

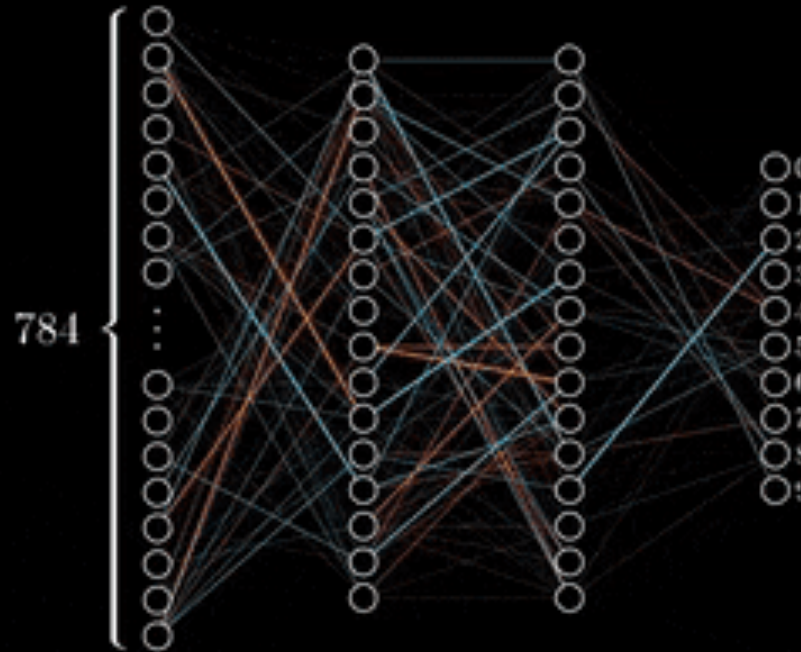
# ENM 3600: Data-driven Modeling and Probabilistic Scientific Computing

## *Lecture #13: Training MLP networks*



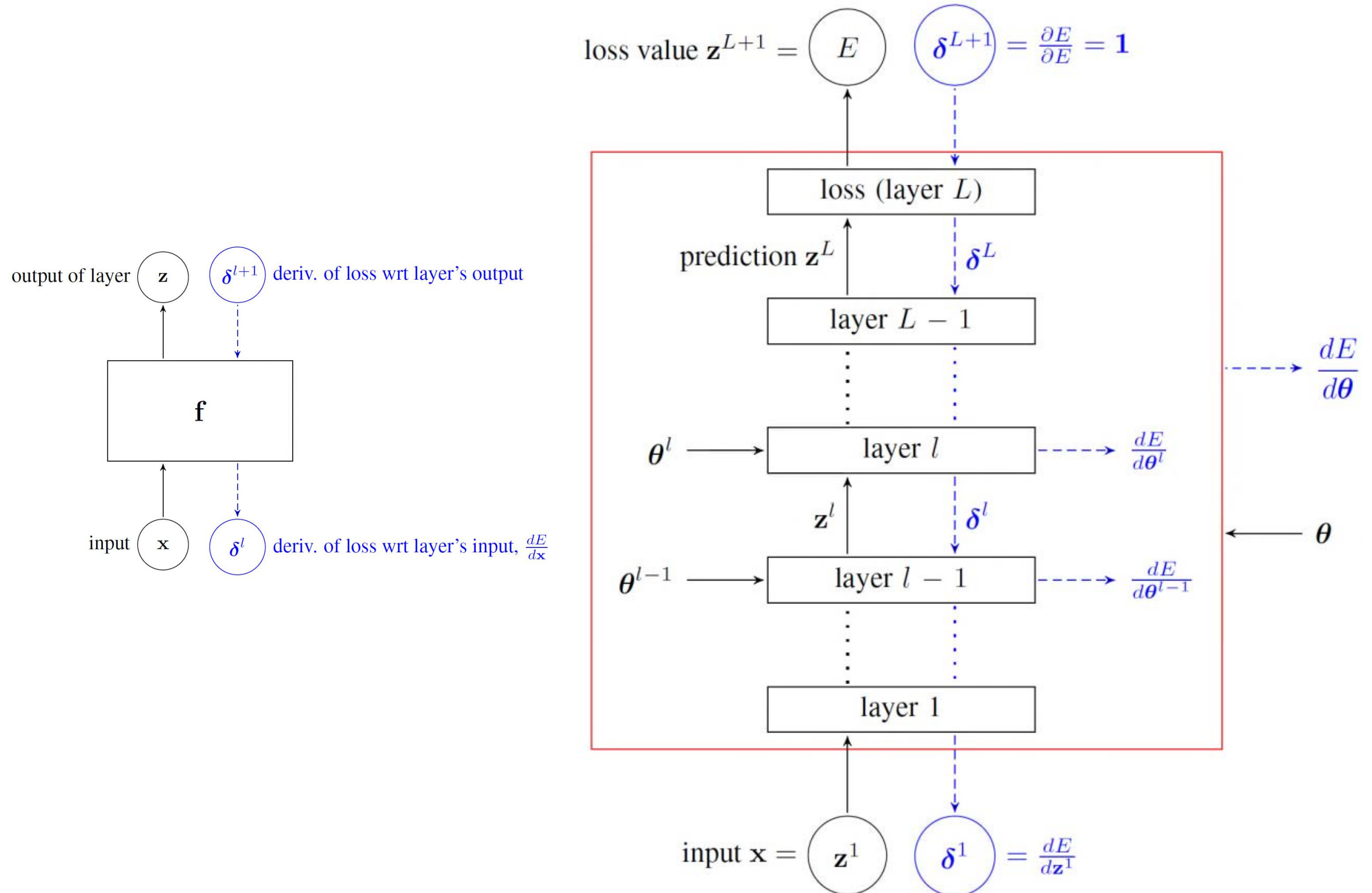
# Backpropagation

Training in progress...



$$\begin{aligned}
 G_\ell &= F - Y \in \mathbb{R}^{n \times p_{\ell+1}}, & \frac{\partial \mathcal{L}}{\partial W_\ell} &= A_\ell^T G_\ell \in \mathbb{R}^{p_\ell \times p_{\ell+1}}, & \frac{\partial \mathcal{L}}{\partial b_\ell} &= \mathbf{1}^T G_\ell \in \mathbb{R}^{1 \times p_{\ell+1}}, \\
 G_{\ell-1} &= (1 - A_\ell \odot A_\ell) \odot (G_\ell W_\ell^T) \in \mathbb{R}^{n \times p_\ell}, & \frac{\partial \mathcal{L}}{\partial W_{\ell-1}} &= A_{\ell-1}^T G_{\ell-1} \in \mathbb{R}^{p_{\ell-1} \times p_\ell}, & \frac{\partial \mathcal{L}}{\partial b_{\ell-1}} &= \mathbf{1}^T G_{\ell-1} \in \mathbb{R}^{1 \times p_\ell}, \\
 G_{\ell-2} &= (1 - A_{\ell-1} \odot A_{\ell-1}) \odot (G_{\ell-1} W_{\ell-1}^T) \in \mathbb{R}^{n \times p_{\ell-1}}, & \frac{\partial \mathcal{L}}{\partial W_{\ell-2}} &= A_{\ell-2}^T G_{\ell-2} \in \mathbb{R}^{p_{\ell-2} \times p_{\ell-1}}, & \frac{\partial \mathcal{L}}{\partial b_{\ell-2}} &= \mathbf{1}^T G_{\ell-2} \in \mathbb{R}^{1 \times p_{\ell-1}}, \\
 &\vdots & &\vdots & &\vdots \\
 G_1 &= (1 - A_2 \odot A_2) \odot (G_2 W_2^T) \in \mathbb{R}^{n \times p_2}, & \frac{\partial \mathcal{L}}{\partial W_1} &= A_1^T G_1 \in \mathbb{R}^{p_1 \times p_2}, & \frac{\partial \mathcal{L}}{\partial b_1} &= \mathbf{1}^T G_1 \in \mathbb{R}^{1 \times p_2}, \\
 G_0 &= (1 - A_1 \odot A_1) \odot (G_1 W_1^T) \in \mathbb{R}^{n \times p_1}, & \frac{\partial \mathcal{L}}{\partial W_0} &= X^T G_0 \in \mathbb{R}^{p_0 \times p_1}, & \frac{\partial \mathcal{L}}{\partial b_0} &= \mathbf{1}^T G_0 \in \mathbb{R}^{1 \times p_1}.
 \end{aligned}$$

# Modular implementation



# 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 [\[edit\]](#)

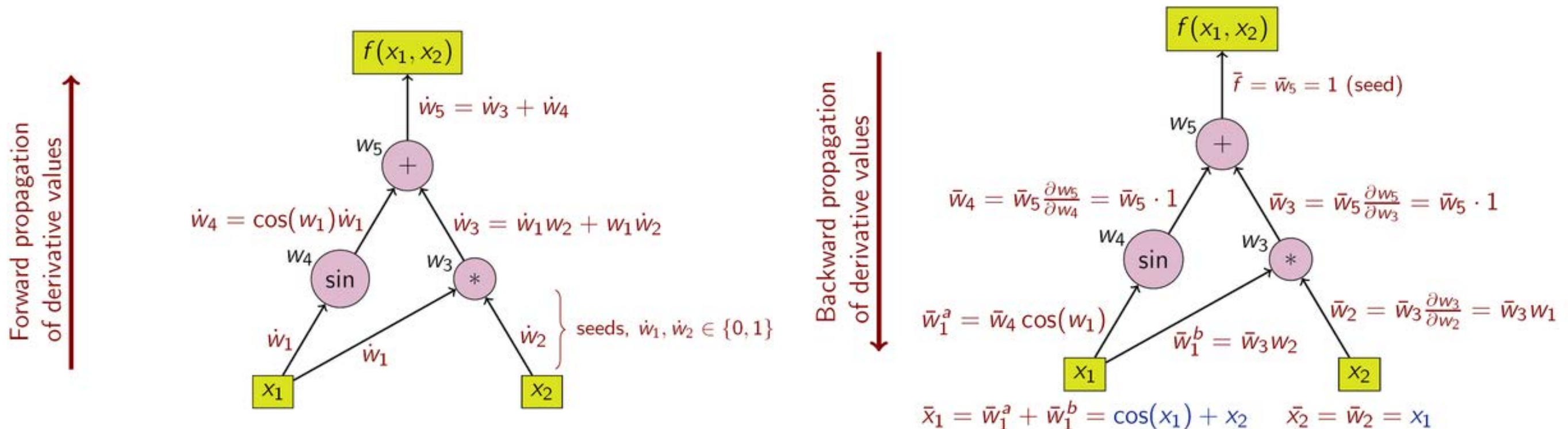
Fundamental to AD is the decomposition of differentials provided by the [chain rule](#). For the simple composition  $y = f(g(h(x))) = f(g(h(w_0))) = f(g(w_1)) = f(w_2) = w_3$  the chain rule gives

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

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  $w_0 = x$ .

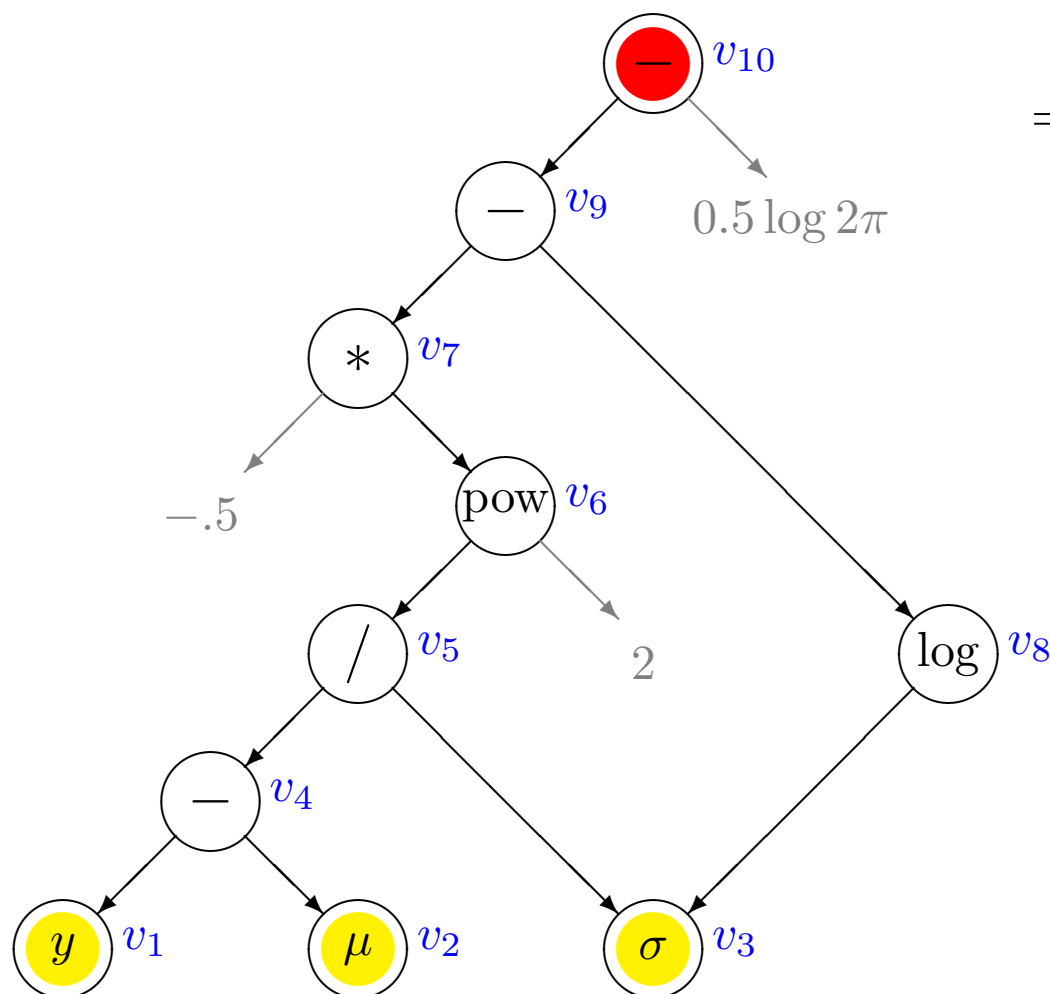
**Example**  $z = f(x_1, x_2) = x_1 x_2 + \sin x_1$



# Automatic differentiation

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

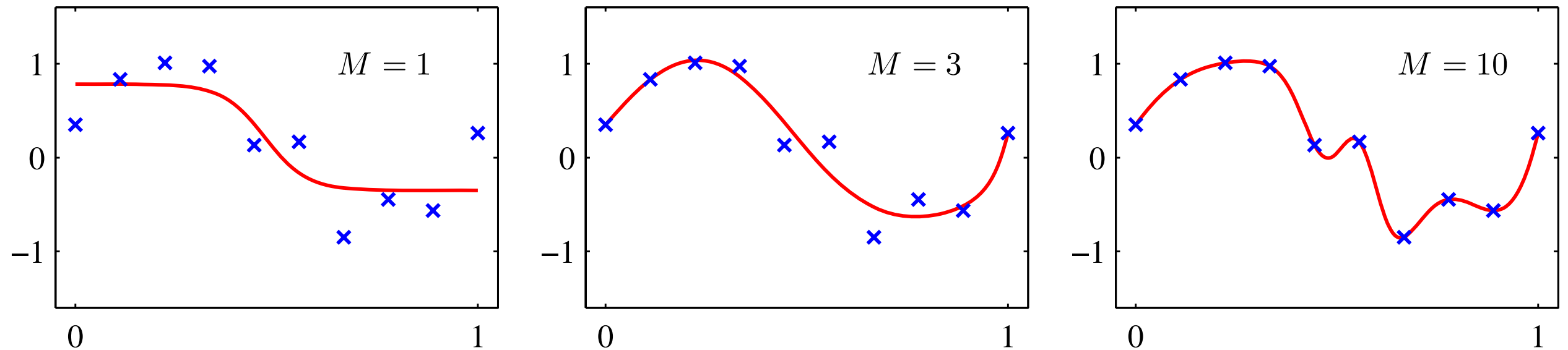
# Automatic differentiation

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.

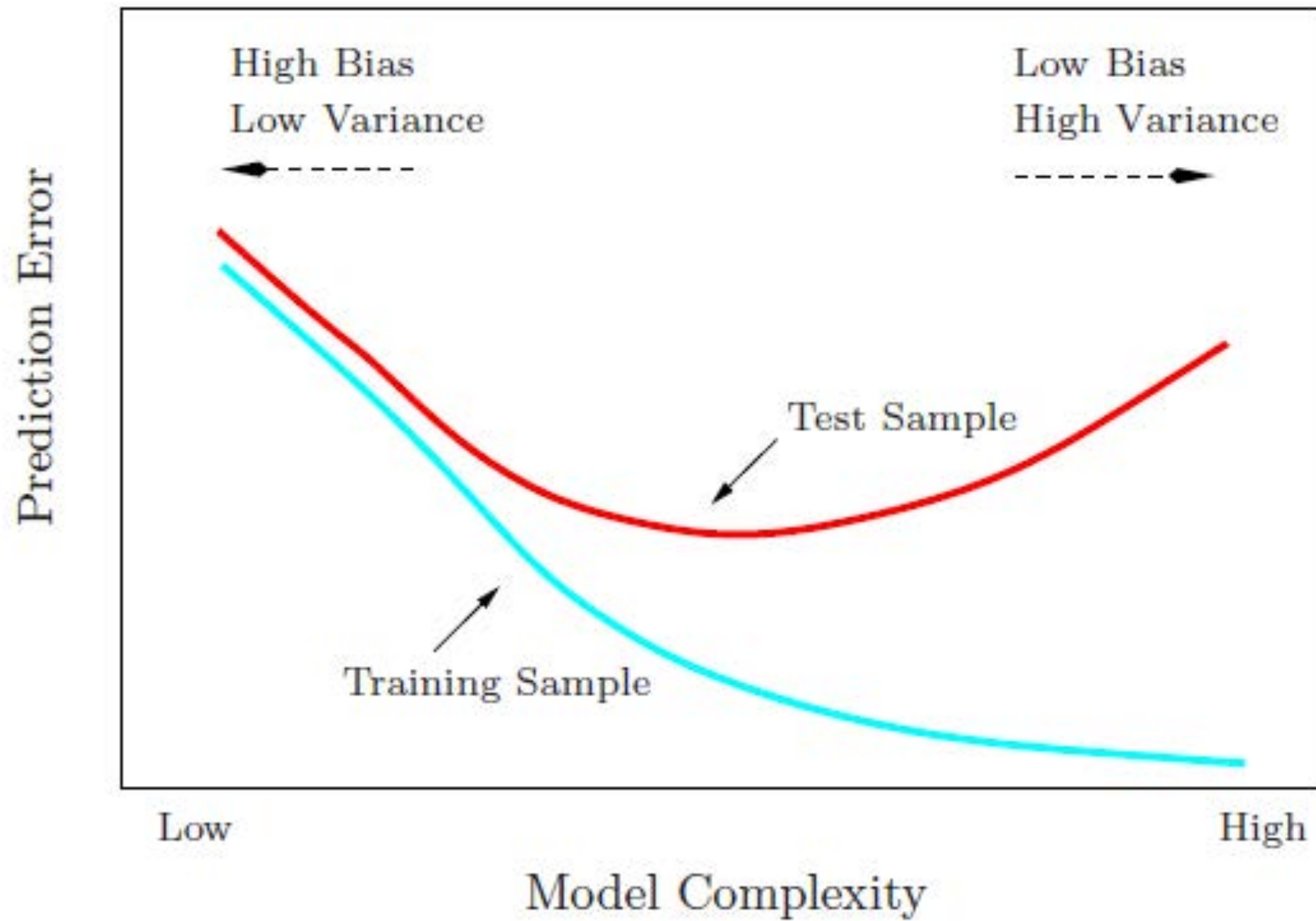
# Overfitting



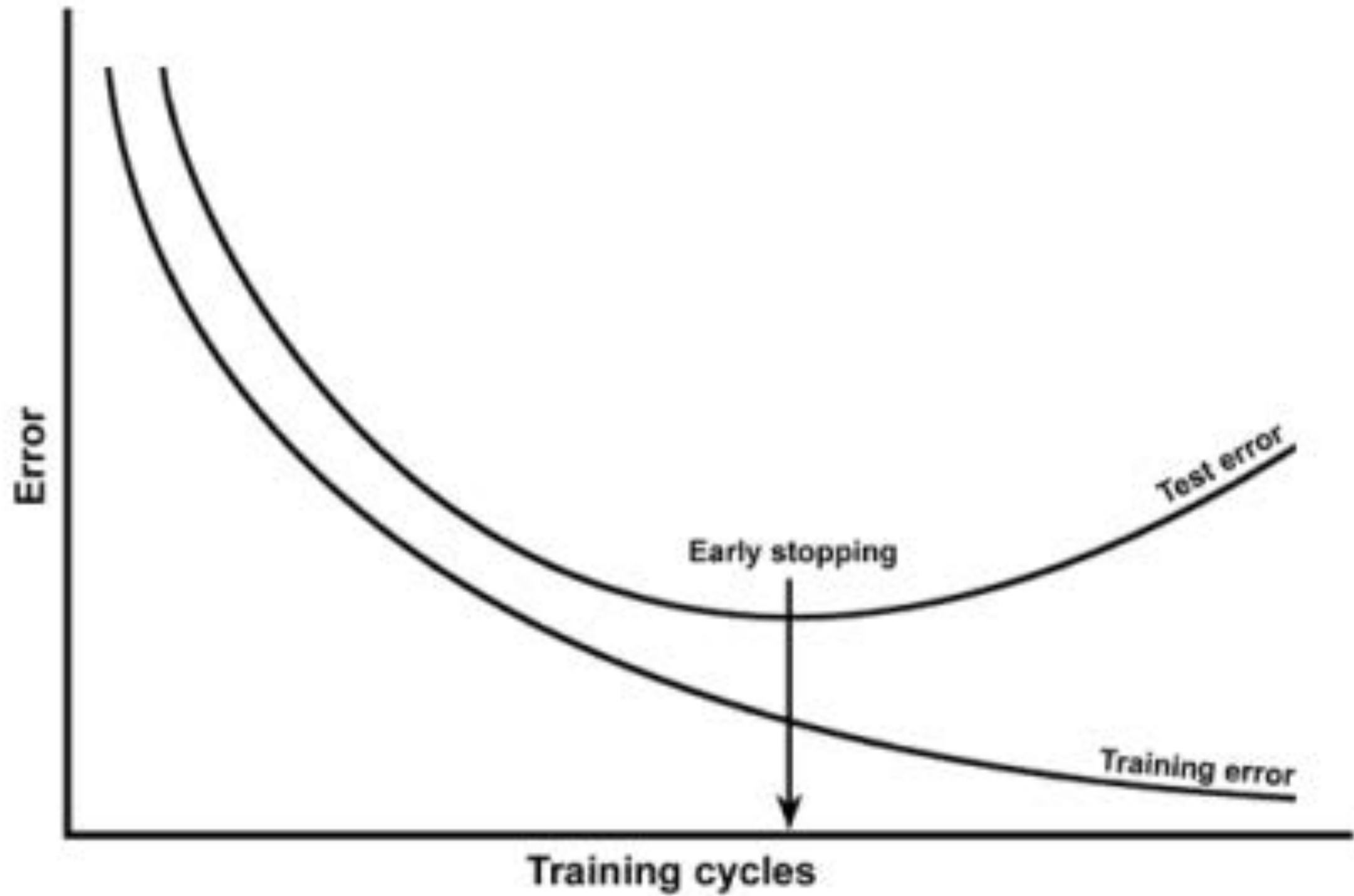
**Figure 5.9** Examples of two-layer networks trained on 10 data points drawn from the sinusoidal data set. The graphs show the result of fitting networks having  $M = 1$ , 3 and 10 hidden units, respectively, by minimizing a sum-of-squares error function using a scaled conjugate-gradient algorithm.



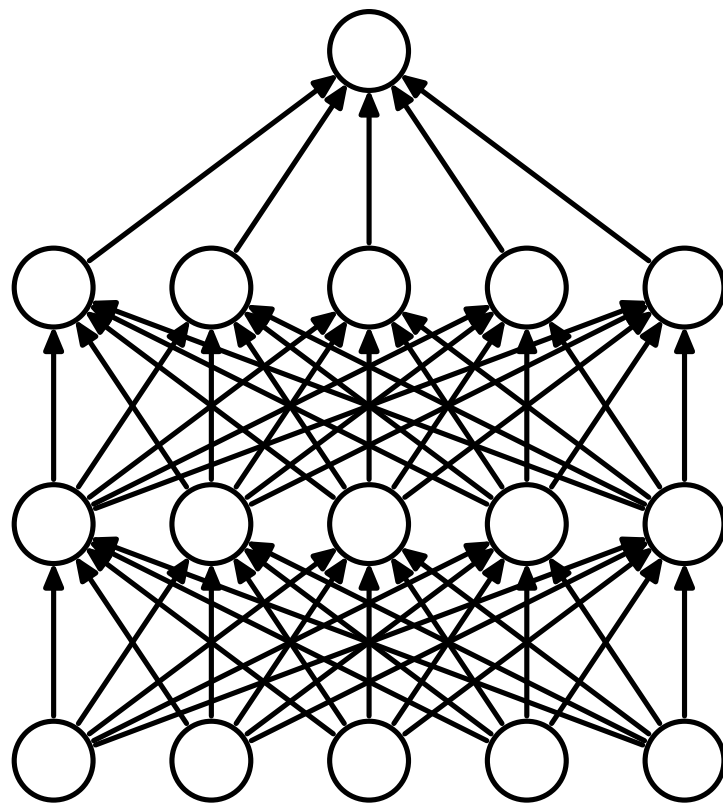
# Overfitting



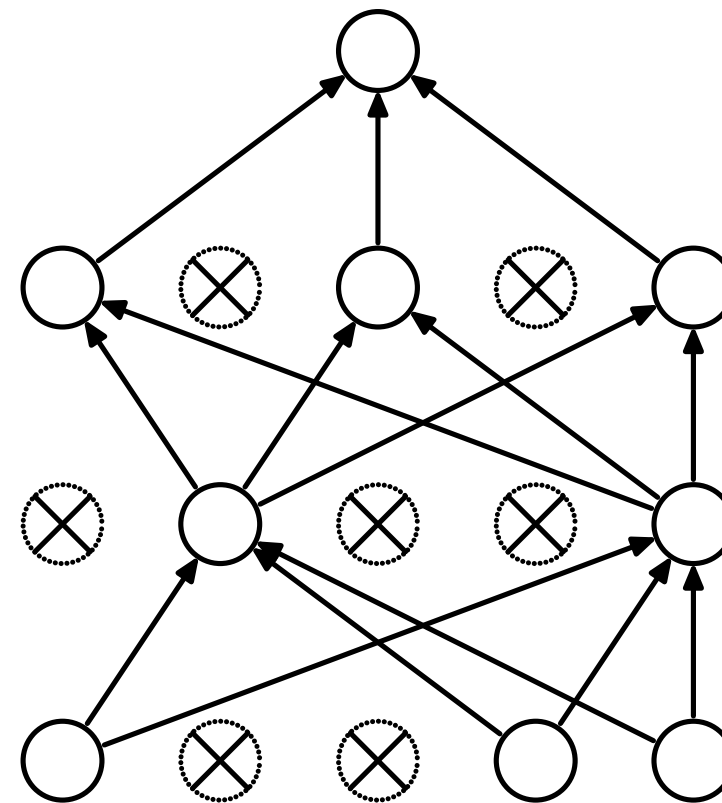
# Early stopping



# Dropout



(a) Standard Neural Net



(b) After applying dropout.

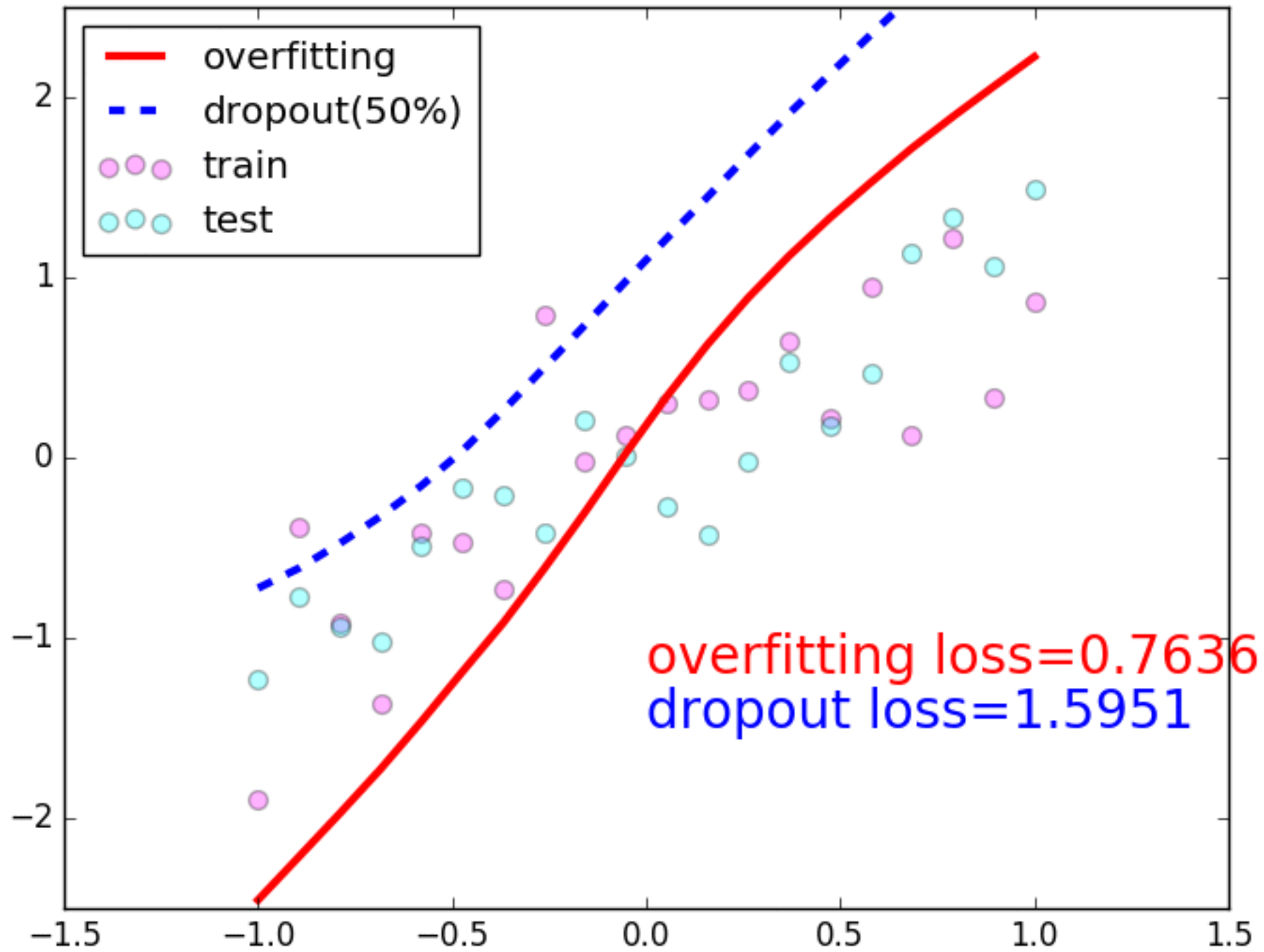
With probability `keep_prob`, outputs the input element scaled up by  $1 / \text{keep\_prob}$ , otherwise outputs 0. The scaling is so that the expected sum is unchanged.

```
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



# Network initialization

---

## Understanding the difficulty of training deep feedforward neural networks

---

**Xavier Glorot**

DIRO, Université de Montréal, Montréal, Québec, Canada

**Yoshua Bengio**

Whereas before 2006 it appears that deep multi-layer neural networks were not successfully trained, since then several algorithms have been shown to successfully train them, with experimental results showing the superiority of deeper vs less deep architectures. All these experimental results were obtained with new initialization or training mechanisms. Our objective here is to understand better why standard gradient descent from random initialization is doing so poorly with deep neural networks, to better understand these recent relative successes and help design better algorithms in the future. We first observe the influence of the non-linear activations functions. We find that the logistic sigmoid activation is unsuited for deep networks with random initialization because of its mean value, which can drive especially the top hidden layer into saturation. Surprisingly, we find that saturated units can move out of saturation by themselves, albeit slowly, and explaining the plateaus sometimes seen when training neural networks. We find that a new non-linearity that saturates less can often be beneficial. Finally, we study how activations and gradients vary across layers and during training, with the idea that training may be more difficult when the singular values of the Jacobian associated with each layer are far from 1. Based on these considerations, we propose a new initialization scheme that brings substantially faster convergence.



# Network initialization

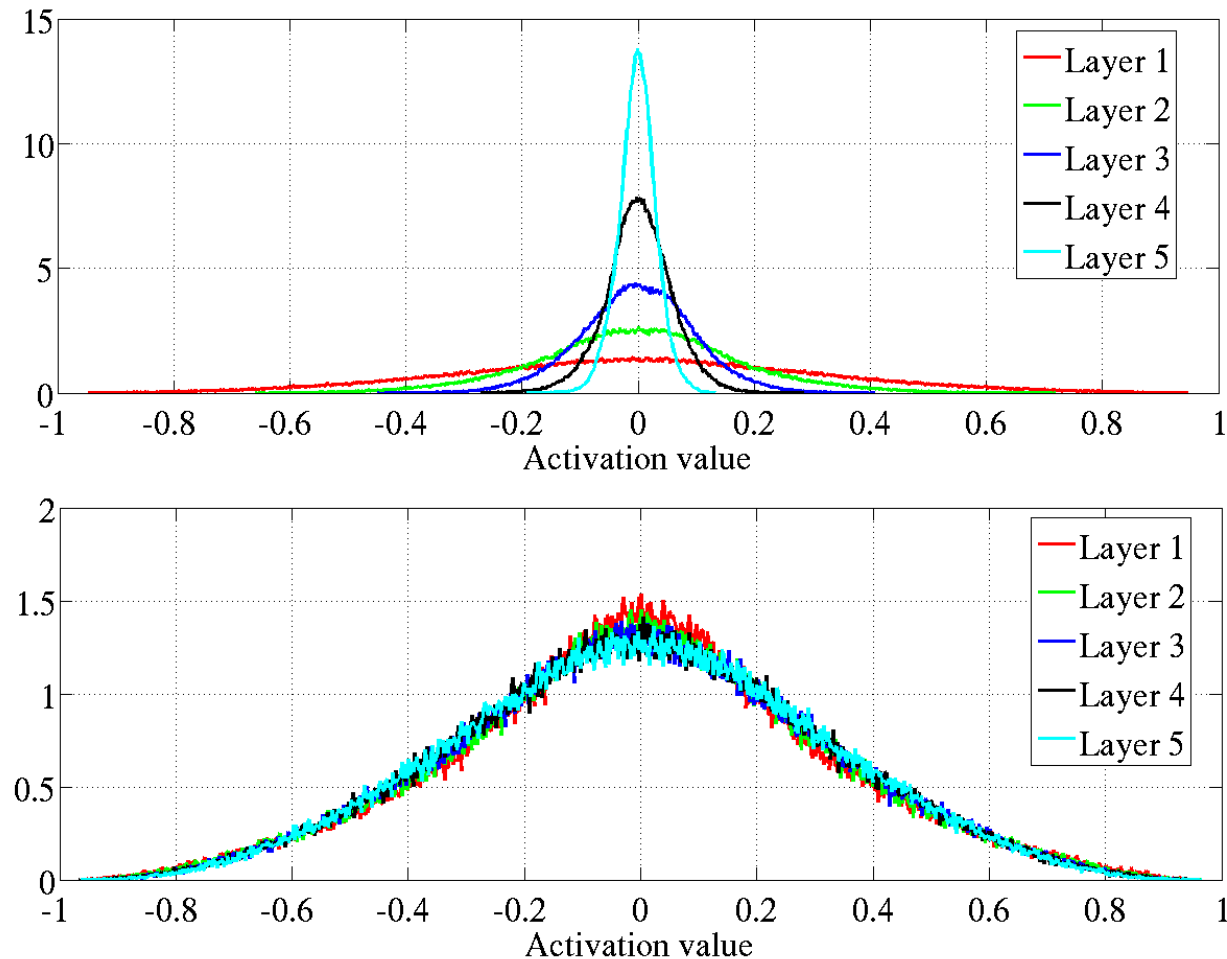


Figure 6: *Activation values normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized initialization (bottom). Top: 0-peak increases for higher layers.*

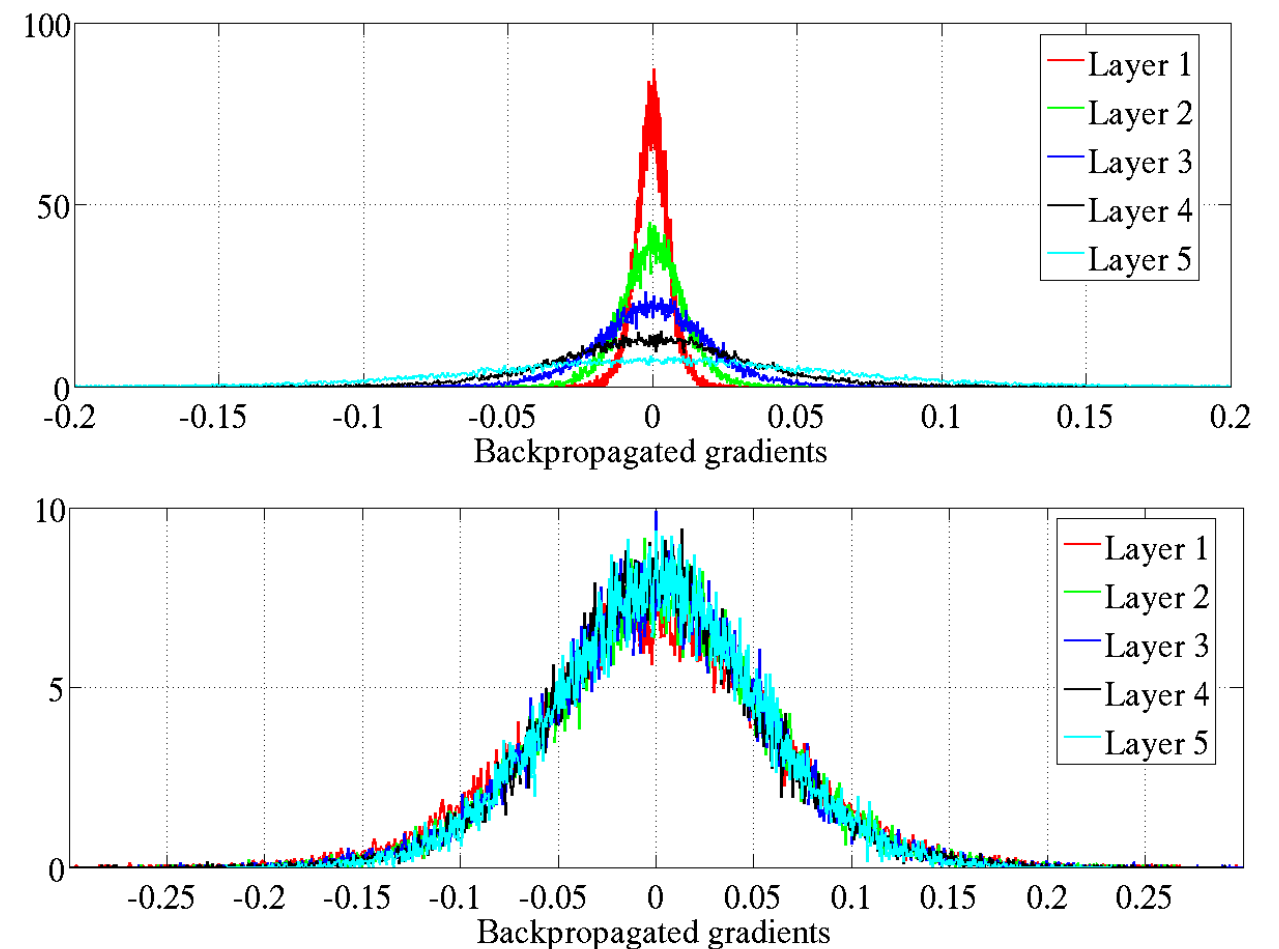


Figure 7: *Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.*

Glorot, X., & Bengio, Y. (2010, March). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (pp. 249-256).

# Tricks of the trade

## Efficient BackProp

Yann LeCun<sup>1</sup>, Leon Bottou<sup>1</sup>, Genevieve B. Orr<sup>2</sup>, and Klaus-Robert Müller<sup>3</sup>

<sup>1</sup> Image Processing Research Department AT&T Labs - Research, 100 Schulz Drive,  
Red Bank, NJ 07701-7033, USA

<sup>2</sup> Willamette University, 900 State Street, Salem, OR 97301, USA

<sup>3</sup> GMD FIRST, Rudower Chaussee 5, 12489 Berlin, Germany  
{yann,leonb}@research.att.com, gorr@willamette.edu, klaus@first.gmd.de

originally published in

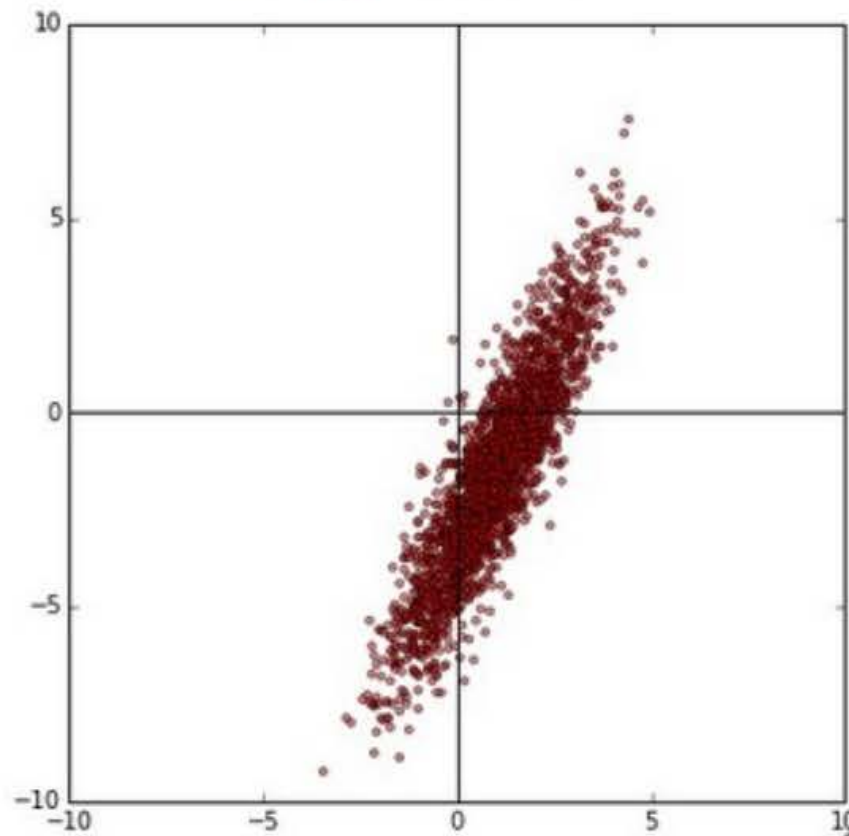
Orr, G. and Müller, K. “Neural Networks: tricks of the trade”,  
Springer, 1998.

**Abstract.** The convergence of back-propagation learning is analyzed so as to explain common phenomenon observed by practitioners. Many undesirable behaviors of backprop can be avoided with tricks that are rarely exposed in serious technical publications. This paper gives some of those tricks, and offers explanations of why they work.

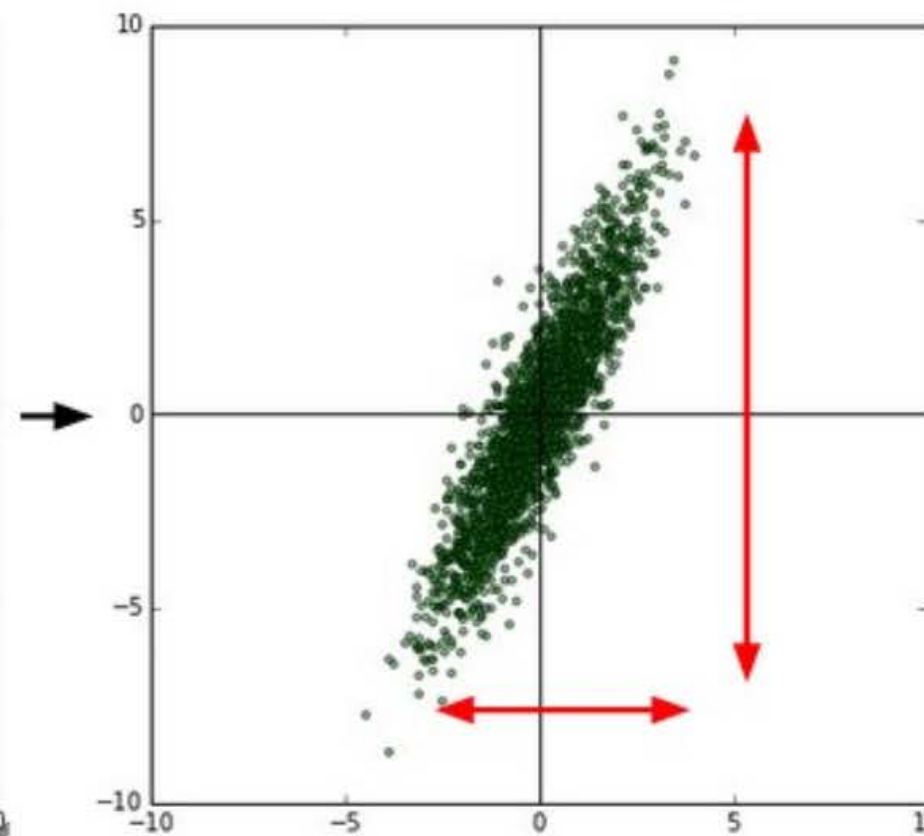
Many authors have suggested that second-order optimization methods are advantageous for neural net training. It is shown that most “classical” second-order methods are impractical for large neural networks. A few methods are proposed that do not have these limitations.

# Normalizing the inputs

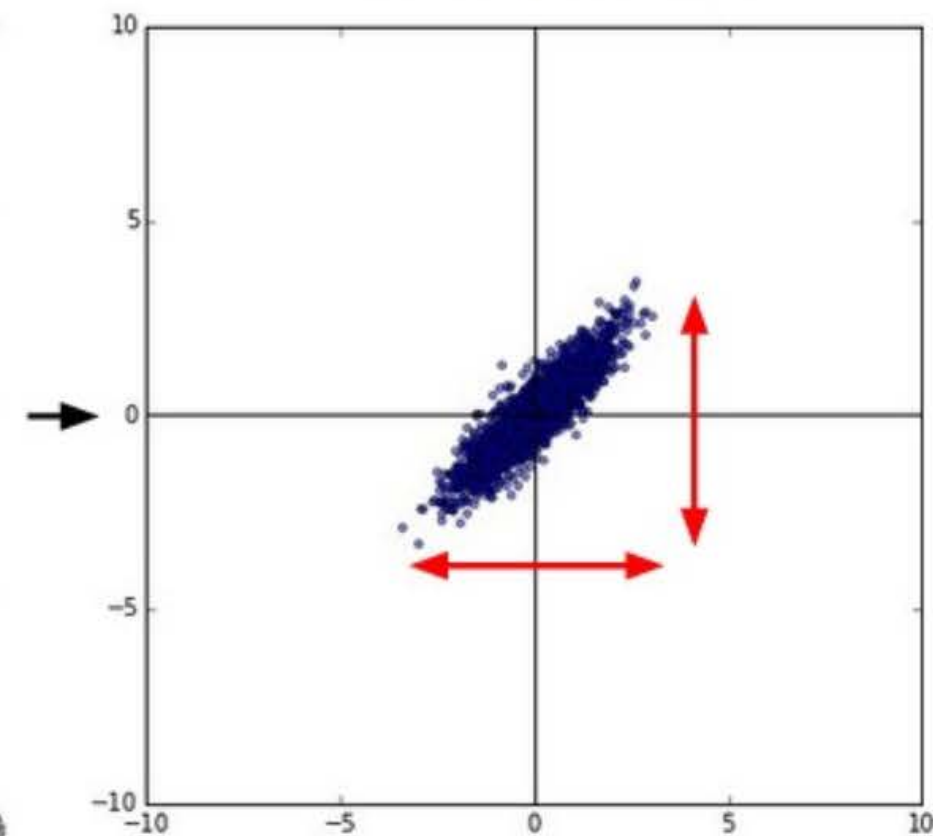
original data



zero-centered data



normalized data



Common data preprocessing pipeline. **Left:** Original toy, 2-dimensional input data. **Middle:** The data is zero-centered by subtracting the mean in each dimension. The data cloud is now centered around the origin. **Right:** Each dimension is additionally scaled by its standard deviation. The red lines indicate the extent of the data - they are of unequal length in the middle, but of equal length on the right.