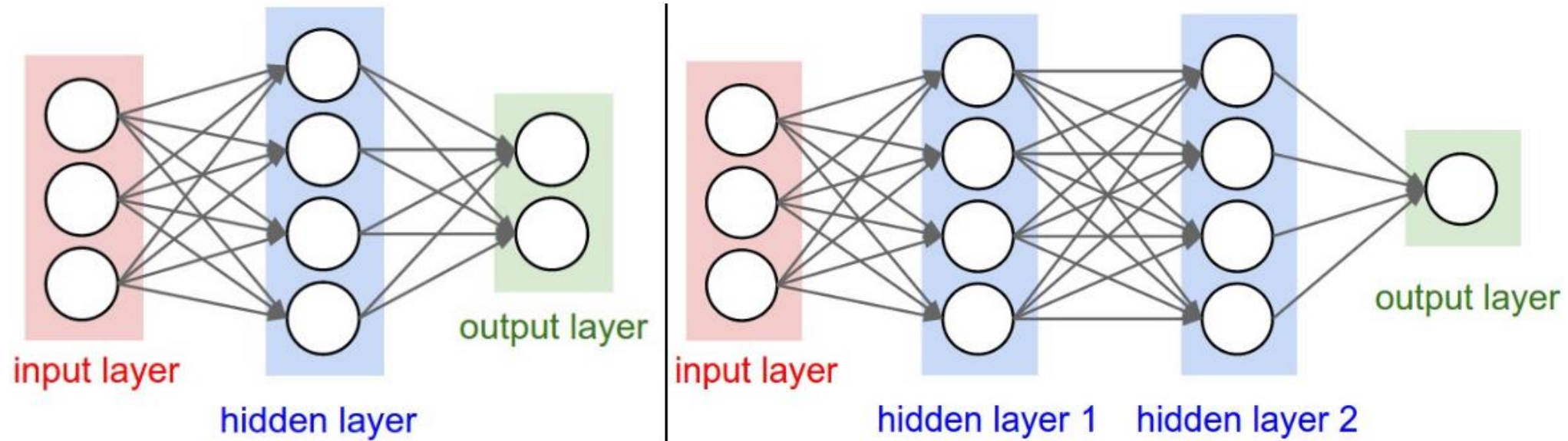
No. of neurons (not counting the inputs):

No. of learnable parameters:

Source:

<http://cs231n.github.io>

SIZING NEURAL NETWORKS



First network (left):

No. of neurons (not counting the inputs): $4 + 2 = 6$

No. of learnable parameters: $[3 \times 4] + [4 \times 2] = 20$ weights +
 $4 + 2 = 6$ biases = 26.

Second network (right):

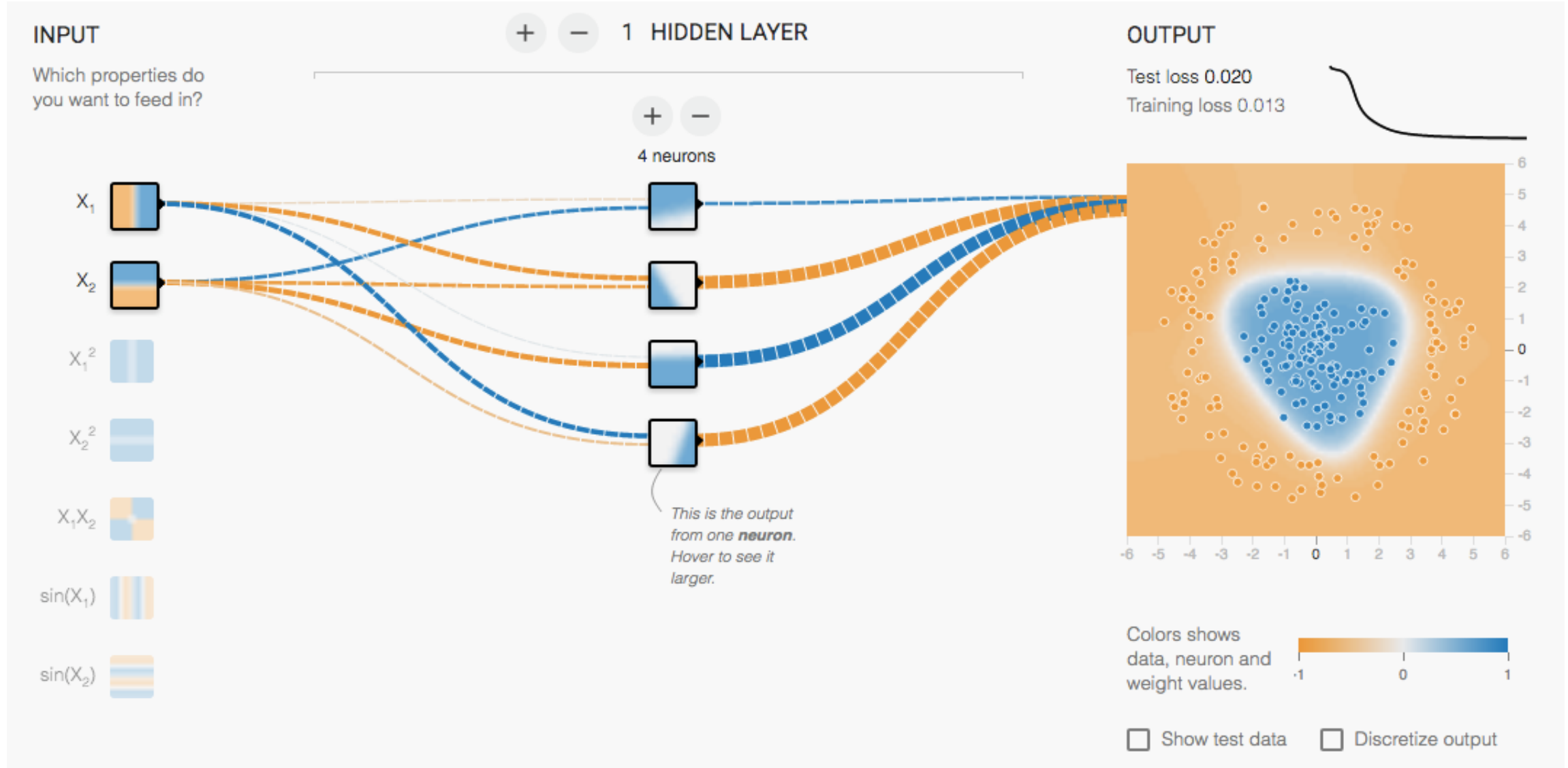
No. of neurons (not counting the inputs): $4 + 4 + 1 = 9$

No. of learnable parameters: $[3 \times 4] + [4 \times 4] + [4 \times 1] = 32$ weights +
 $4 + 4 + 1 = 9$ biases = 41.

Source:

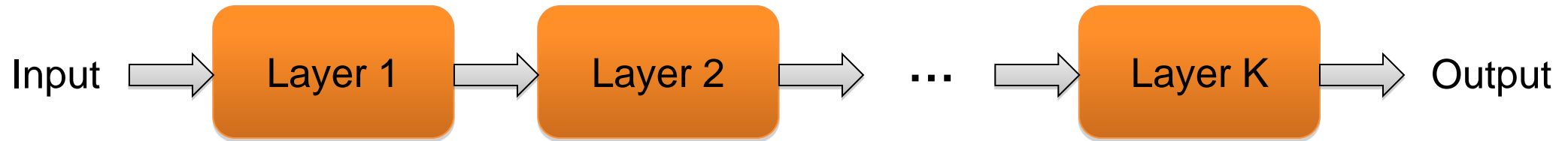
<http://cs231n.github.io>

MULTI-LAYER NETWORK DEMO



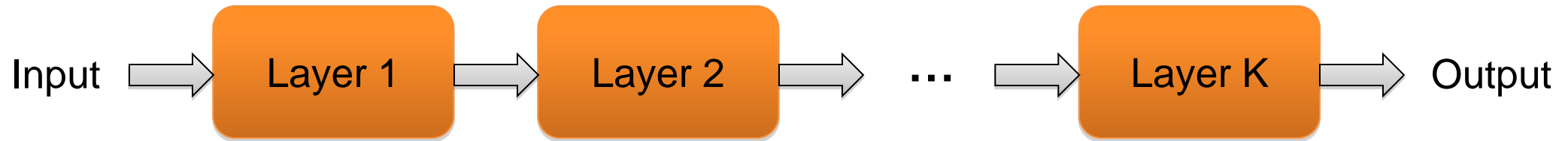
RECALL: MULTI-LAYER NEURAL NETWORKS

- The function computed by the network is a composition of the functions computed by individual layers (e.g., linear layers and nonlinearities):



RECALL: MULTI-LAYER NEURAL NETWORKS

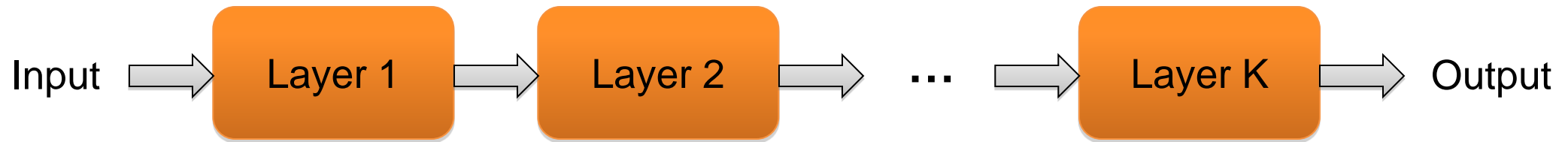
- The function computed by the network is a composition of the functions computed by individual layers (e.g., linear layers and nonlinearities):



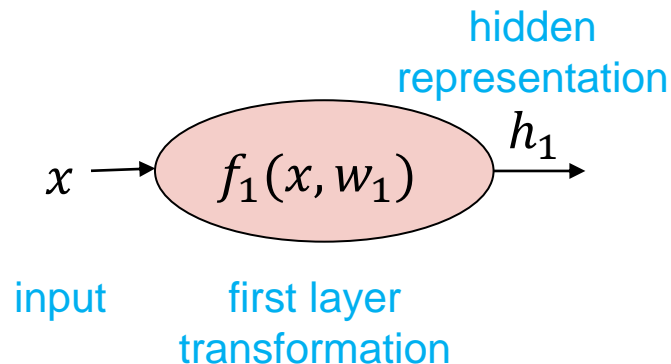
- More precisely:

RECALL: MULTI-LAYER NEURAL NETWORKS

- The function computed by the network is a composition of the functions computed by individual layers (e.g., linear layers and nonlinearities):

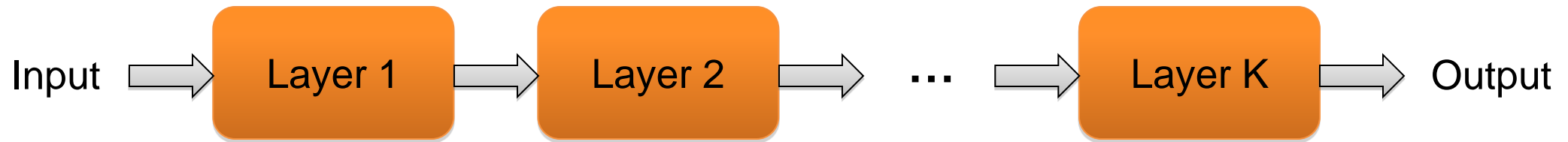


- More precisely:

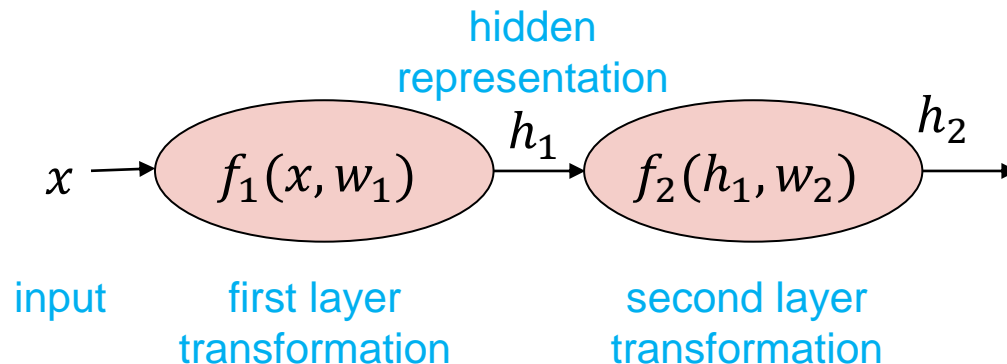


RECALL: MULTI-LAYER NEURAL NETWORKS

- The function computed by the network is a composition of the functions computed by individual layers (e.g., linear layers and nonlinearities):

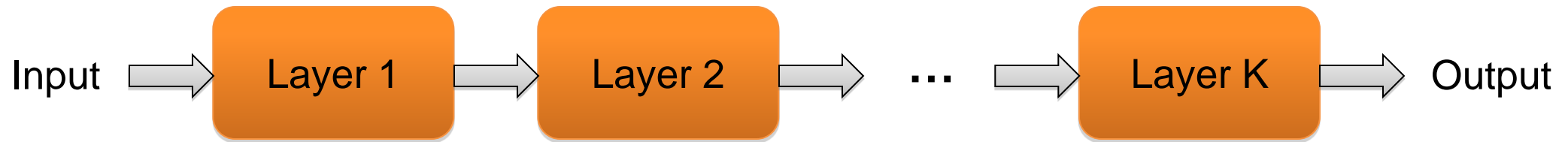


- More precisely:

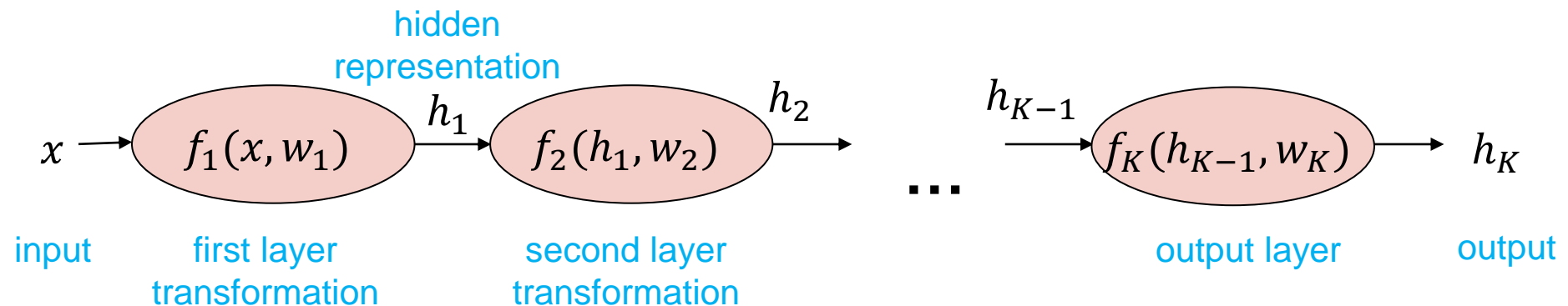


RECALL: MULTI-LAYER NEURAL NETWORKS

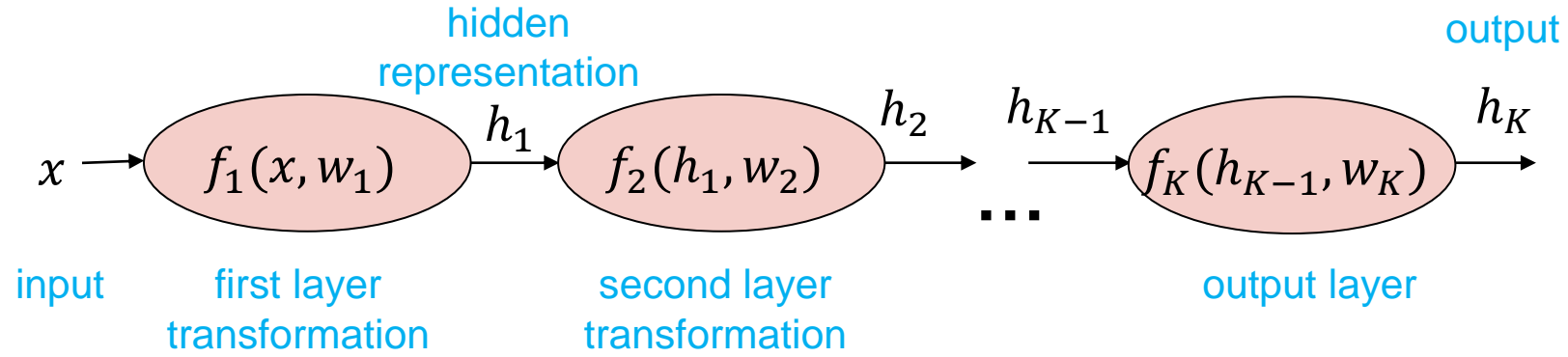
- The function computed by the network is a composition of the functions computed by individual layers (e.g., linear layers and nonlinearities):



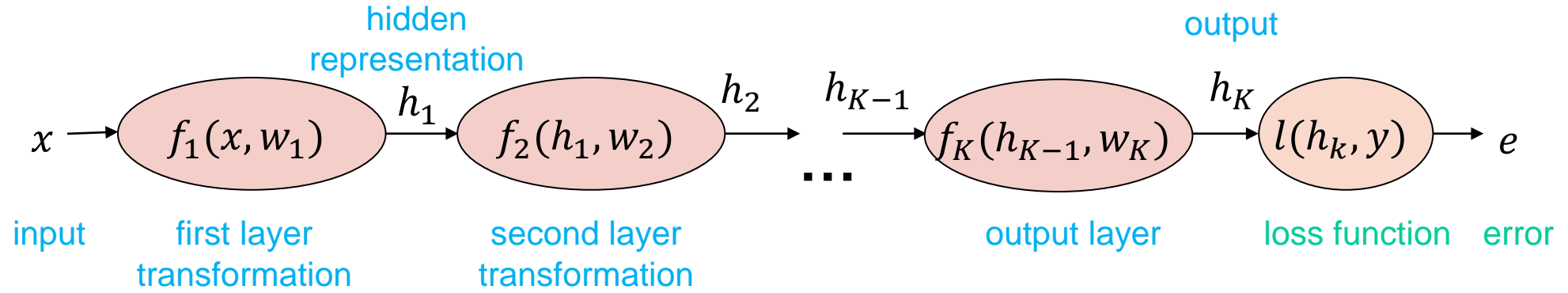
- More precisely:



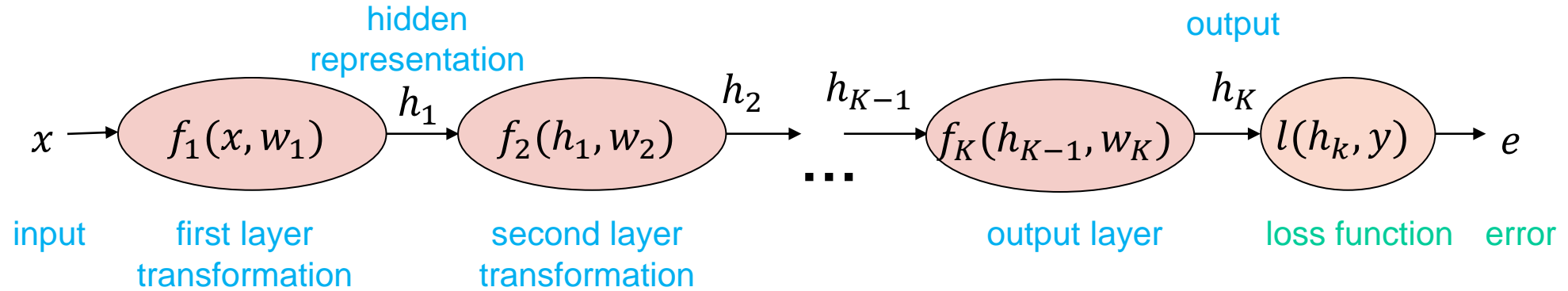
TRAINING A MULTI-LAYER NETWORK



TRAINING A MULTI-LAYER NETWORK

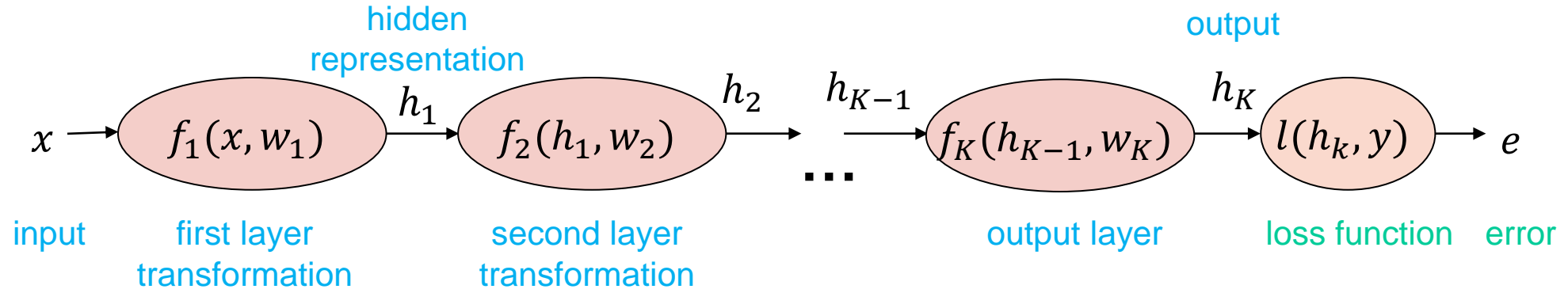


TRAINING A MULTI-LAYER NETWORK



- What is the SGD update for the parameters w_k of the k th layer?

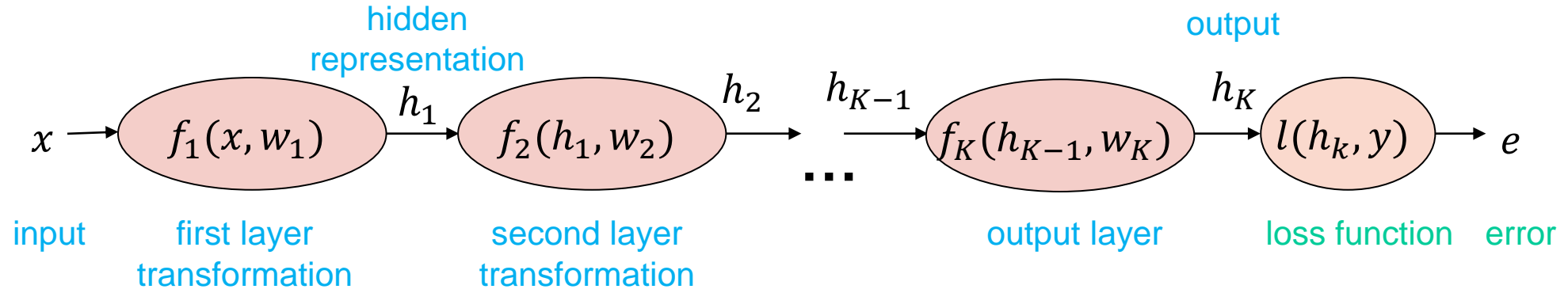
TRAINING A MULTI-LAYER NETWORK



- What is the SGD update for the parameters w_k of the k th layer?

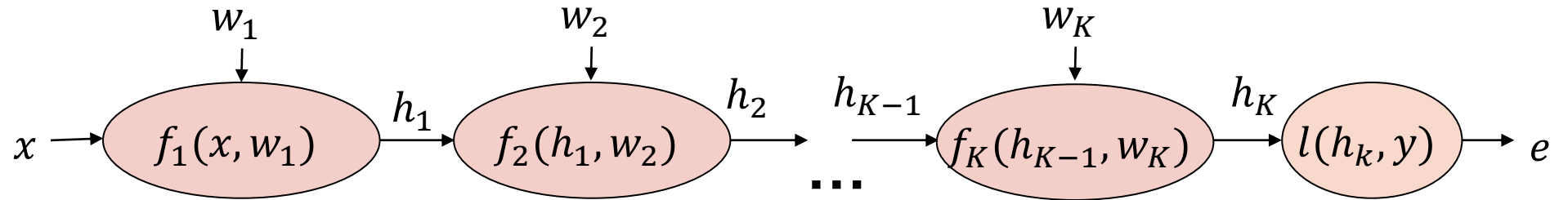
$$w_k \leftarrow w_k - \eta \frac{\partial e}{\partial w_k}$$

TRAINING A MULTI-LAYER NETWORK



- What is the SGD update for the parameters w_k of the k th layer?
$$w_k \leftarrow w_k - \eta \frac{\partial e}{\partial w_k}$$
- To train the network, we need to find the **gradient of the error w.r.t. the parameters of each layer, $\frac{\partial e}{\partial w_k}$**

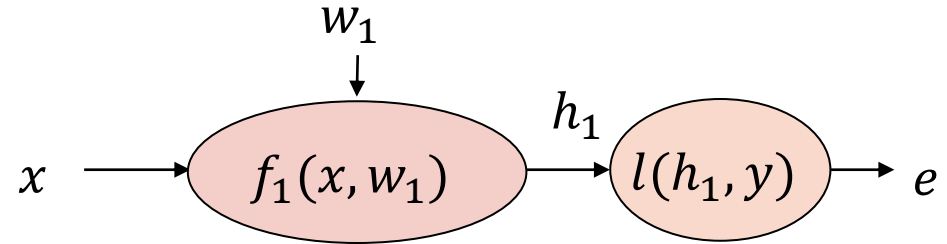
COMPUTATION GRAPH



CHAIN RULE

Let's start with $k = 1$

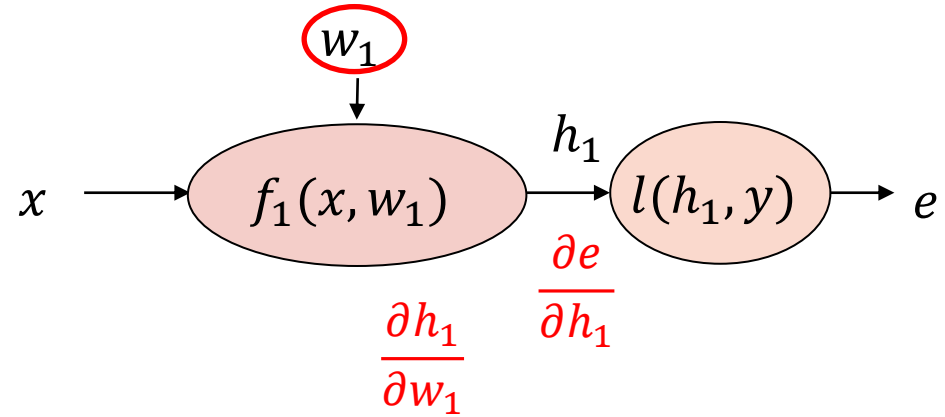
$$e = l(f_1(x, w_1), y)$$



CHAIN RULE

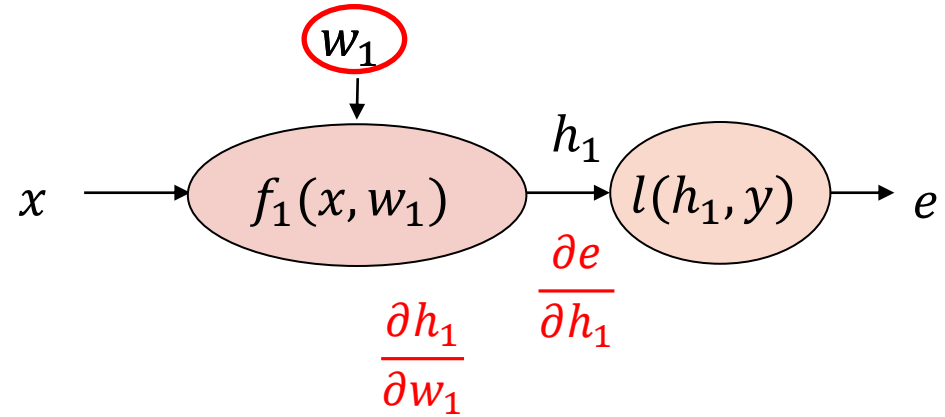
Let's start with $k = 1$

$$e = l(f_1(x, w_1), y)$$



CHAIN RULE

Let's start with $k = 1$

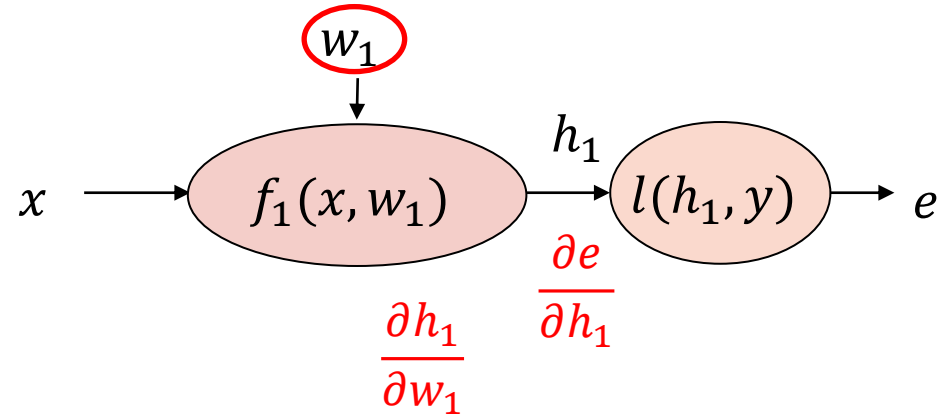


$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

CHAIN RULE

Let's start with $k = 1$



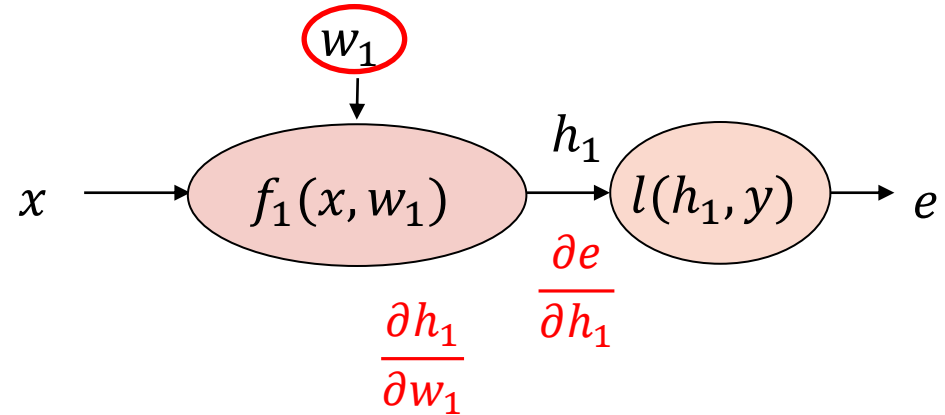
$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

Example: $e = (y - w_1^T x)^2$

CHAIN RULE

Let's start with $k = 1$



$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

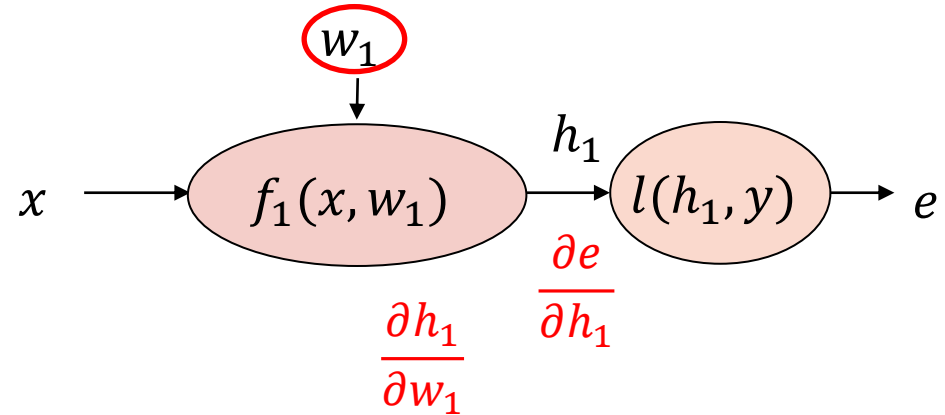
Example: $e = (y - w_1^T x)^2$

$$h_1 = f_1(x, w_1) = w_1^T x$$

$$e = l(h_1, y) = (y - h_1)^2$$

CHAIN RULE

Let's start with $k = 1$



$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

Example: $e = (y - w_1^T x)^2$

$$h_1 = f_1(x, w_1) = w_1^T x$$

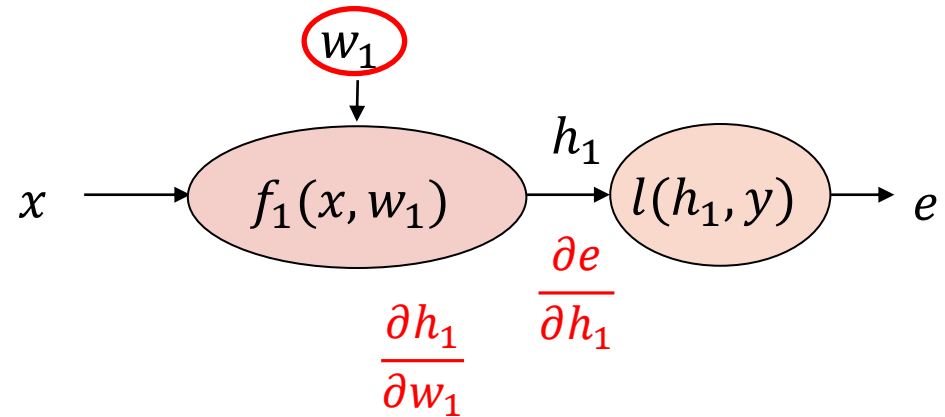
$$e = l(h_1, y) = (y - h_1)^2$$

$$\frac{\partial h_1}{\partial w_1} = \text{??????}$$

$$\frac{\partial e}{\partial h_1} = \text{??????}$$

CHAIN RULE

Let's start with $k = 1$



$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

Example: $e = (y - w_1^T x)^2$

$$h_1 = f_1(x, w_1) = w_1^T x$$

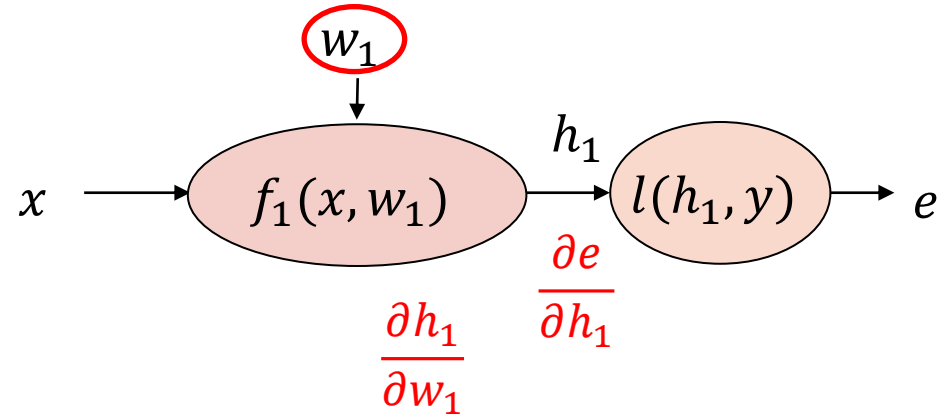
$$e = l(h_1, y) = (y - h_1)^2$$

$$\frac{\partial h_1}{\partial w_1} = x$$

$$\frac{\partial e}{\partial h_1} = -2(y - h_1) = -2(y - w_1^T x)$$

CHAIN RULE

Let's start with $k = 1$



$$e = l(f_1(x, w_1), y)$$

$$\frac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

Example: $e = (y - w_1^T x)^2$

$$h_1 = f_1(x, w_1) = w_1^T x$$

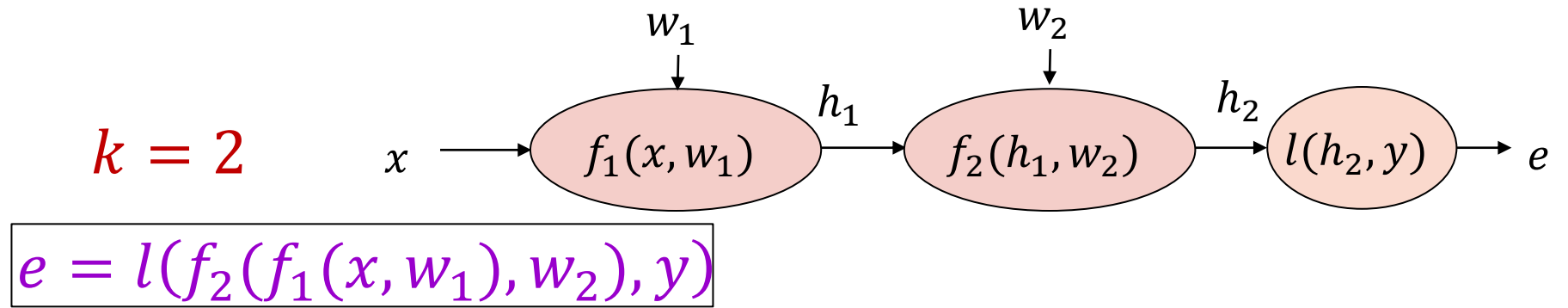
$$e = l(h_1, y) = (y - h_1)^2$$

$$\frac{\partial h_1}{\partial w_1} = x$$

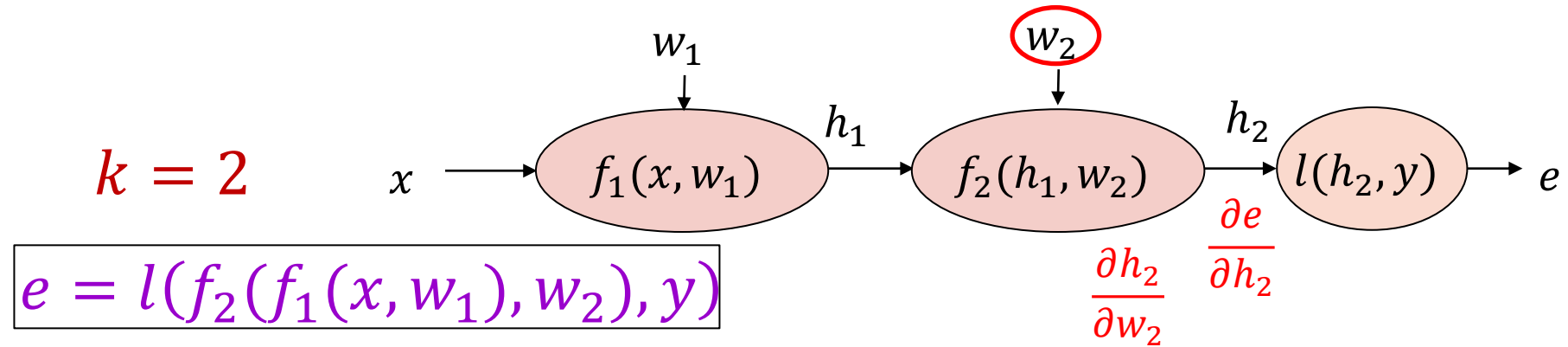
$$\frac{\partial e}{\partial h_1} = -2(y - h_1) = -2(y - w_1^T x)$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_1} \frac{\partial h_1}{\partial w_1} = -2x(y - w_1^T x)$$

CHAIN RULE

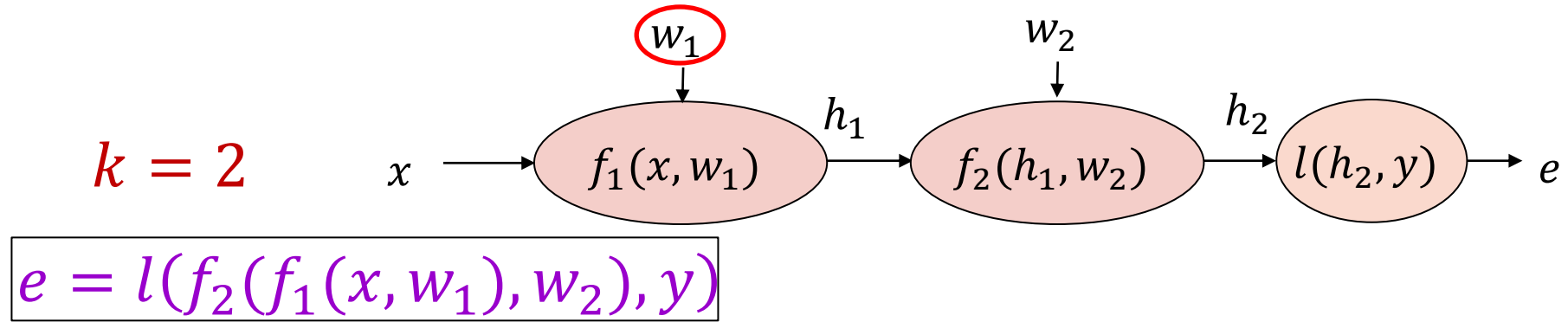


CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

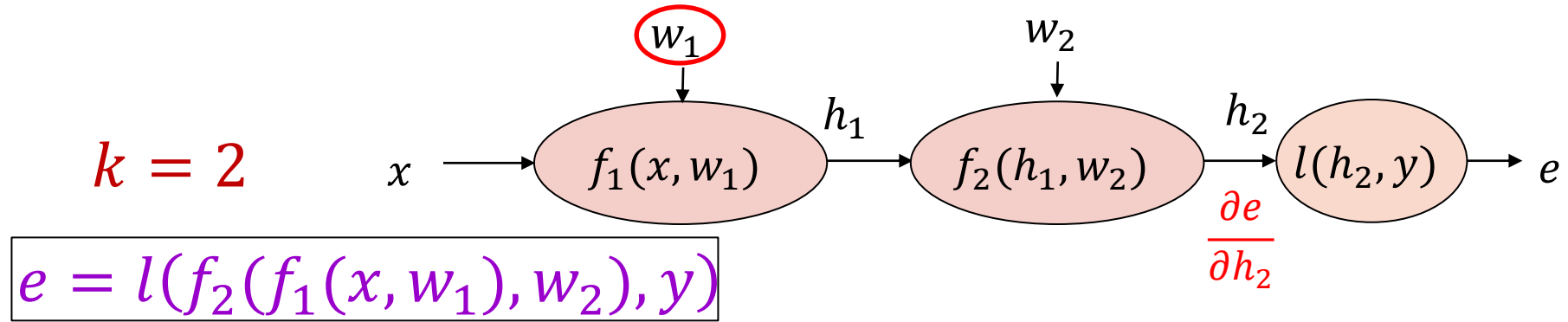
CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} =$$

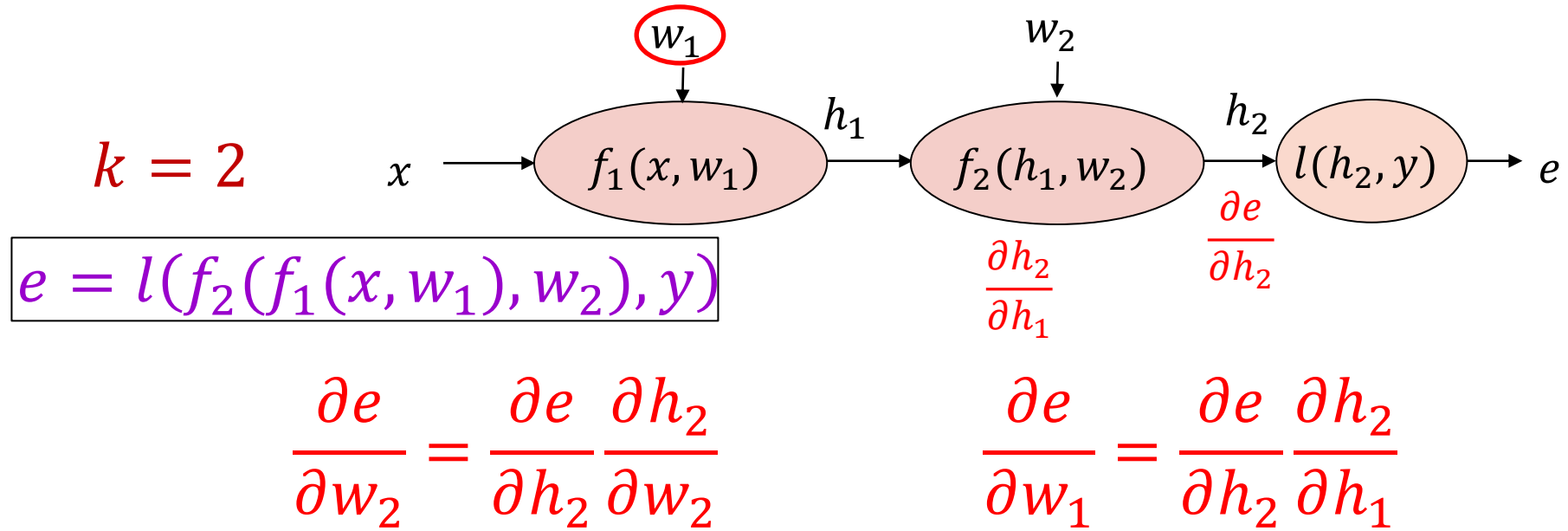
CHAIN RULE



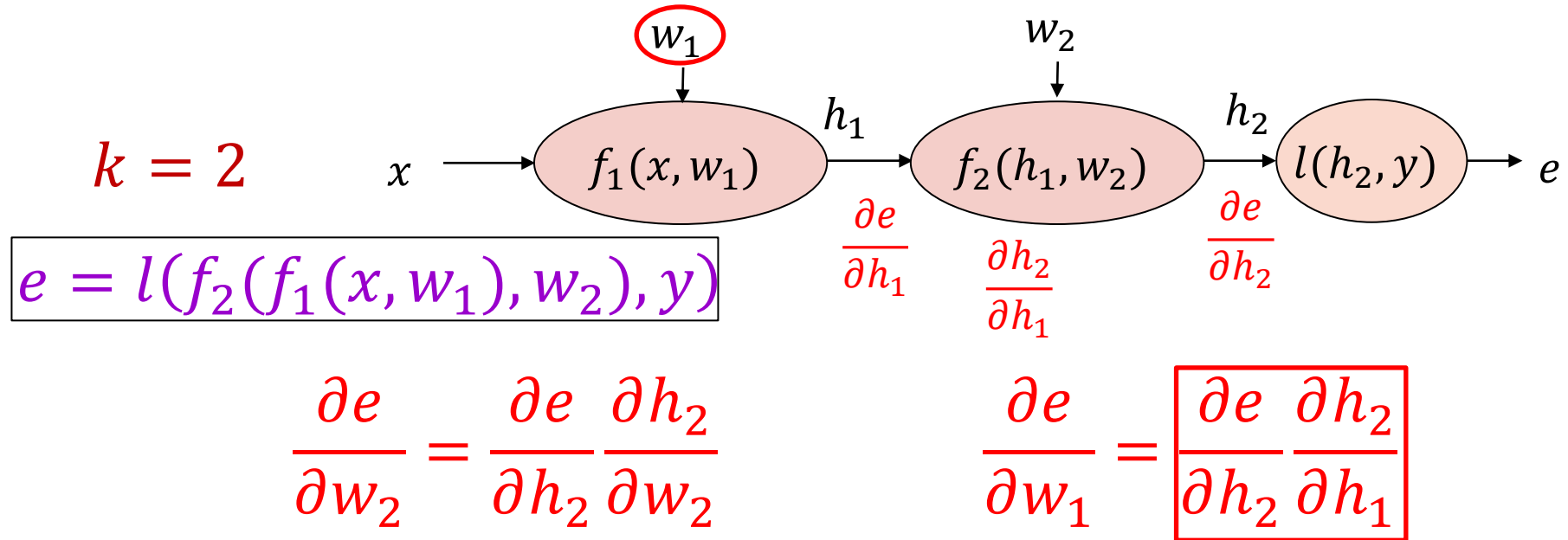
$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2}$$

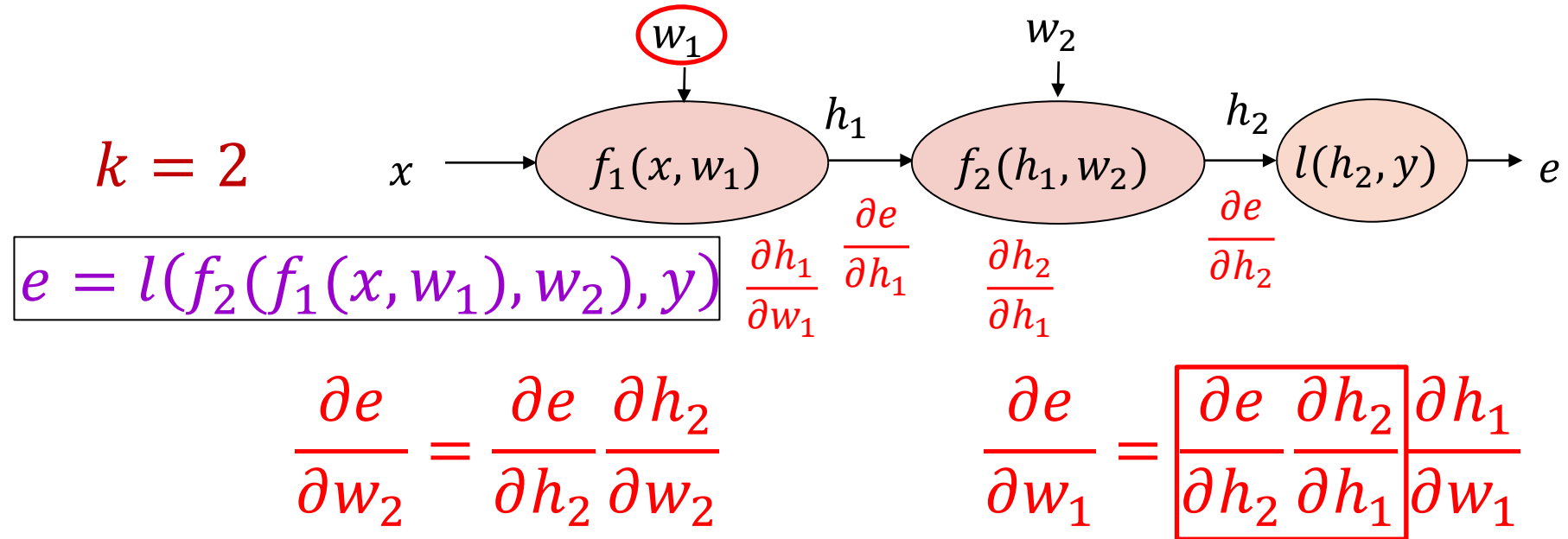
CHAIN RULE



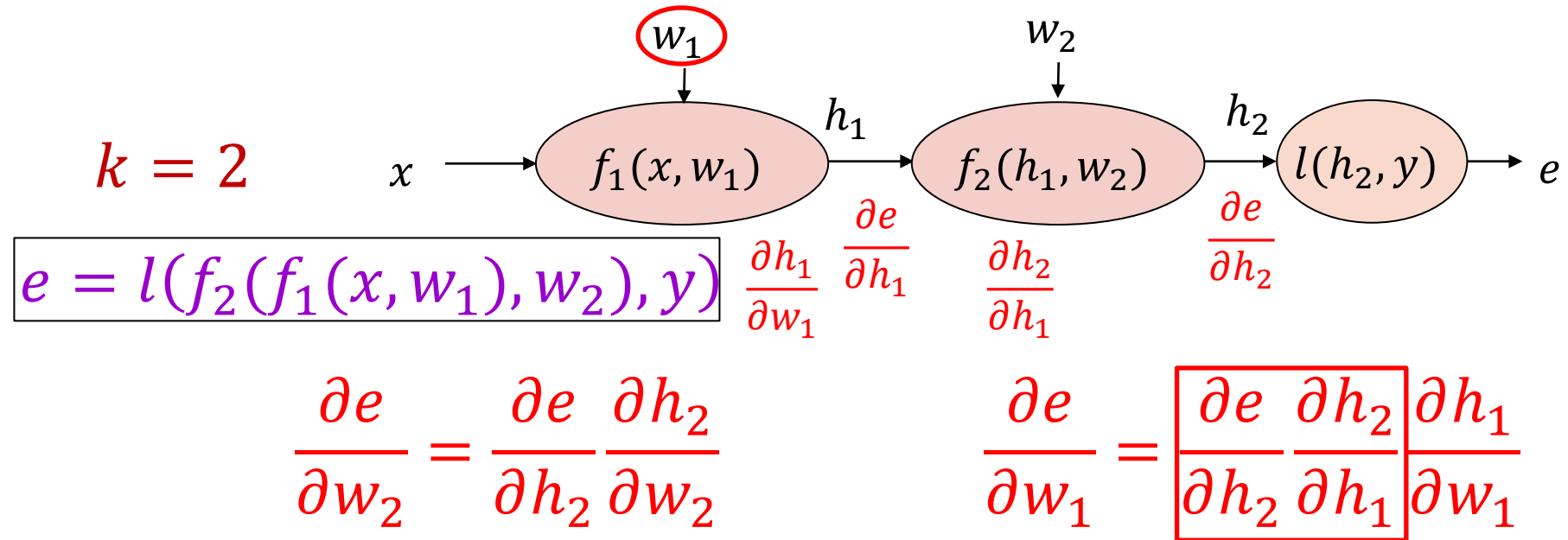
CHAIN RULE



CHAIN RULE

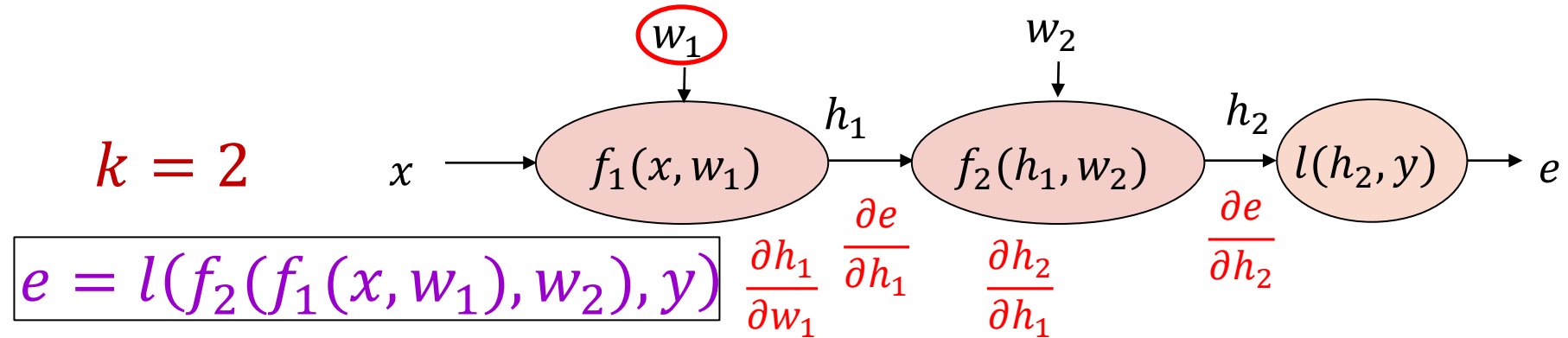


CHAIN RULE



Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} = \boxed{\frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1}} \frac{\partial h_1}{\partial w_1}$$

Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

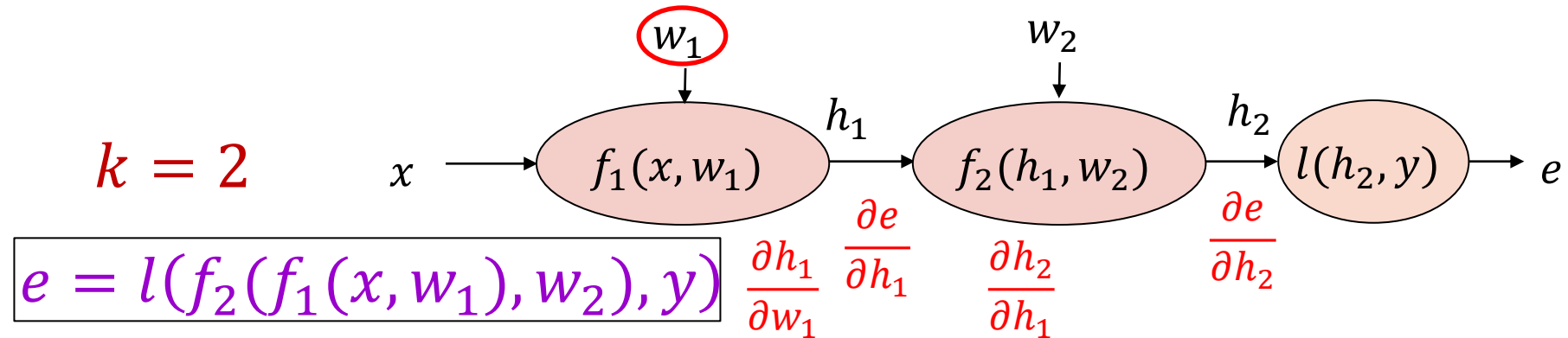
$$h_1 = f_1(x, w_1) = w_1^T x$$

$$h_2 = f_2(h_1) = \sigma(h_1)$$

$$e = l(h_2, 1) = -\log(h_2)$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} = \text{??????}$$

CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} = \boxed{\frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1}} \frac{\partial h_1}{\partial w_1}$$

Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

$$h_1 = f_1(x, w_1) = w_1^T x$$

$$h_2 = f_2(h_1) = \sigma(h_1)$$

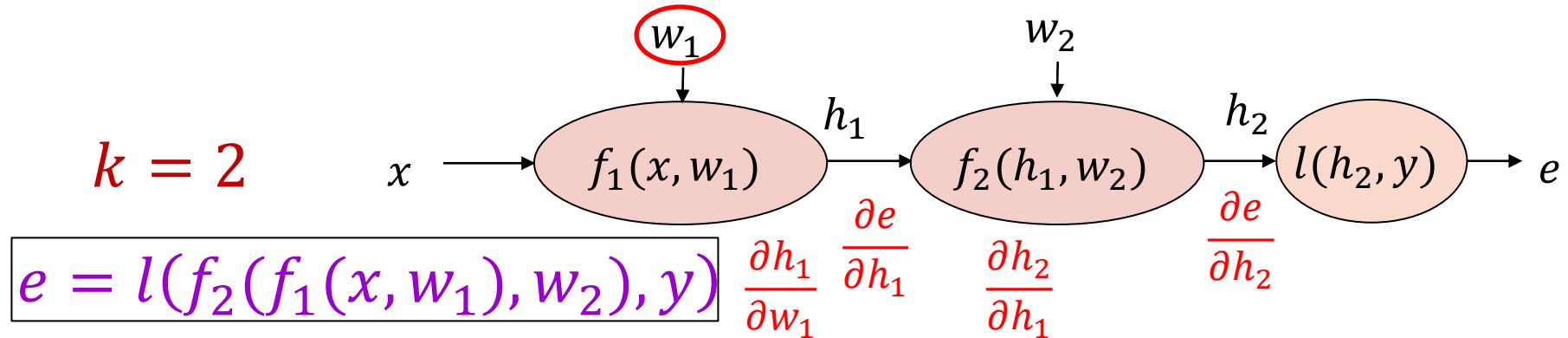
$$e = l(h_2, 1) = -\log(h_2)$$

$$\frac{\partial h_1}{\partial w_1} = ??? \quad \frac{\partial h_2}{\partial h_1} = ???$$

$$\frac{\partial e}{\partial h_2} = ???$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} = ???????$$

CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} = \boxed{\frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1}} \frac{\partial h_1}{\partial w_1}$$

Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

$$h_1 = f_1(x, w_1) = w_1^T x$$

$$h_2 = f_2(h_1) = \sigma(h_1)$$

$$e = l(h_2, 1) = -\log(h_2)$$

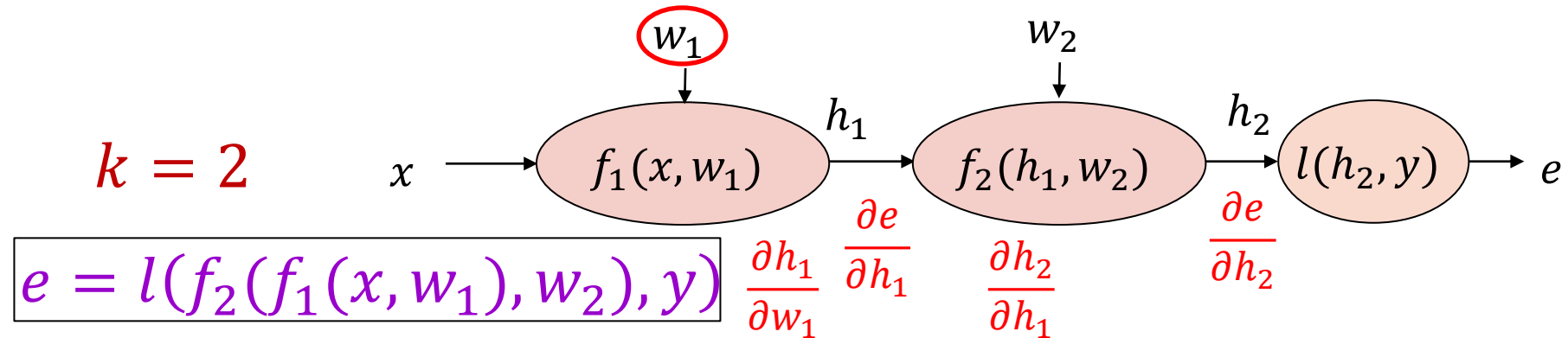
$$\frac{\partial h_1}{\partial w_1} = x$$

$$\frac{\partial h_2}{\partial h_1} = \sigma'(h_1) = \sigma(h_1)(1 - \sigma(h_1))$$

$$\frac{\partial e}{\partial h_2} = -\frac{1}{h_2}$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} = \text{??????}$$

CHAIN RULE



$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$\frac{\partial e}{\partial w_1} = \boxed{\frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1}} \frac{\partial h_1}{\partial w_1}$$

Example: $e = -\log(\sigma(w_1^T x))$ (assume $y = 1$)

$$h_1 = f_1(x, w_1) = w_1^T x$$

$$h_2 = f_2(h_1) = \sigma(h_1)$$

$$e = l(h_2, 1) = -\log(h_2)$$

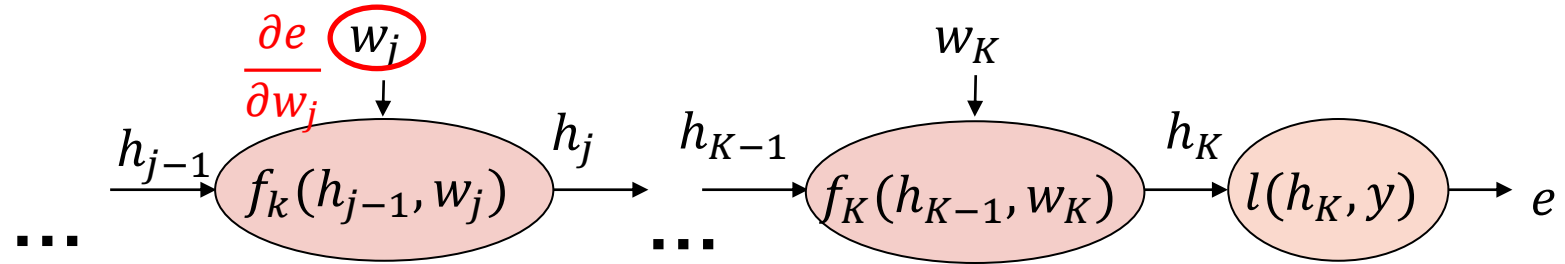
$$\frac{\partial h_1}{\partial w_1} =$$

$$\frac{\partial h_2}{\partial h_1} = \sigma'(h_1) = \sigma(h_1)(1 - \sigma(h_1))$$

$$\frac{\partial e}{\partial h_2} = -\frac{1}{h_2}$$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} = -\frac{1}{\sigma(w_1^T x)} \sigma(w_1^T x) (1 - \sigma(w_1^T x)) x = -\sigma(-w_1^T x) x$$

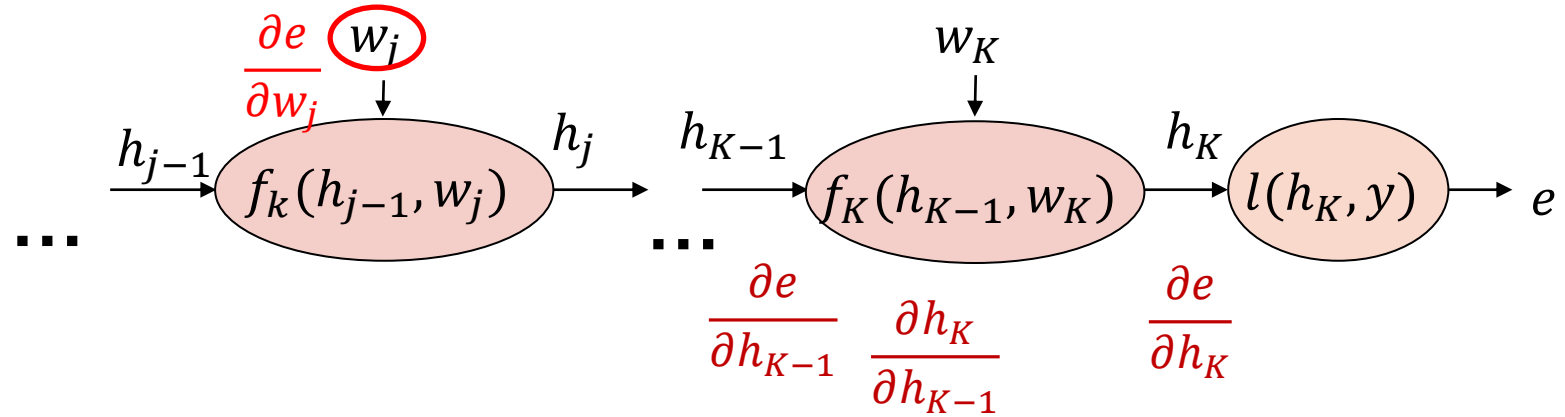
CHAIN RULE



General case:

$$\frac{\partial e}{\partial w_j} = \frac{\partial e}{\partial h_K} \frac{\partial h_K}{\partial h_{K-1}} \cdots \frac{\partial h_{j+1}}{\partial h_j} \frac{\partial h_j}{\partial w_j}$$

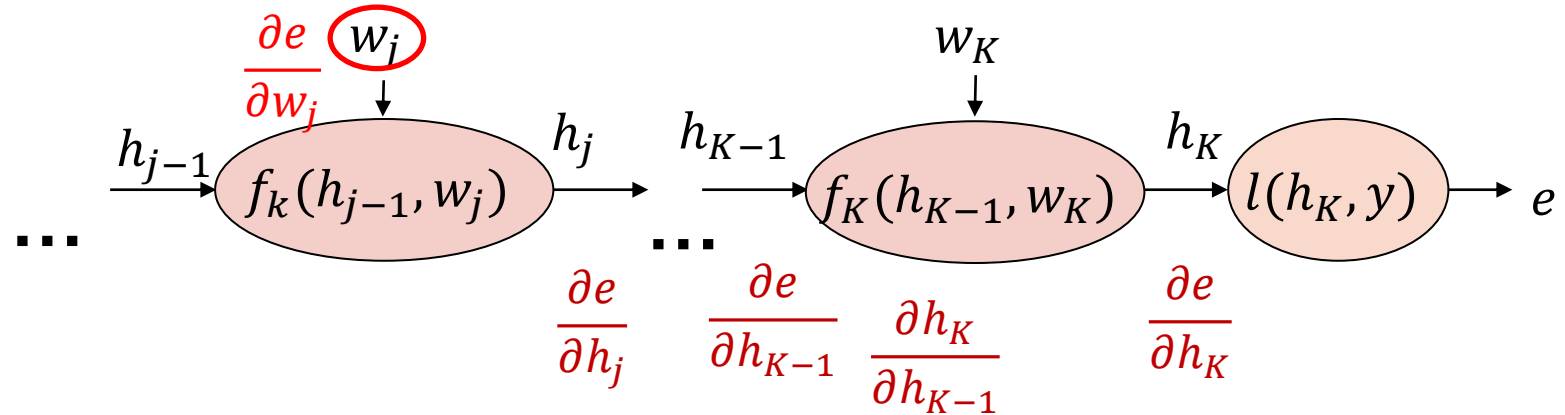
CHAIN RULE



General case:

$$\frac{\partial e}{\partial w_j} = \frac{\partial e}{\partial h_K} \frac{\partial h_K}{\partial h_{K-1}} \dots \frac{\partial h_{j+1}}{\partial h_j} \frac{\partial h_j}{\partial w_j}$$

CHAIN RULE



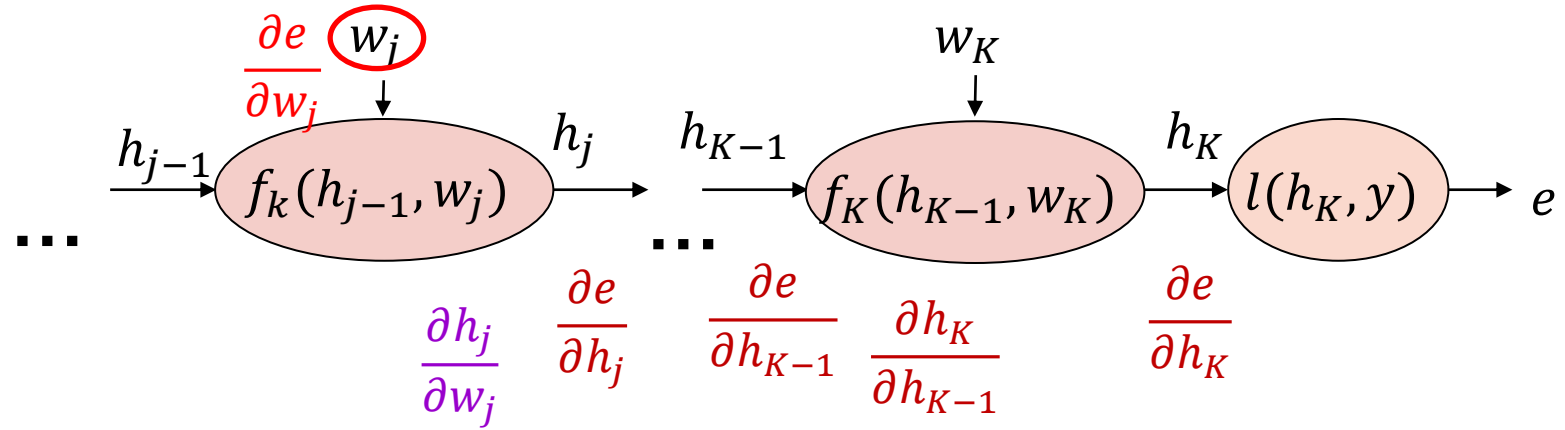
General case:

$$\frac{\partial e}{\partial w_j} = \underbrace{\frac{\partial e}{\partial h_K} \frac{\partial h_K}{\partial h_{K-1}} \cdots \frac{\partial h_{j+1}}{\partial h_j}}_{\text{Upstream gradient}} \frac{\partial h_j}{\partial w_j}$$

Upstream gradient

$\frac{\partial e}{\partial h_j}$

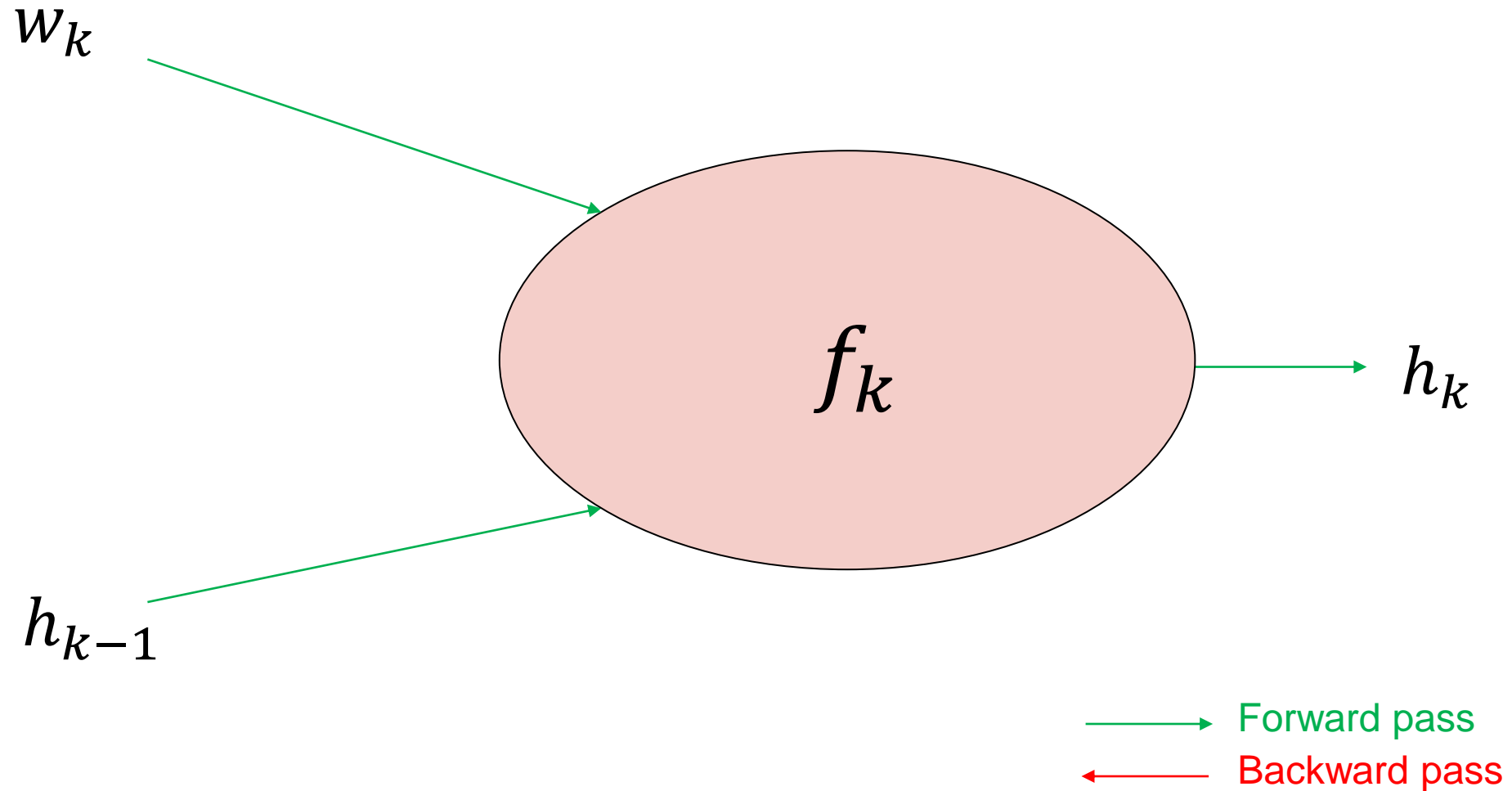
CHAIN RULE



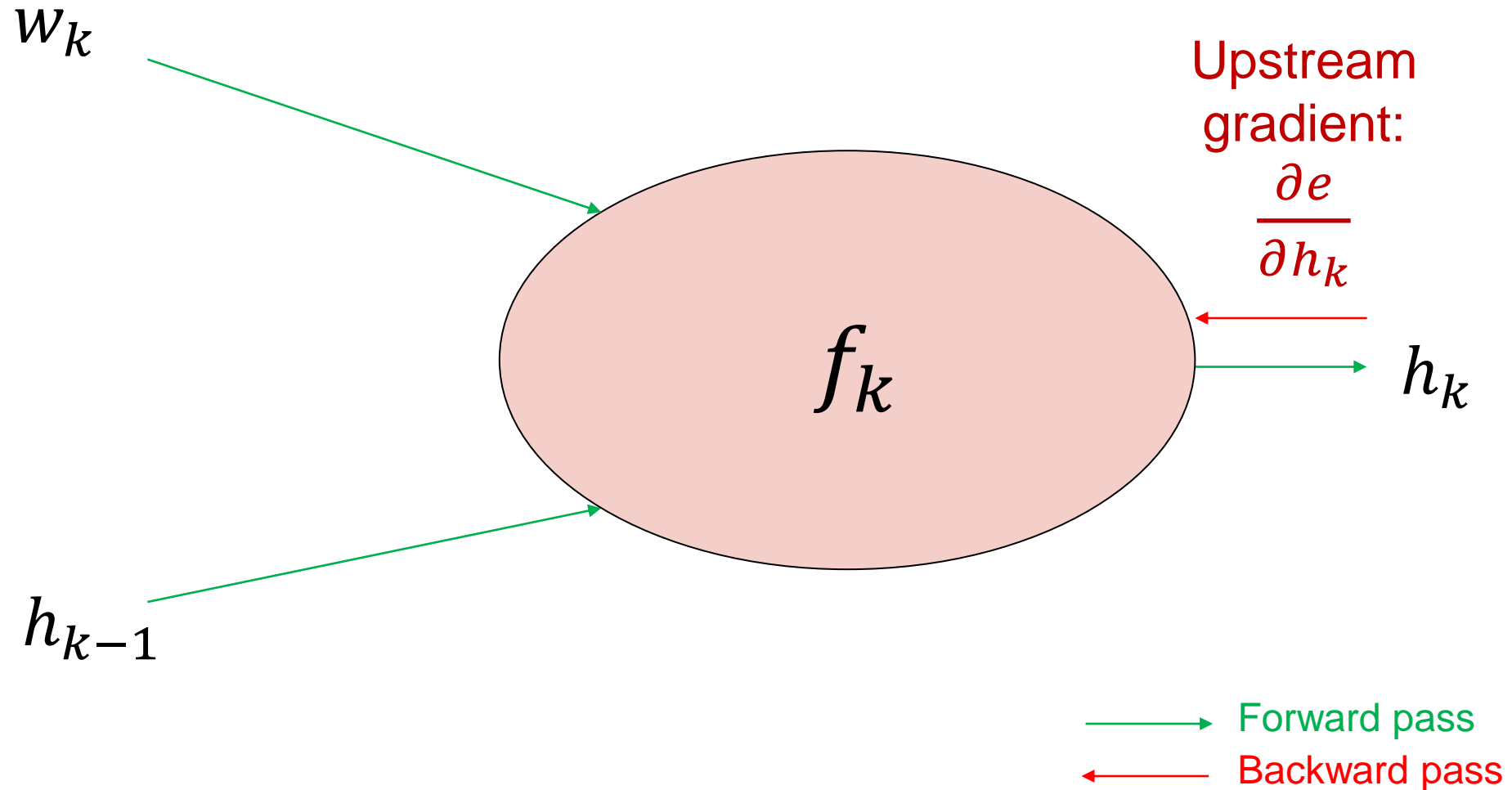
General case:

$$\frac{\partial e}{\partial w_j} = \underbrace{\frac{\partial e}{\partial h_K} \frac{\partial h_K}{\partial h_{K-1}} \cdots \frac{\partial h_{j+1}}{\partial h_j}}_{\text{Upstream gradient}} \underbrace{\frac{\partial h_j}{\partial w_j}}_{\text{Local gradient}}$$

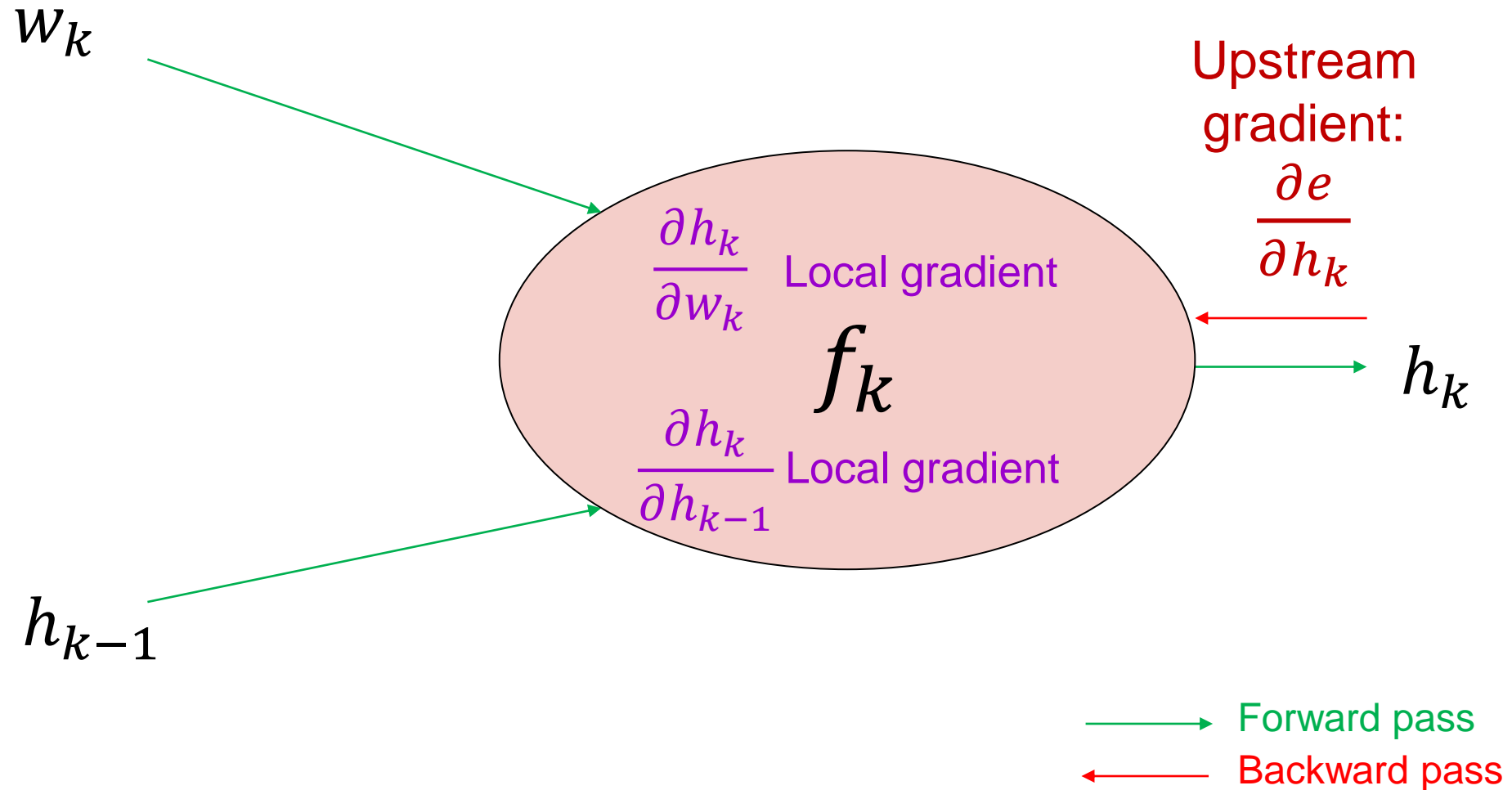
BACKPROPAGATION SUMMARY



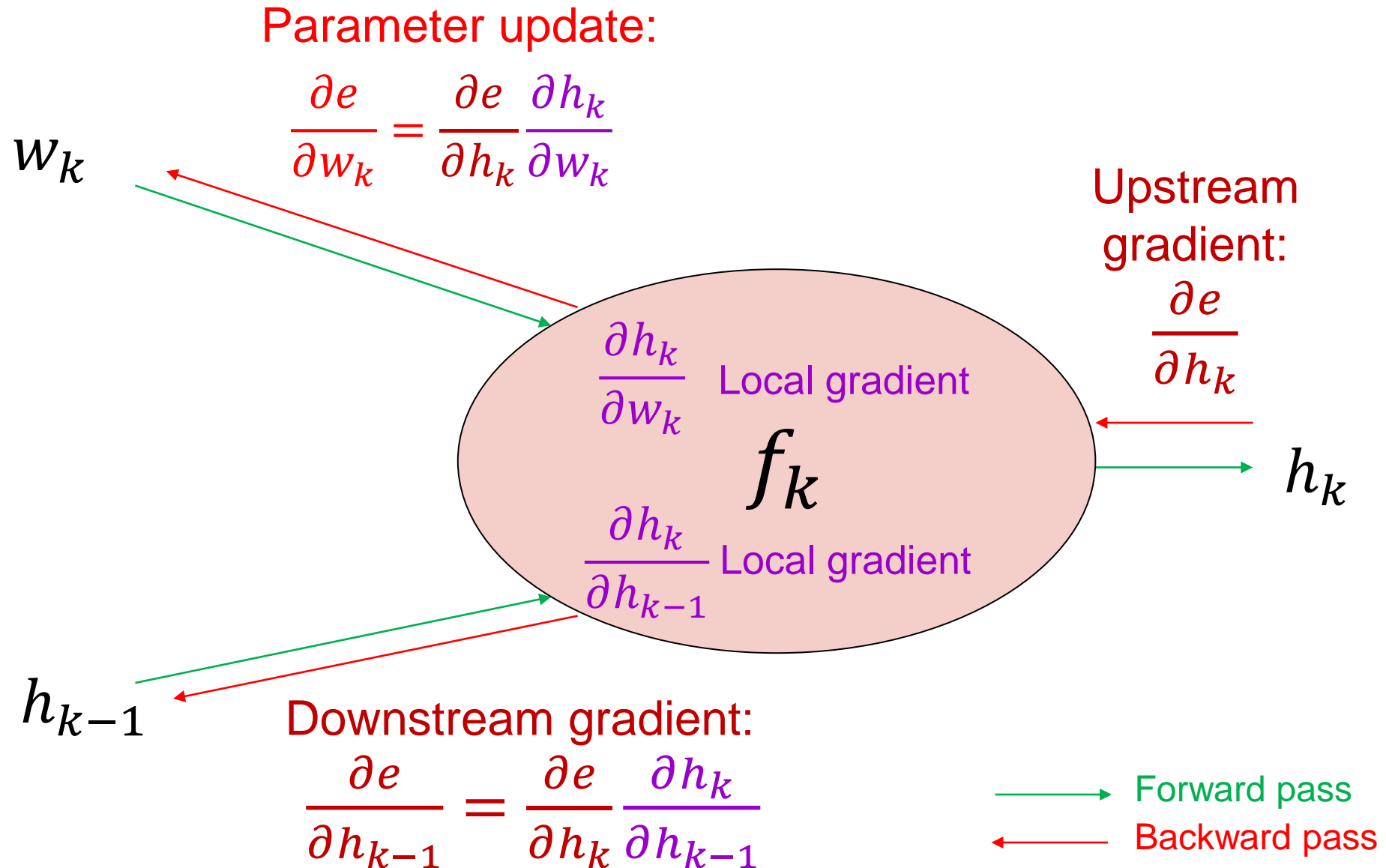
BACKPROPAGATION SUMMARY



BACKPROPAGATION SUMMARY

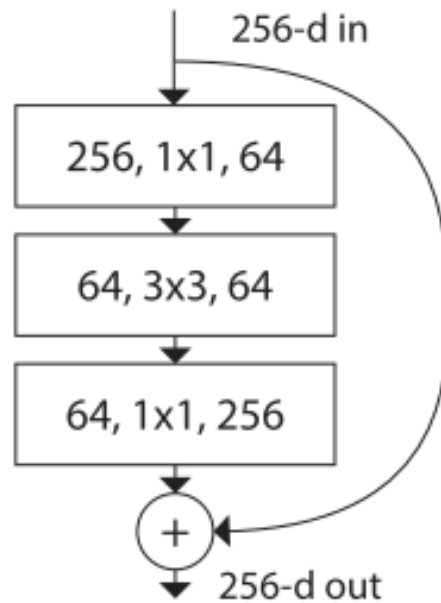


BACKPROPAGATION SUMMARY

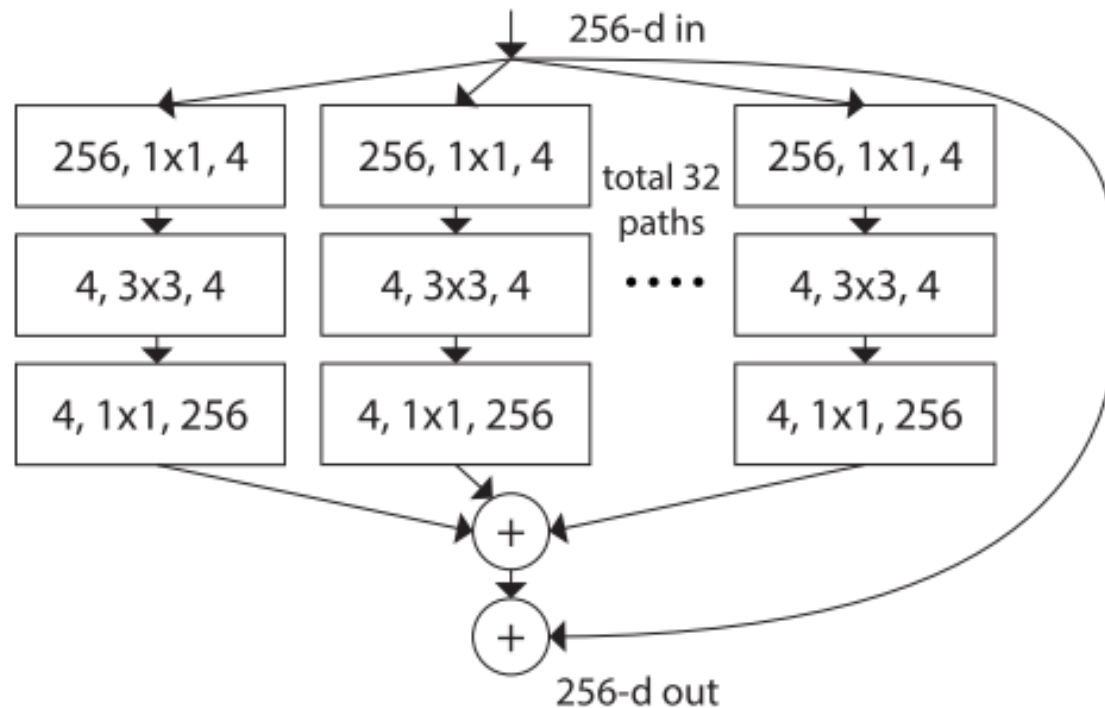


WHAT ABOUT MORE GENERAL COMPUTATION GRAPHS?

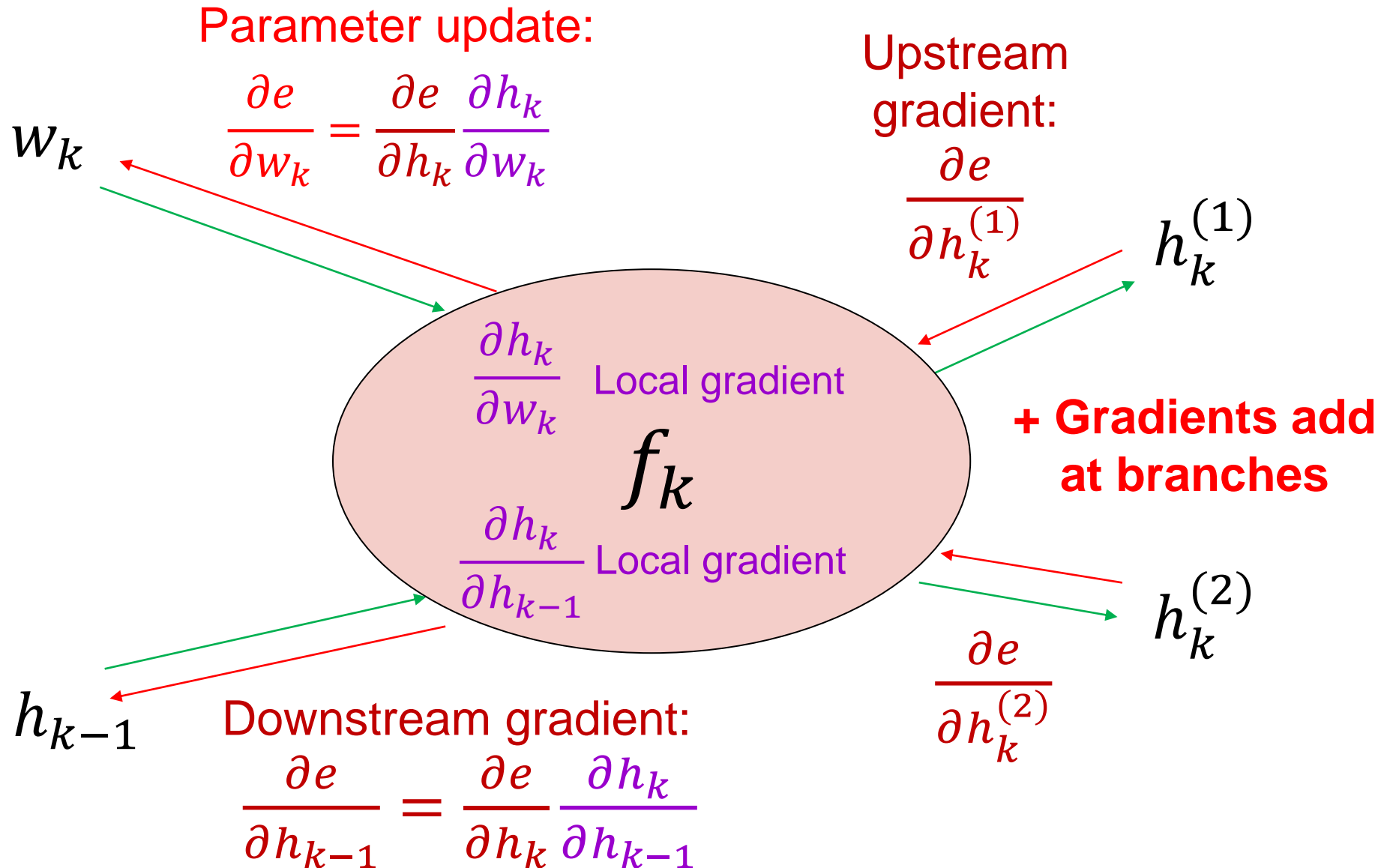
ResNet



ResNeXt



WHAT ABOUT MORE GENERAL COMPUTATION GRAPHS?



TRAINING OF MULTI-LAYER NETWORKS

- Find network weights to minimize the error between true and estimated outputs of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N \left(y_j - f_{\mathbf{w}}(\mathbf{x}_j) \right)^2$$

TRAINING OF MULTI-LAYER NETWORKS

- Find network weights to minimize the error between true and estimated outputs of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N \left(y_j - f_{\mathbf{w}}(\mathbf{x}_j) \right)^2$$

- Update weights by **gradient descent**: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$

TRAINING OF MULTI-LAYER NETWORKS

- Find network weights to minimize the error between true and estimated outputs of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N \left(y_j - f_{\mathbf{w}}(\mathbf{x}_j) \right)^2$$

- Update weights by **gradient descent**: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$
- **Back-propagation**: gradients are computed in the direction from output to input layers and combined using chain rule

NEURAL NETWORKS: PROS AND CONS

Pros

- Flexible and general function approximation framework
- Can build extremely powerful models by adding more layers

Cons

- Hard to analyze theoretically (e.g., training is prone to local optima)
- Huge amount of training data, computing power may be required to get good performance
- The space of implementation choices are huge (network architectures, parameters)

ACKNOWLEDGEMENT

Thanks to the following courses and corresponding researchers for making their teaching/research material online

- Deep Learning, Stanford University
- Introduction to Deep Learning, University of Illinois at Urbana-Champaign
- Introduction to Deep Learning, Carnegie Mellon University
- Convolutional Neural Networks for Visual Recognition, Stanford University
- Natural Language Processing with Deep Learning, Stanford University
- And Many More

Thank You