

# ENM 360: Introduction to Data-driven Modeling

## *Lecture #14: Neural networks*

Paris Perdikaris  
October 15, 2020



# Feed-forward neural networks

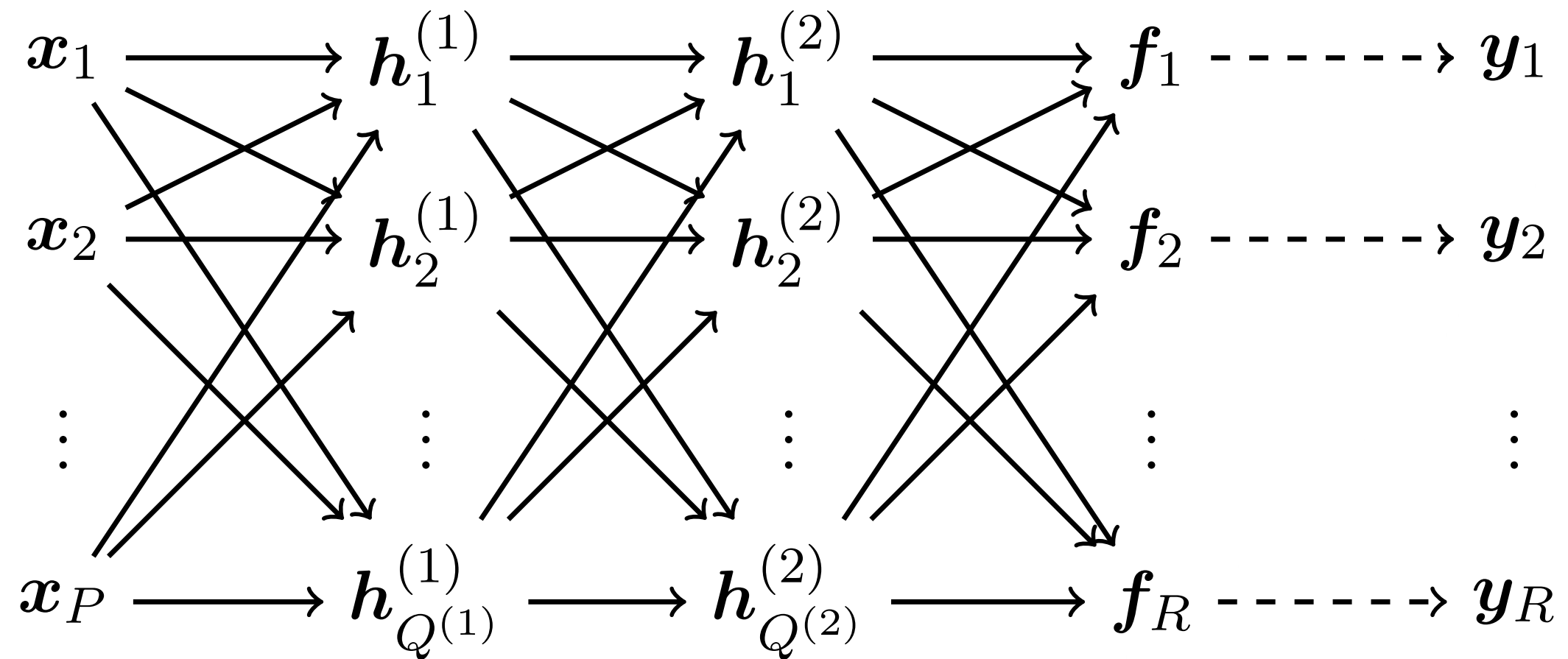
## **Pros:**

- Adaptive features/basis functions (parametric)
- Flexible non-linear regression models that can approximate any function.
- Scalability to high dimensions.

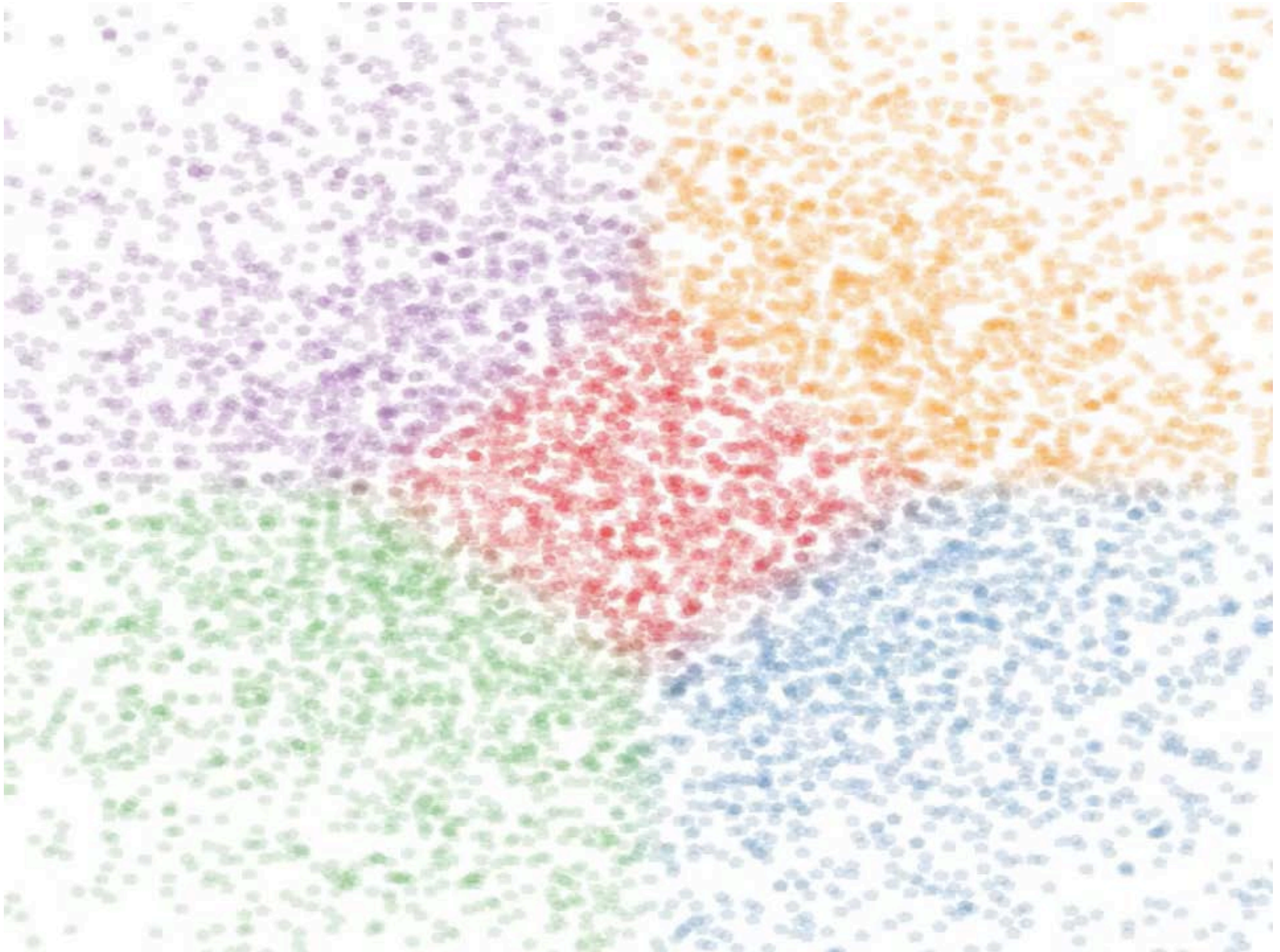
## **Cons:**

- The likelihood function is no longer a convex function of the model parameters.
- Over-fitting in data-scarce scenarios.
- Results are hard to interpret.

# Feed-forward neural networks



# Feed-forward neural networks





# Universal approximation theorem

**Theorem 1.** *Let  $\sigma$  be any continuous discriminatory function. Then finite sums of the form*

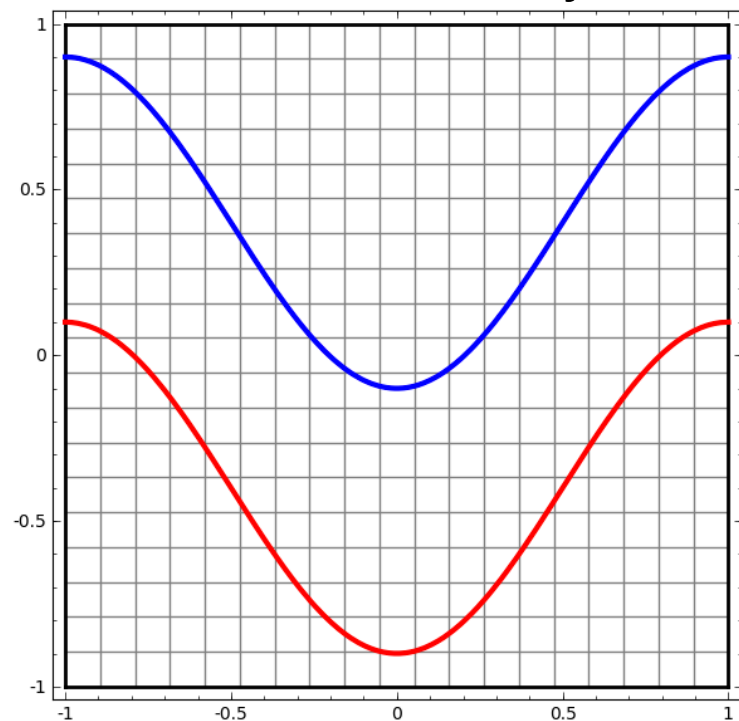
$$G(x) = \sum_{j=1}^N \alpha_j \sigma(y_j^T x + \theta_j) \quad (2)$$

*are dense in  $C(I_n)$ . In other words, given any  $f \in C(I_n)$  and  $\varepsilon > 0$ , there is a sum,  $G(x)$ , of the above form, for which*

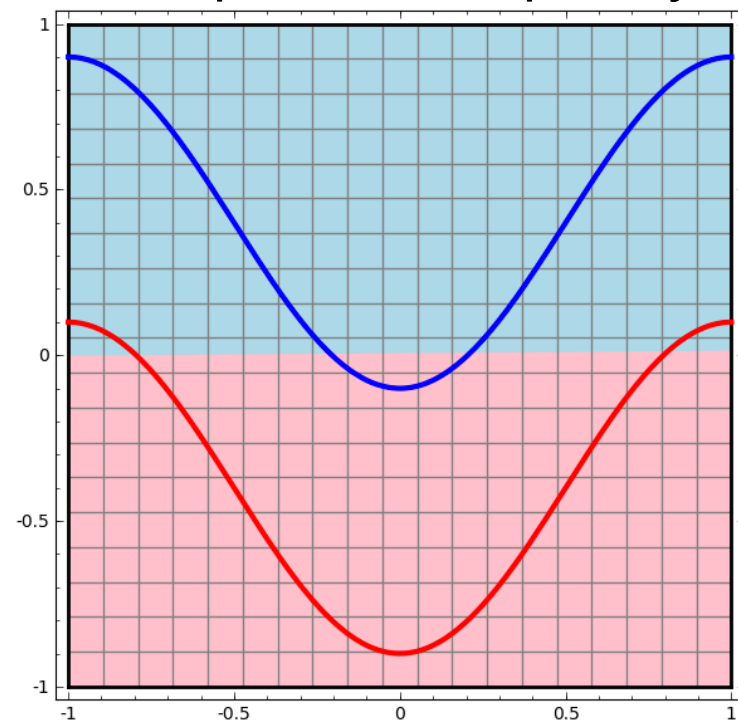
$$|G(x) - f(x)| < \varepsilon \quad \text{for all } x \in I_n.$$

# Some intuition

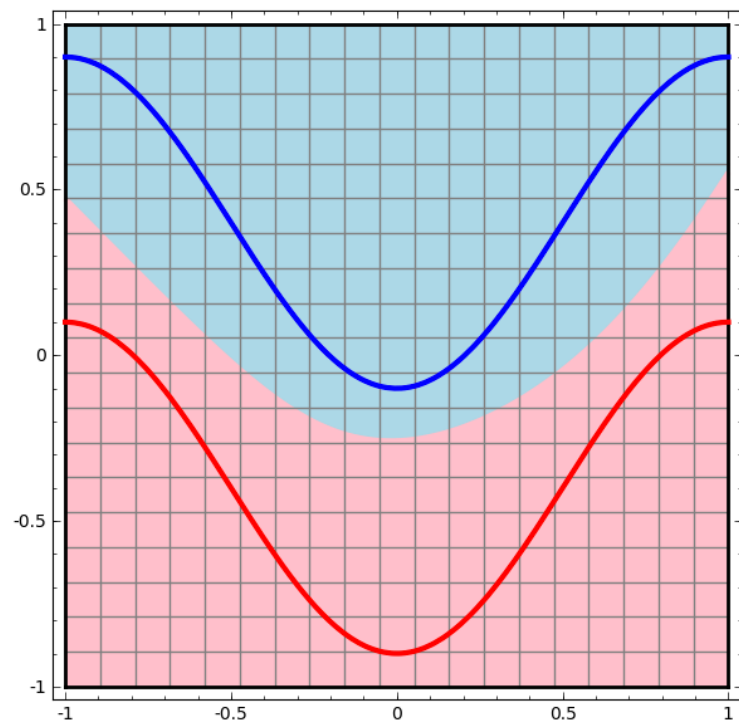
Data to classify



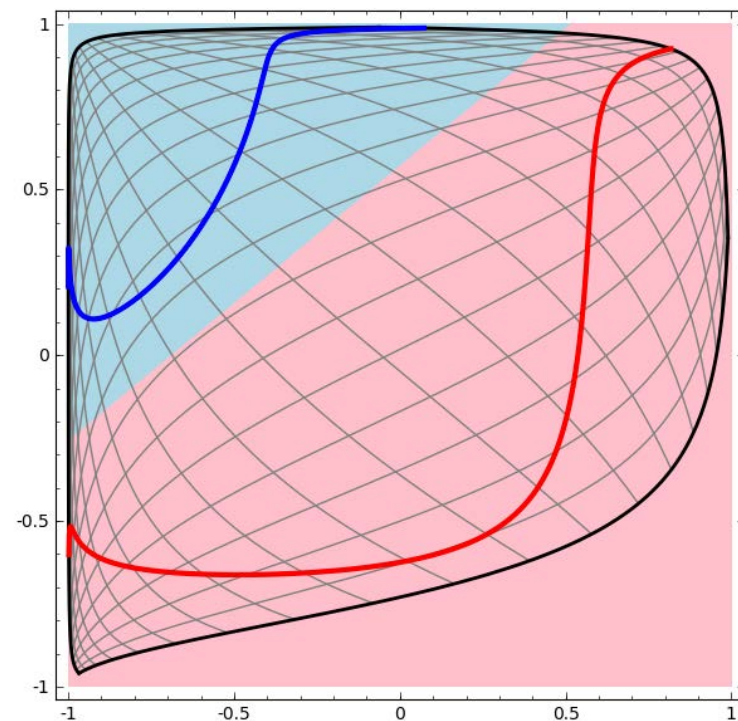
One input, one output layer



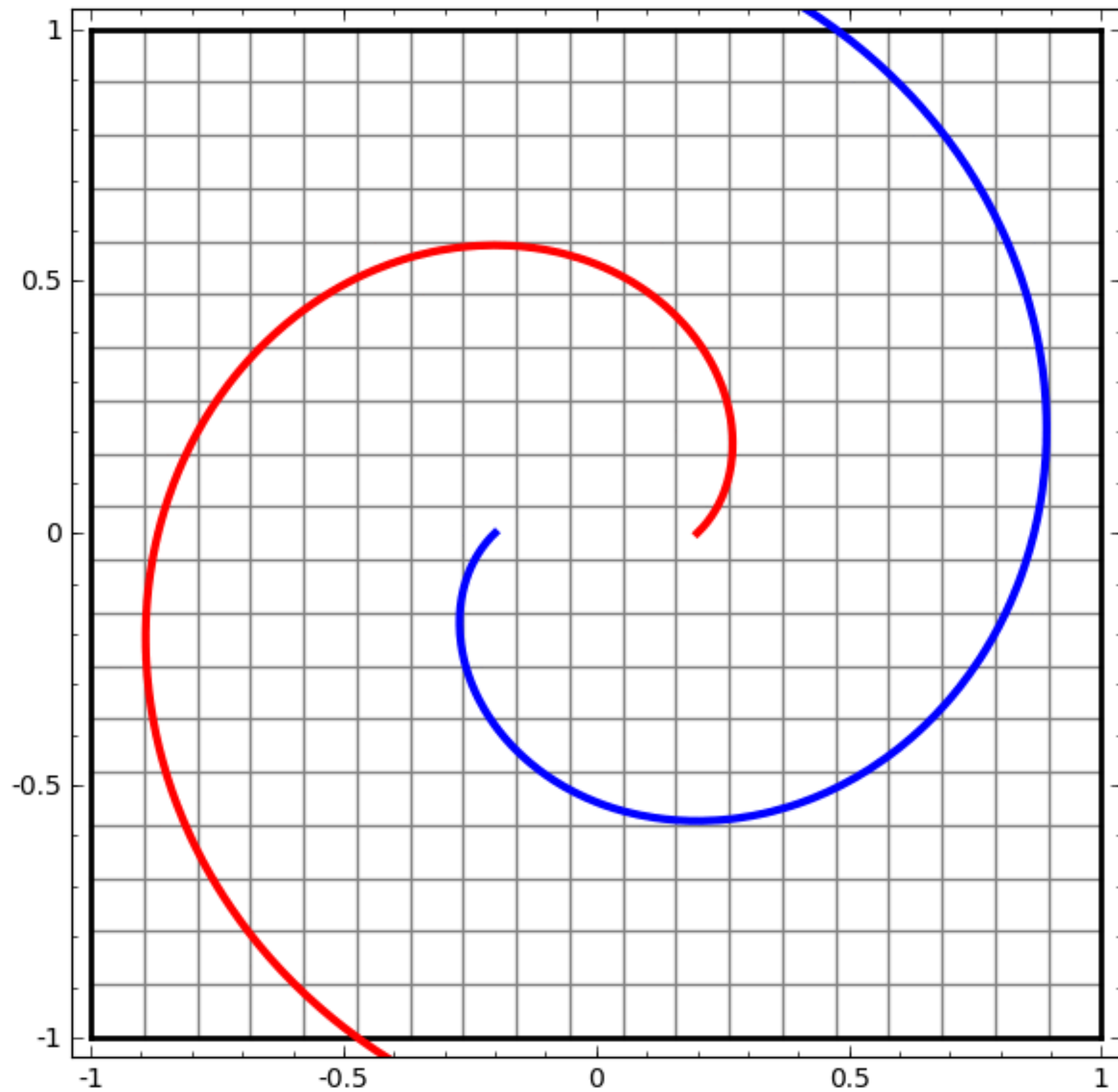
One input, one hidden, one output layer



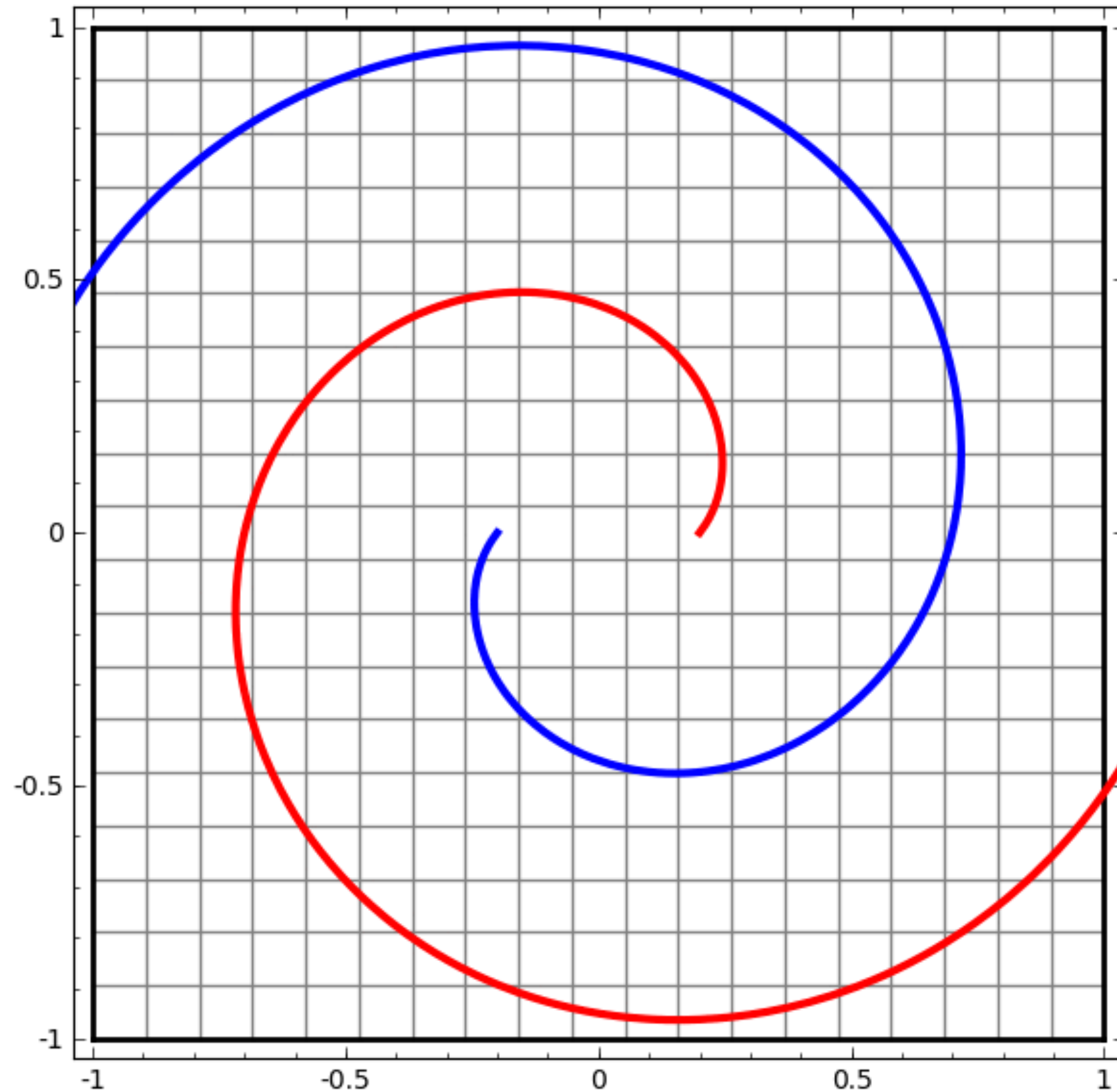
Visualizing the hidden layer



## Some intuition



## Some intuition





## Some intuition

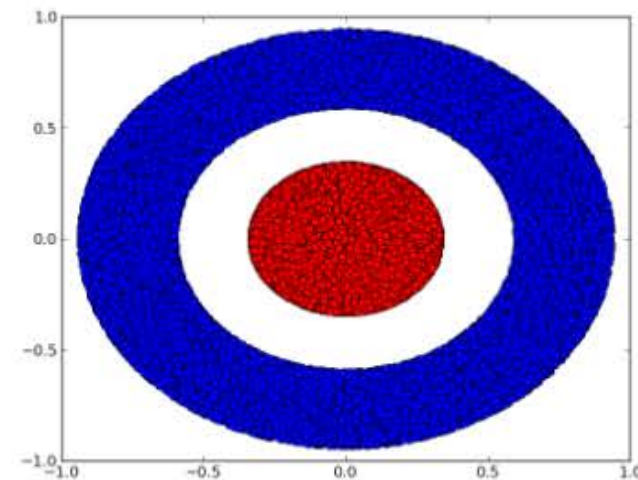
# Topology and Classification

Consider a two dimensional dataset with two classes  $A, B \subset \mathbb{R}^2$ :

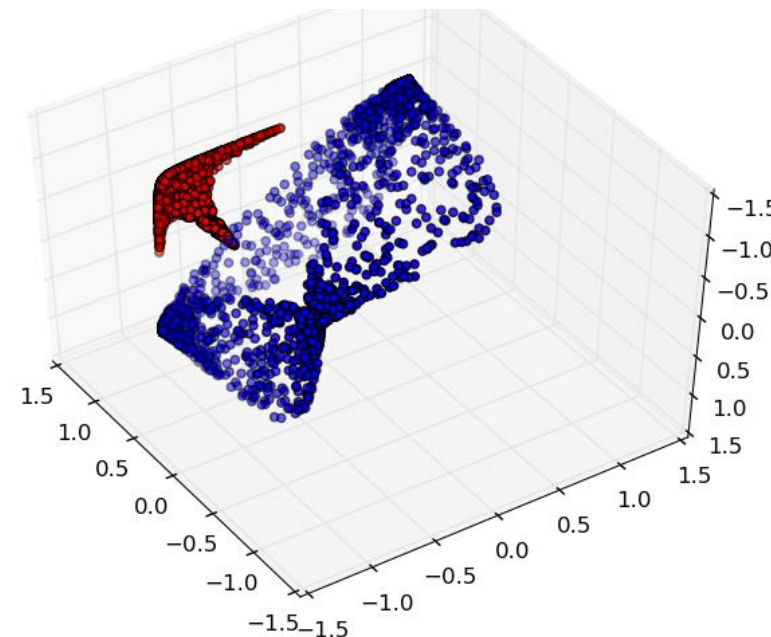
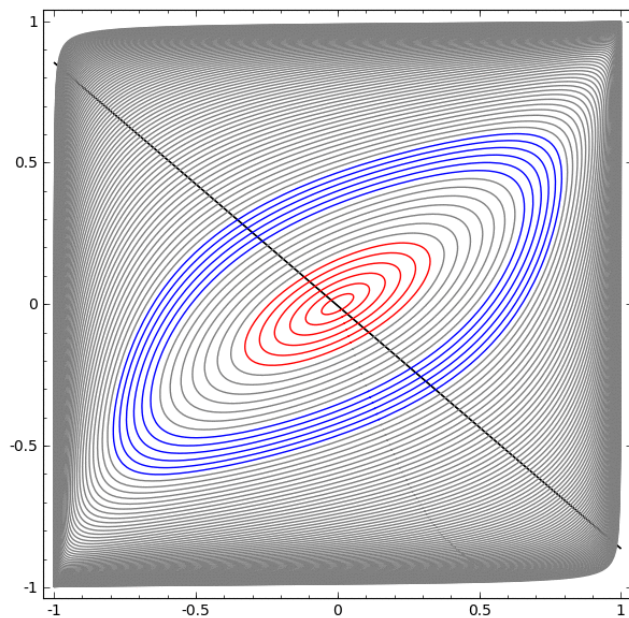
$$A = \{x | d(x, 0) < 1/3\}$$

$$B = \{x | 2/3 < d(x, 0) < 1\}$$





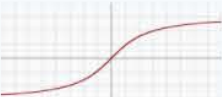
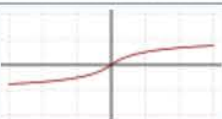




**Claim:** It is impossible for a neural network to classify this dataset without having a layer that has 3 or more hidden units, regardless of depth.



$A$  is red,  $B$  is blue



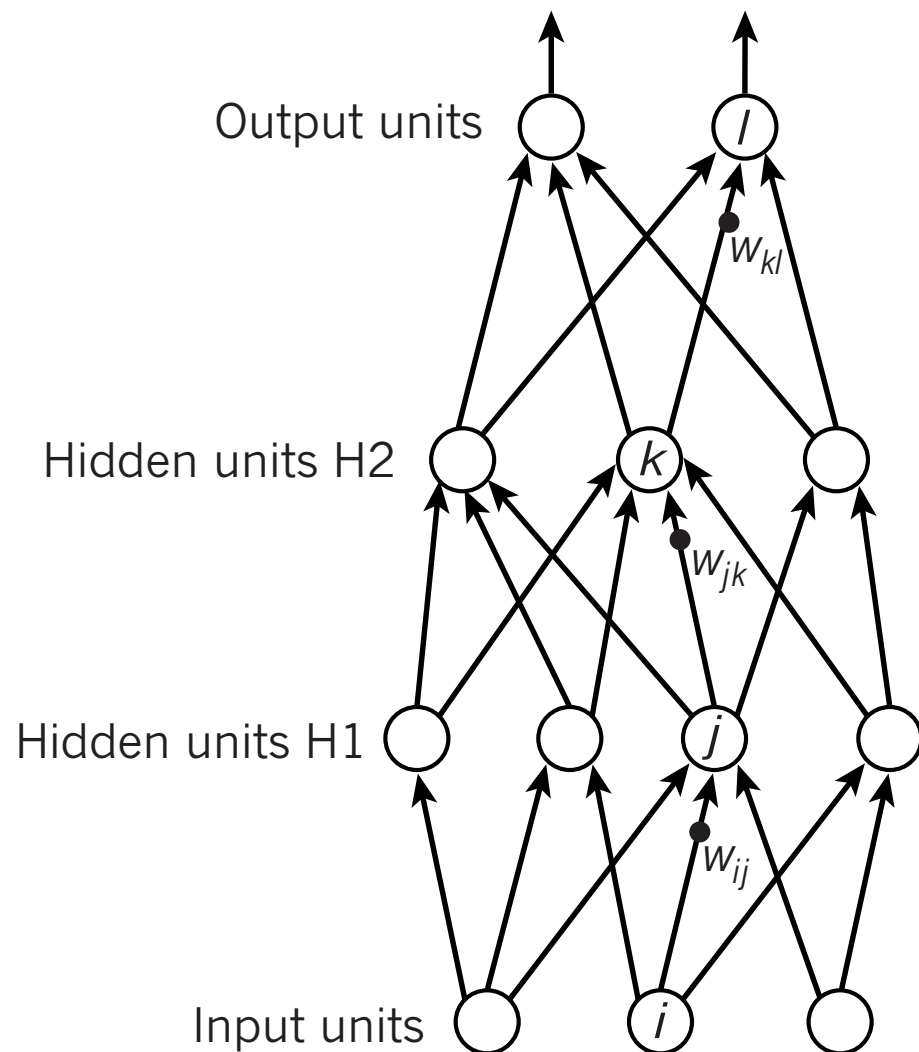
# Activation functions

Name ↕	Plot ↕	Equation ↕	Derivative (with respect to $x$ ) ↕	Range ↕	Order of continuity ↕	Monotonic ↕	Derivative Monotonic ↕	Approximates identity near the origin ↕
Identity		$f(x) = x$	$f'(x) = 1$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$	$\{0, 1\}$	$C^{-1}$	Yes	No	No
Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$	$C^\infty$	Yes	No	No
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$	$C^\infty$	Yes	No	Yes
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	$\left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$	$C^\infty$	Yes	No	Yes
Softsign <sup>[7][8]</sup>		$f(x) = \frac{x}{1 +  x }$	$f'(x) = \frac{1}{(1 +  x )^2}$	$(-1, 1)$	$C^1$	Yes	No	Yes
Inverse square root unit (ISRU) <sup>[9]</sup>		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$	$f'(x) = \left(\frac{1}{\sqrt{1 + \alpha x^2}}\right)^3$	$\left(-\frac{1}{\sqrt{\alpha}}, \frac{1}{\sqrt{\alpha}}\right)$	$C^\infty$	Yes	No	Yes
Rectified linear unit (ReLU) <sup>[10]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$	$C^0$	Yes	Yes	No
Leaky rectified linear unit (Leaky ReLU) <sup>[11]</sup>		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Parametric rectified linear unit (PReLU) <sup>[12]</sup>		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes iff $\alpha \geq 0$	Yes	Yes iff $\alpha = 1$



# Back-propagation

## Forward pass



$$y_l = f(z_l)$$

$$z_l = \sum_{k \in H2} w_{kl} y_k$$

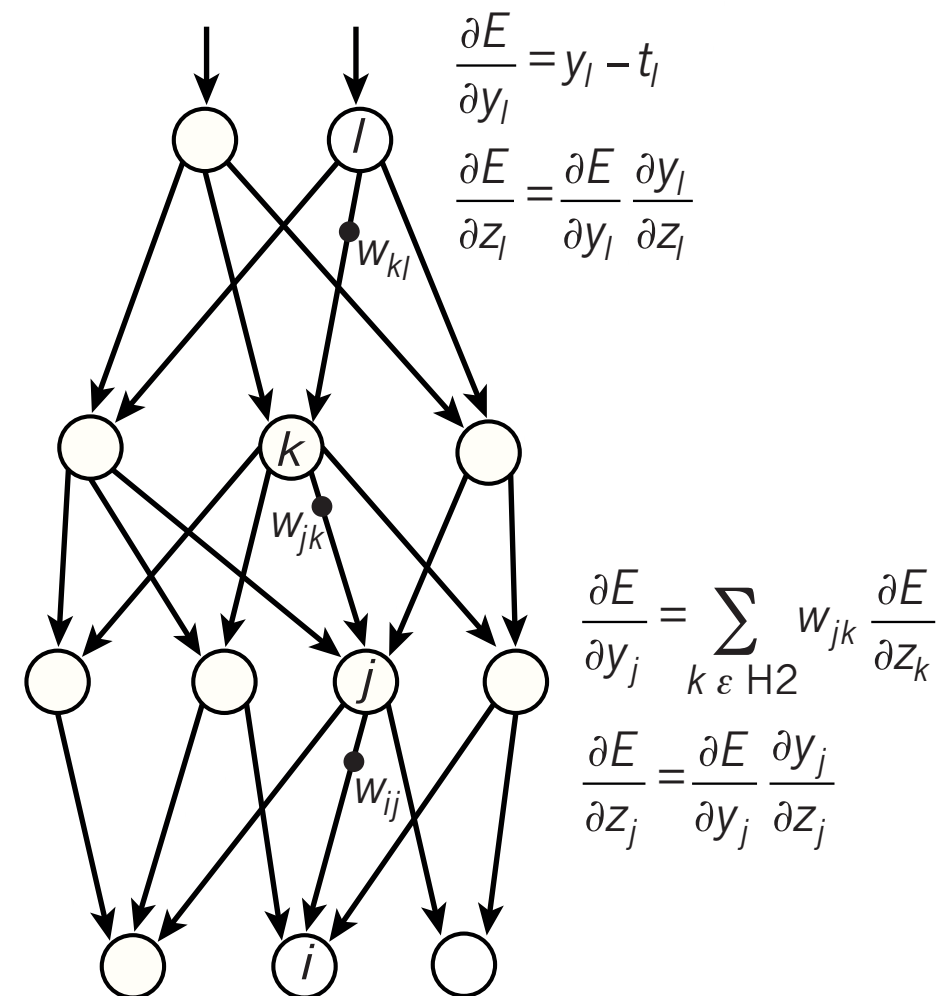
$$y_k = f(z_k)$$

$$z_k = \sum_{j \in H1} w_{jk} y_j$$

$$y_j = f(z_j)$$

$$z_j = \sum_{i \in \text{Input}} w_{ij} x_i$$

## Backward pass



$$\frac{\partial E}{\partial y_l} = y_l - t_l$$

$$\frac{\partial E}{\partial z_l} = \frac{\partial E}{\partial y_l} \frac{\partial y_l}{\partial z_l}$$

$$\frac{\partial E}{\partial y_k} = \sum_{l \in \text{out}} w_{kl} \frac{\partial E}{\partial z_l}$$

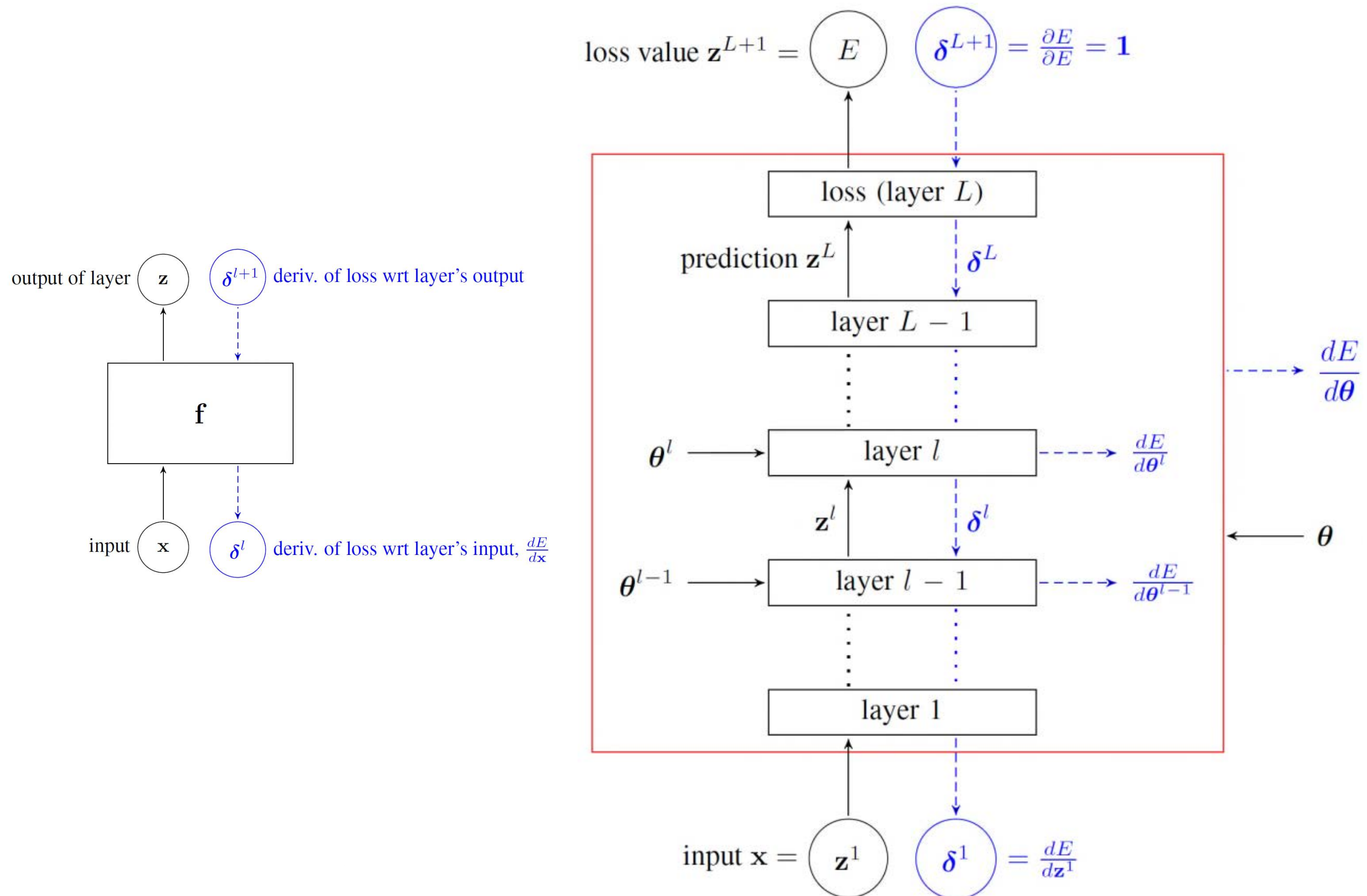
$$\frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k}$$

$$\frac{\partial E}{\partial y_j} = \sum_{k \in H2} w_{jk} \frac{\partial E}{\partial z_k}$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_j}$$



# Back-propagation



## 6.5.2 Chain Rule of Calculus

The chain rule of calculus (not to be confused with the chain rule of probability) is used to compute the derivatives of functions formed by composing other functions whose derivatives are known. Back-propagation is an algorithm that computes the chain rule, with a specific order of operations that is highly efficient.

Let  $x$  be a real number, and let  $f$  and  $g$  both be functions mapping from a real number to a real number. Suppose that  $y = g(x)$  and  $z = f(g(x)) = f(y)$ . Then the chain rule states that

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}. \quad (6.44)$$

We can generalize this beyond the scalar case. Suppose that  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $g$  maps from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ , and  $f$  maps from  $\mathbb{R}^n$  to  $\mathbb{R}$ . If  $\mathbf{y} = g(\mathbf{x})$  and  $z = f(\mathbf{y})$ , then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}. \quad (6.45)$$

In vector notation, this may be equivalently written as

$$\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z, \quad (6.46)$$

where  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$  is the  $n \times m$  Jacobian matrix of  $g$ .

From this we see that the gradient of a variable  $\mathbf{x}$  can be obtained by multiplying a Jacobian matrix  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$  by a gradient  $\nabla_{\mathbf{y}} z$ . The back-propagation algorithm consists of performing such a Jacobian-gradient product for each operation in the graph.



# Automatic differentiation

Many contemporary algorithms require the evaluation of a derivative of a given differentiable function,  $f$ , at a given input value,  $(x_1, \dots, x_N)$ , for example a gradient,

$$\left( \frac{\partial f}{\partial x_1} (x_1, \dots, x_N), \dots, \frac{\partial f}{\partial x_N} (x_1, \dots, x_N) \right),$$

or a directional derivative,<sup>1</sup>

$$\vec{v}(f) (x_1, \dots, x_N) = \sum_{n=1}^N v_n \frac{\partial f}{\partial x_n} (x_1, \dots, x_N).$$

In its most basic description, automatic differentiation relies on the fact that all numerical computations are ultimately compositions of a finite set of elementary operations for which derivatives are known. Combining the derivatives of the constituent operations through the chain rule gives the derivative of the overall composition. This allows accurate evaluation of derivatives at machine precision with ideal asymptotic efficiency and only a small constant factor of overhead.

# Automatic differentiation

## The chain rule, forward and reverse accumulation [\[edit\]](#)

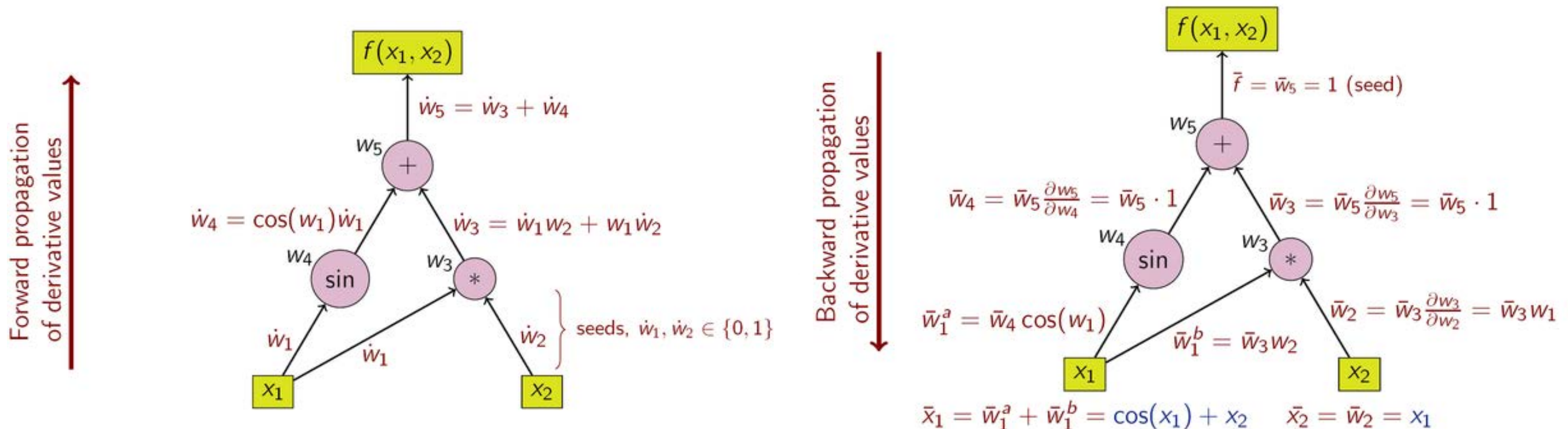
Fundamental to AD is the decomposition of differentials provided by the [chain rule](#). For the simple composition  $y = f(g(h(x))) = f(g(h(w_0))) = f(g(w_1)) = f(w_2) = w_3$  the chain rule gives

$$\frac{dy}{dx} = \frac{dy}{dw_2} \frac{dw_2}{dw_1} \frac{dw_1}{dx}$$

Usually, two distinct modes of AD are presented, **forward accumulation** (or **forward mode**) and **reverse accumulation** (or **reverse mode**). Forward accumulation specifies that one traverses the chain rule from inside to outside (that is, first compute  $dw_1/dx$  and then  $dw_2/dx$  and at last  $dy/dx$ ), while reverse accumulation has the traversal from outside to inside (first compute  $dy/dw_2$  and then  $dy/dw_1$  and at last  $dy/dx$ ). More succinctly,

1. **forward accumulation** computes the recursive relation:  $\frac{dw_i}{dx} = \frac{dw_i}{dw_{i-1}} \frac{dw_{i-1}}{dx}$  with  $w_3 = y$ , and,
2. **reverse accumulation** computes the recursive relation:  $\frac{dy}{dw_i} = \frac{dy}{dw_{i+1}} \frac{dw_{i+1}}{dw_i}$  with  $w_0 = x$ .

**Example**  $z = f(x_1, x_2) = x_1 x_2 + \sin x_1$



# Automatic differentiation

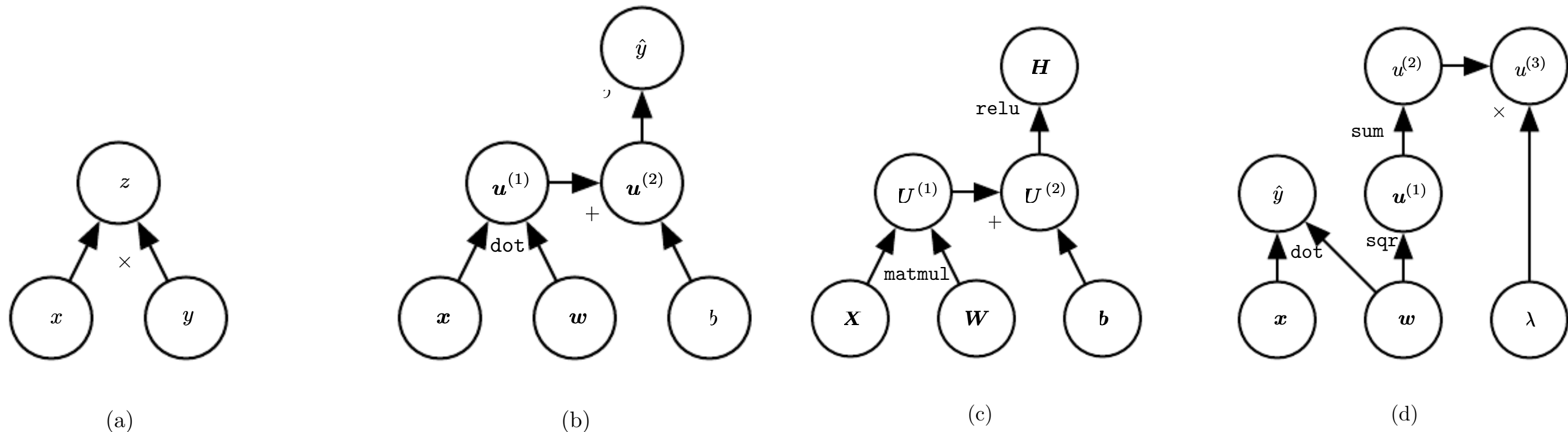


Figure 6.8: Examples of computational graphs. (a) The graph using the  $\times$  operation to compute  $z = xy$ . (b) The graph for the logistic regression prediction  $\hat{y} = \sigma(\mathbf{x}^\top \mathbf{w} + b)$ . Some of the intermediate expressions do not have names in the algebraic expression but need names in the graph. We simply name the  $i$ -th such variable  $\mathbf{u}^{(i)}$ . (c) The computational graph for the expression  $\mathbf{H} = \max\{0, \mathbf{XW} + \mathbf{b}\}$ , which computes a design matrix of rectified linear unit activations  $\mathbf{H}$  given a design matrix containing a minibatch of inputs  $\mathbf{X}$ . (d) Examples a–c applied at most one operation to each variable, but it is possible to apply more than one operation. Here we show a computation graph that applies more than one operation to the weights  $\mathbf{w}$  of a linear regression model. The weights are used to make both the prediction  $\hat{y}$  and the weight decay penalty  $\lambda \sum_i w_i^2$ .

# Automatic differentiation

It is one of the most useful - and perhaps underused - tools in modern scientific computing!

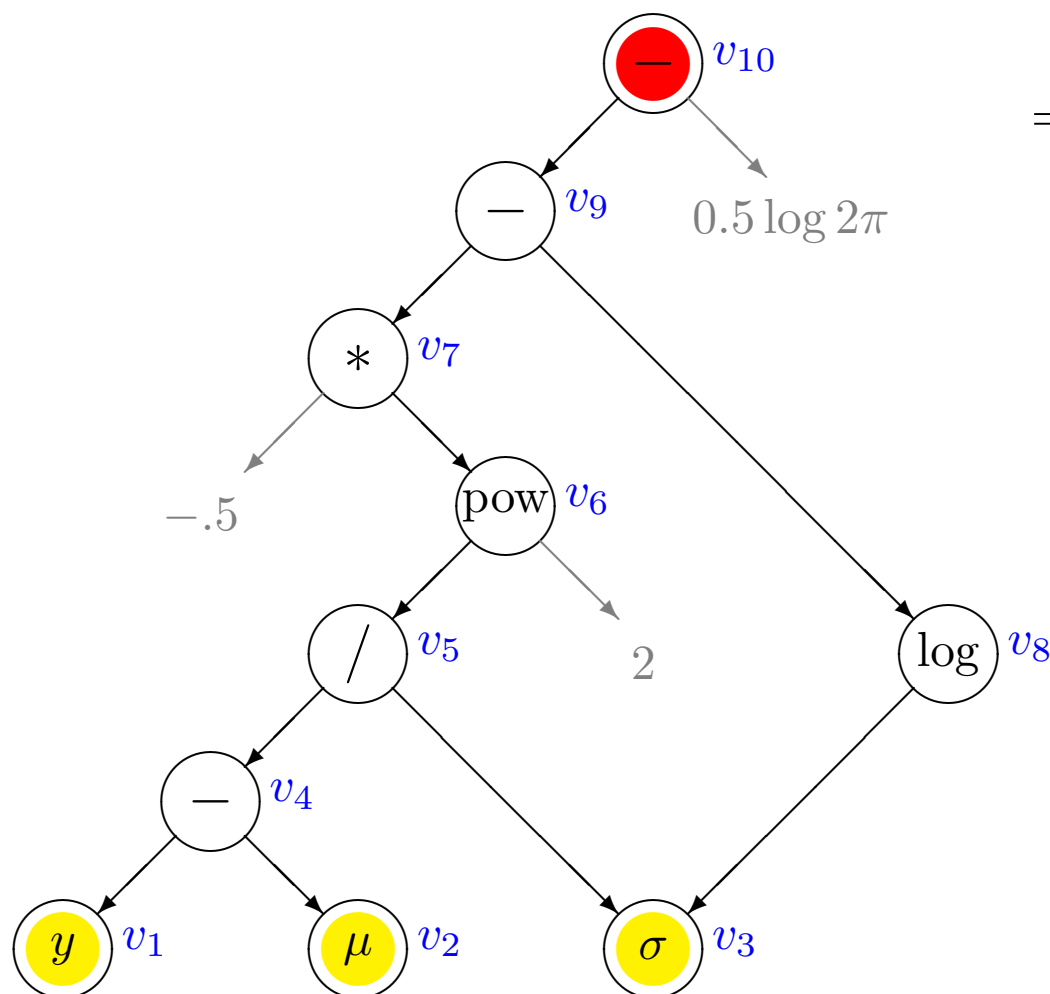
## **Applications:**

- real-parameter optimization (many good methods are gradient-based)
- sensitivity analysis (local sensitivity =  $\partial(\text{result})/\partial(\text{input})$ )
- physical modeling (forces are derivatives of potentials; equations of motion are derivatives of Lagrangians and Hamiltonians; etc.)
- probabilistic inference (e.g., Hamiltonian Monte Carlo)
- machine learning
- and who knows how many other scientific computing applications.

# Automatic differentiation

As an example, consider the log of the normal probability density function for a variable  $y$  with a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ ,

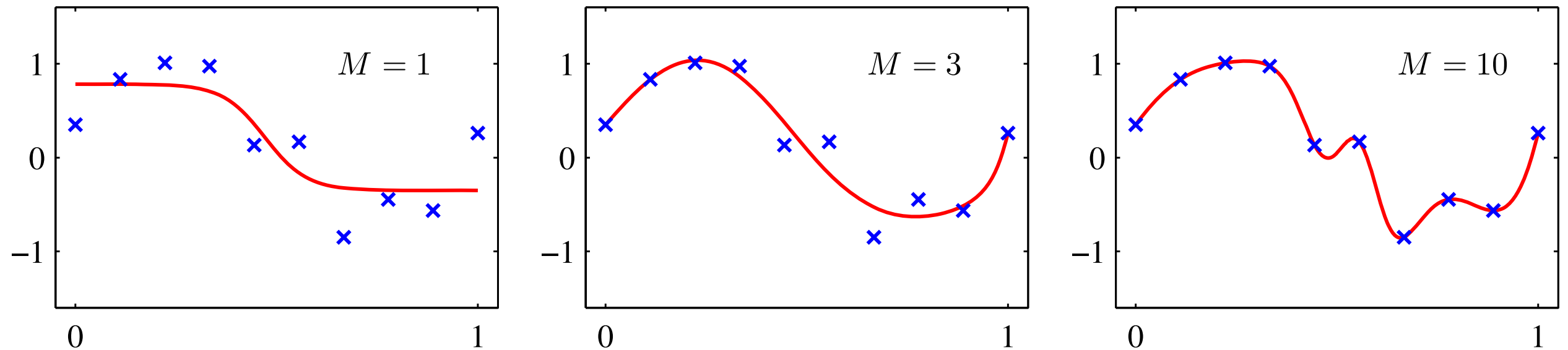
$$f(y, \mu, \sigma) = \log(\text{Normal}(y|\mu, \sigma)) = -\frac{1}{2} \left( \frac{y - \mu}{\sigma} \right)^2 - \log \sigma - \frac{1}{2} \log(2\pi) \quad (1)$$



<i>var</i>	<i>value</i>	<i>partials</i>
$v_1$	$y$	
$v_2$	$\mu$	
$v_3$	$\sigma$	
$v_4$	$v_1 - v_2$	$\partial v_4 / \partial v_1 = 1 \quad \partial v_4 / \partial v_2 = -1$
$v_5$	$v_4 / v_3$	$\partial v_5 / \partial v_4 = 1 / v_3 \quad \partial v_5 / \partial v_3 = -v_4 v_3^{-2}$
$v_6$	$(v_5)^2$	$\partial v_6 / \partial v_5 = 2v_5$
$v_7$	$(-0.5)v_6$	$\partial v_7 / \partial v_6 = -0.5$
$v_8$	$\log v_3$	$\partial v_8 / \partial v_3 = 1 / v_3$
$v_9$	$v_7 - v_8$	$\partial v_9 / \partial v_7 = 1 \quad \partial v_9 / \partial v_8 = -1$
$v_{10}$	$v_9 - (0.5 \log 2\pi)$	$\partial v_{10} / \partial v_9 = 1$

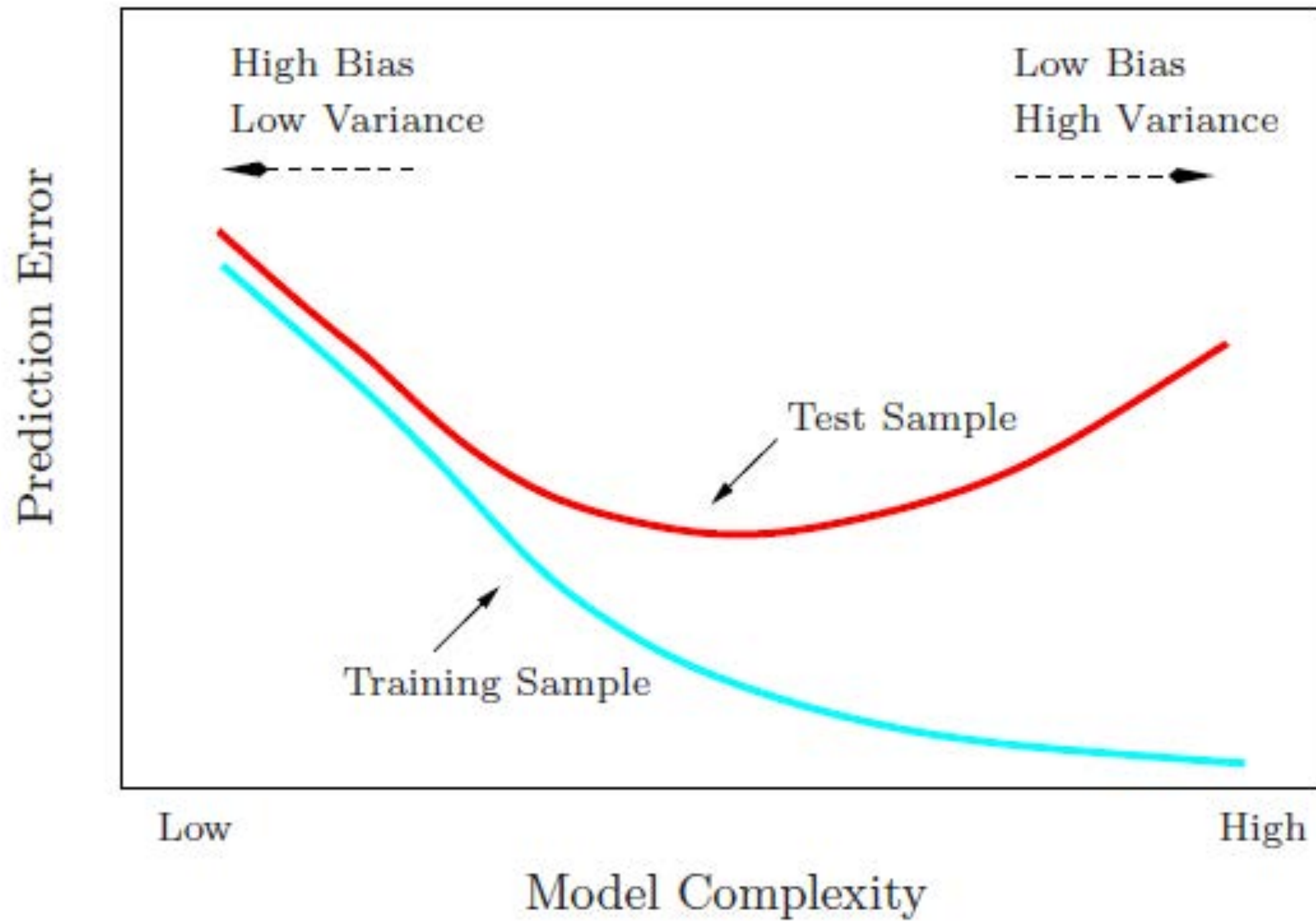


# Overfitting

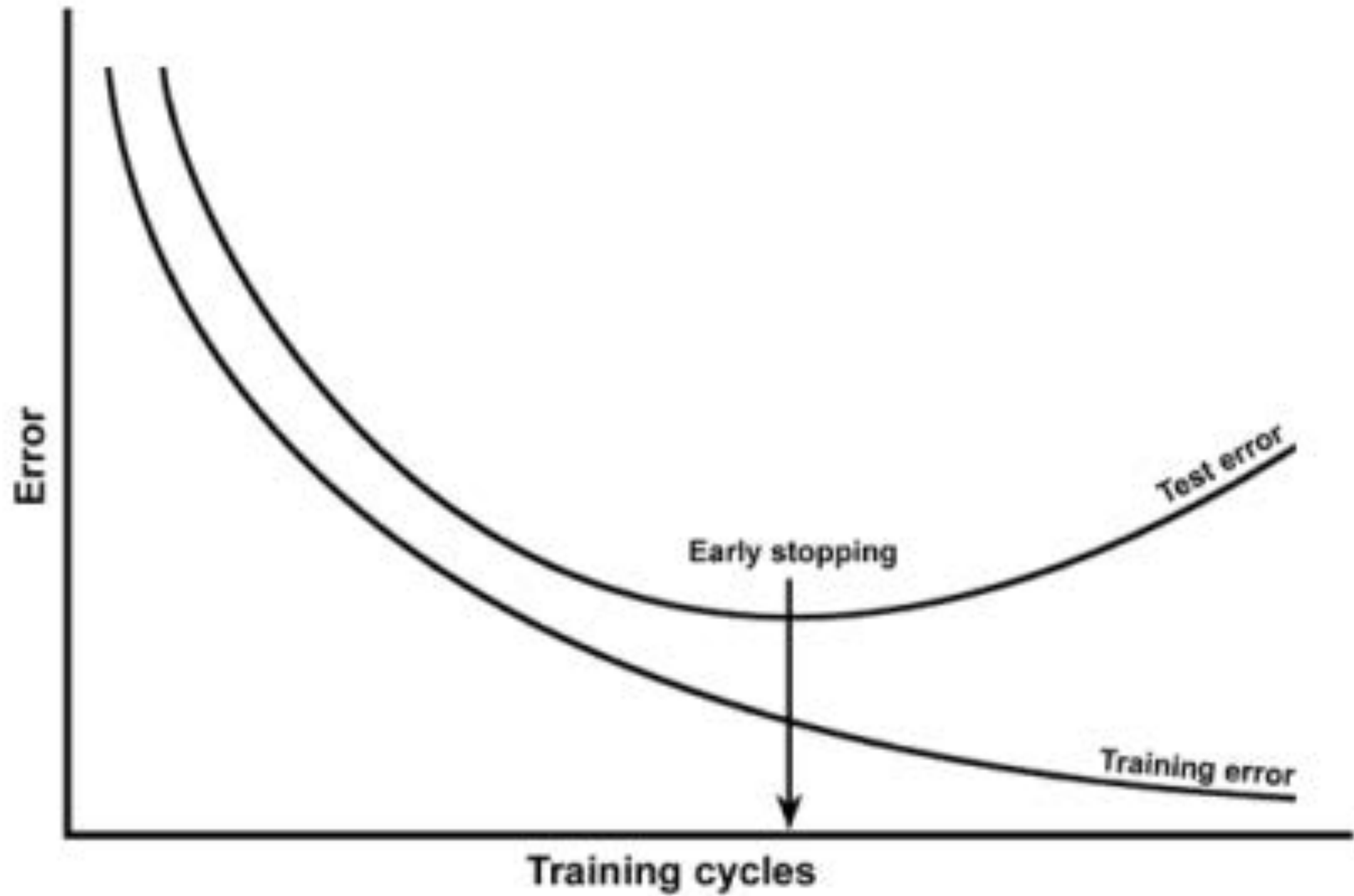


**Figure 5.9** Examples of two-layer networks trained on 10 data points drawn from the sinusoidal data set. The graphs show the result of fitting networks having  $M = 1$ , 3 and 10 hidden units, respectively, by minimizing a sum-of-squares error function using a scaled conjugate-gradient algorithm.

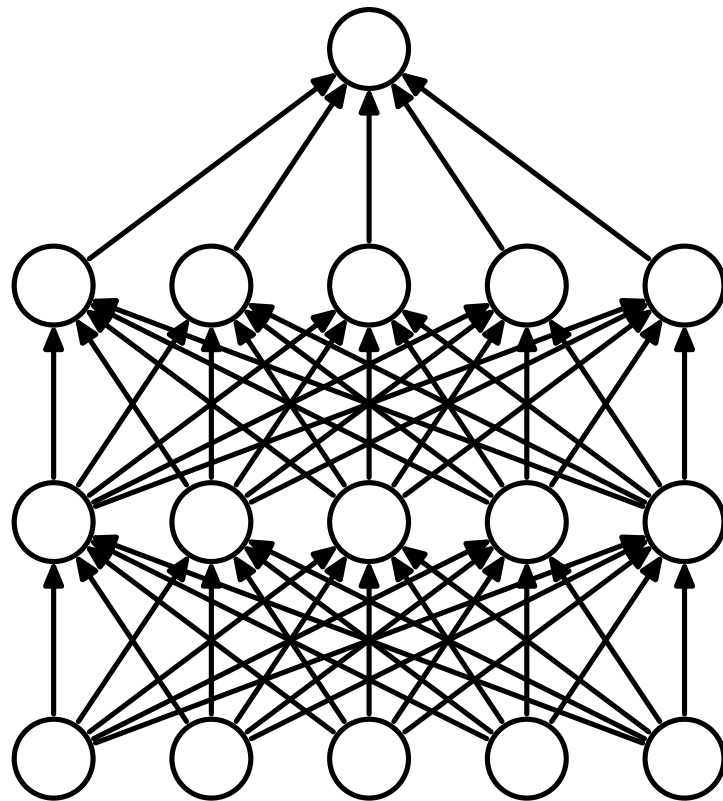
# Overfitting



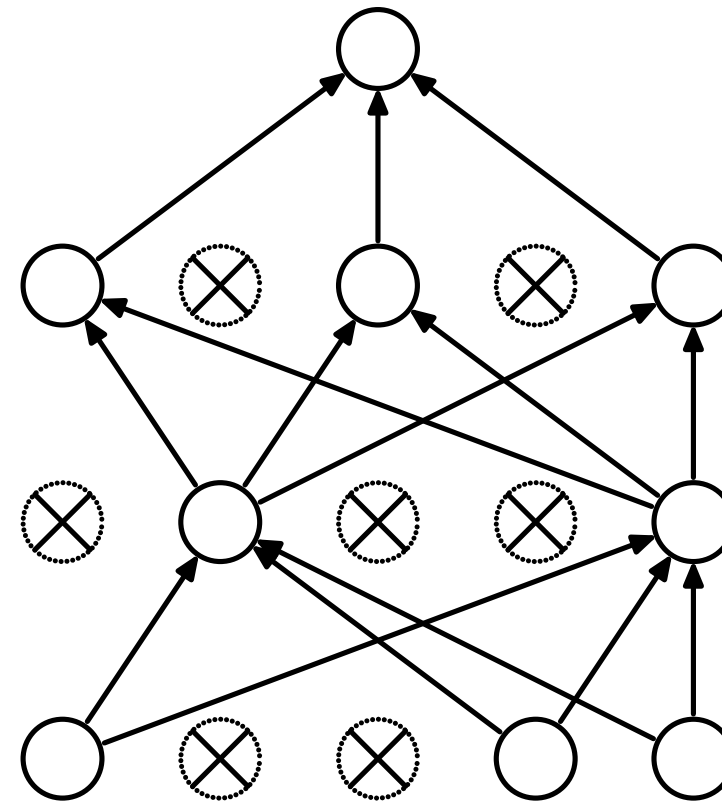
# Early stopping



# Dropout



(a) Standard Neural Net



(b) After applying dropout.

With probability `keep_prob`, outputs the input element scaled up by  $1 / \text{keep\_prob}$ , otherwise outputs 0. The scaling is so that the expected sum is unchanged.

```
for W, b in params:
    outputs = np.dot(inputs, W) + b
    inputs = np.tanh(outputs)
    if dropout_train: inputs *= np.random.binomial([np.ones_like(inputs)], (1-
keep_prob))[0]/(1-keep_prob)
```

*Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research, 15(1), 1929-1958.*

*Gal, Y., & Ghahramani, Z. (2016, June). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In international conference on machine learning (pp. 1050-1059).*

# Dropout

