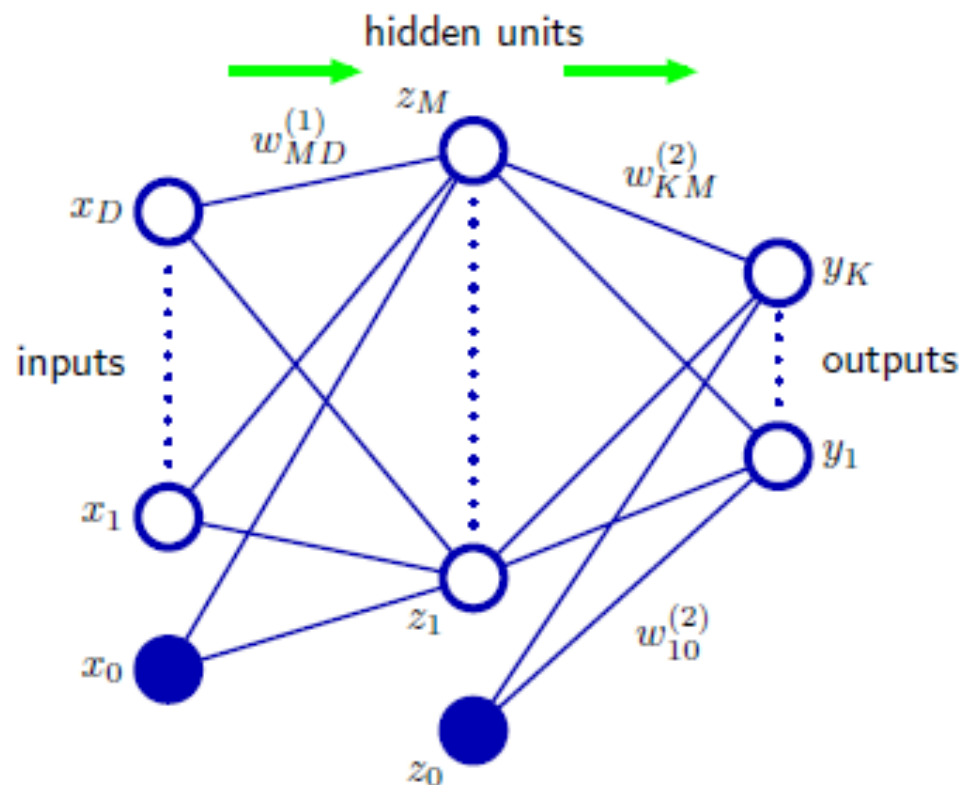
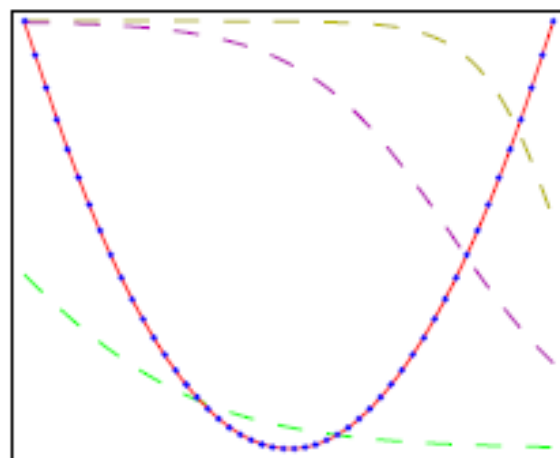


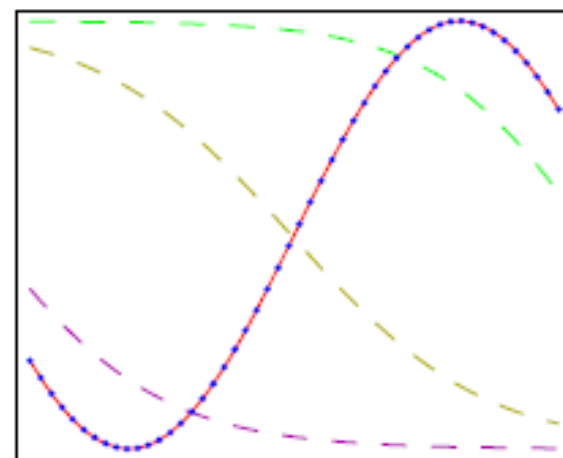
**Figure 5.1** Network diagram for the two-layer neural network corresponding to (5.7). The input, hidden, and output variables are represented by nodes, and the weight parameters are represented by links between the nodes, in which the bias parameters are denoted by links coming from additional input and hidden variables  $x_0$  and  $z_0$ . Arrows denote the direction of information flow through the network during forward propagation.



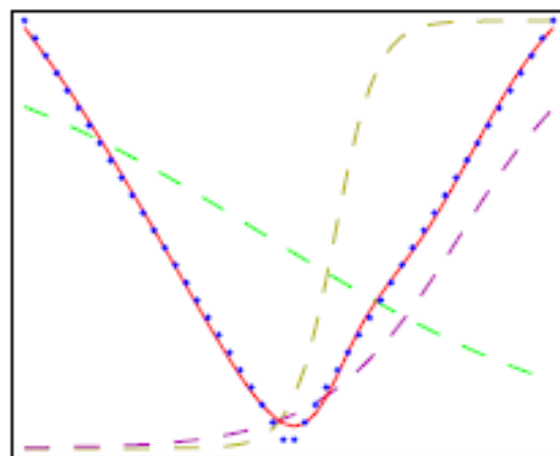
**Figure 5.3** Illustration of the capability of a multilayer perceptron to approximate four different functions comprising (a)  $f(x) = x^2$ , (b)  $f(x) = \sin(x)$ , (c),  $f(x) = |x|$ , and (d)  $f(x) = H(x)$  where  $H(x)$  is the Heaviside step function. In each case,  $N = 50$  data points, shown as blue dots, have been sampled uniformly in  $x$  over the interval  $(-1, 1)$  and the corresponding values of  $f(x)$  evaluated. These data points are then used to train a two-layer network having 3 hidden units with 'tanh' activation functions and linear output units. The resulting network functions are shown by the red curves, and the outputs of the three hidden units are shown by the three dashed curves.



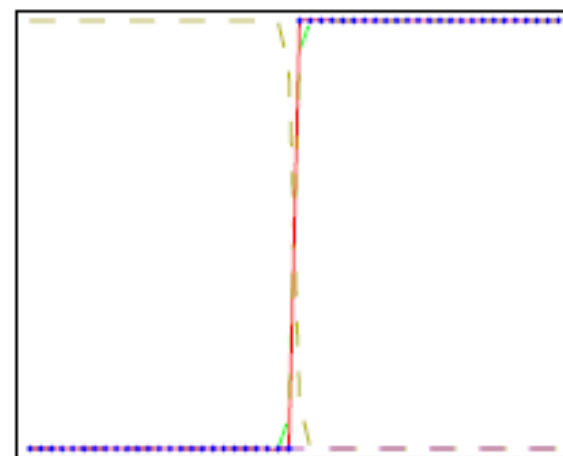
(a)



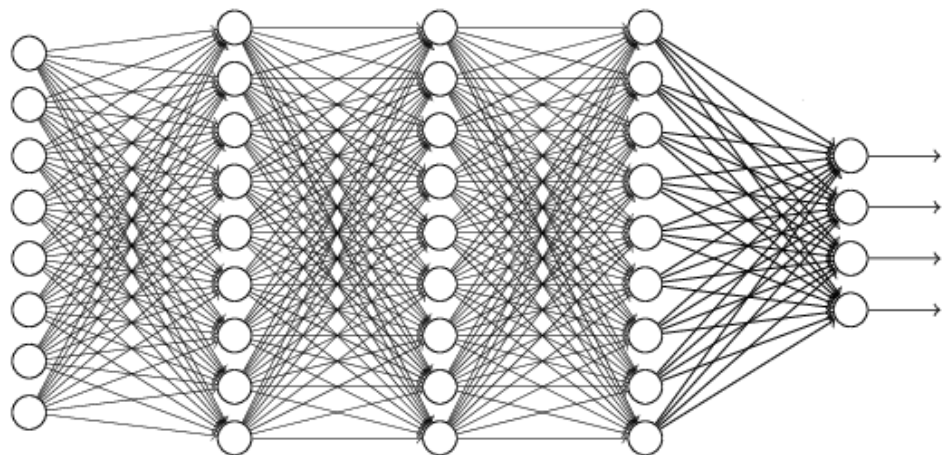
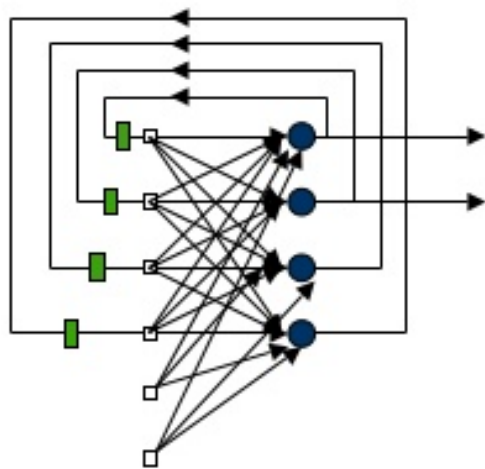
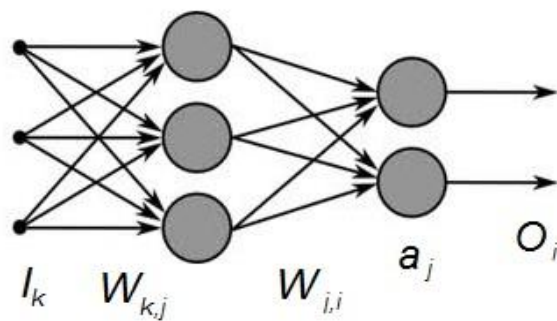
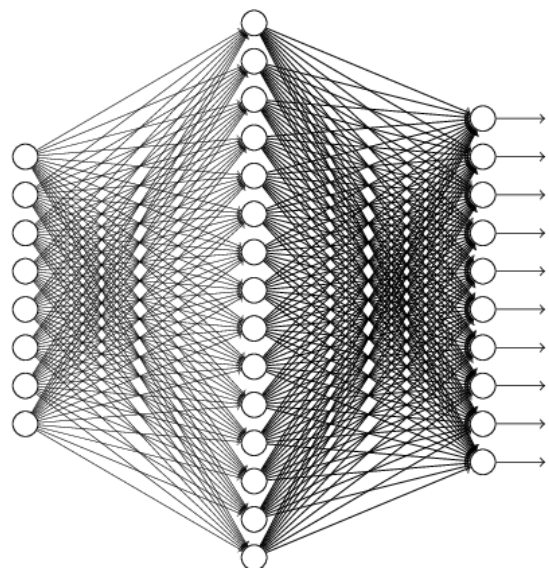
(b)



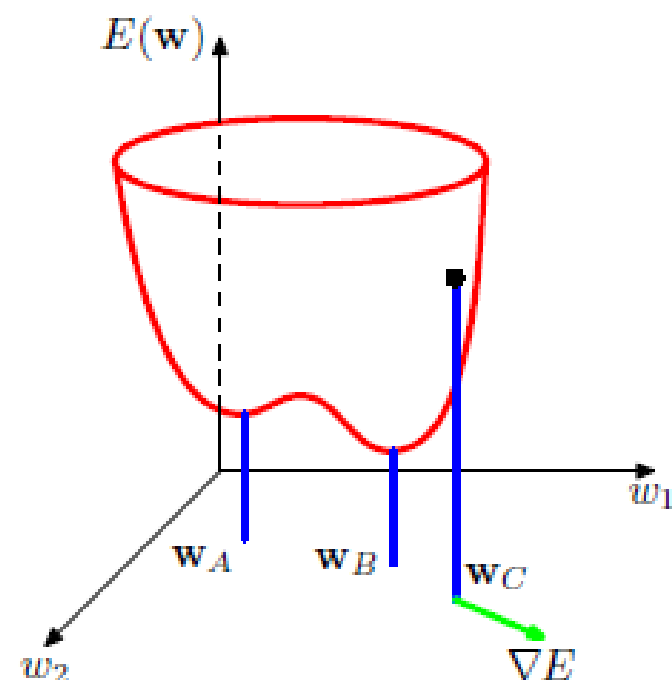
(c)



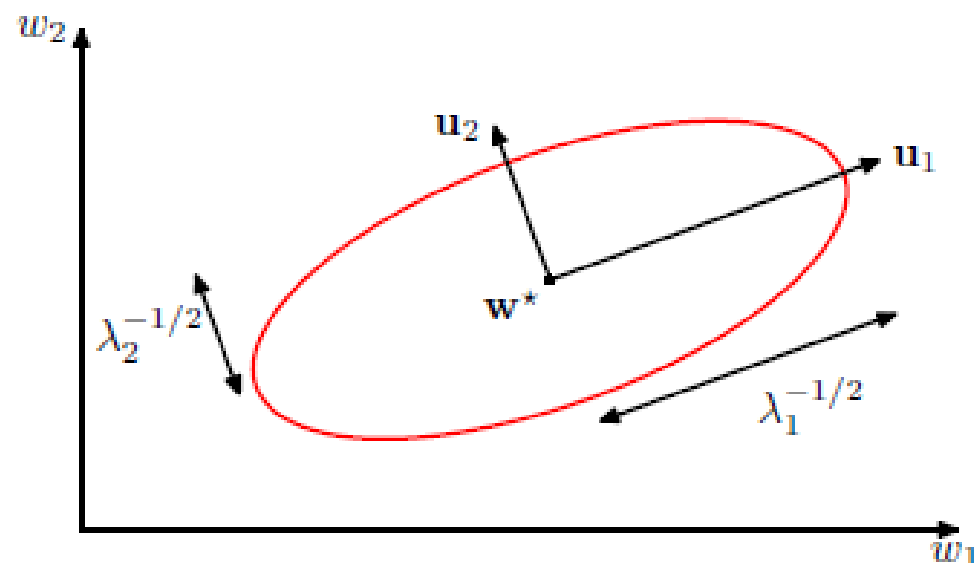
(d)

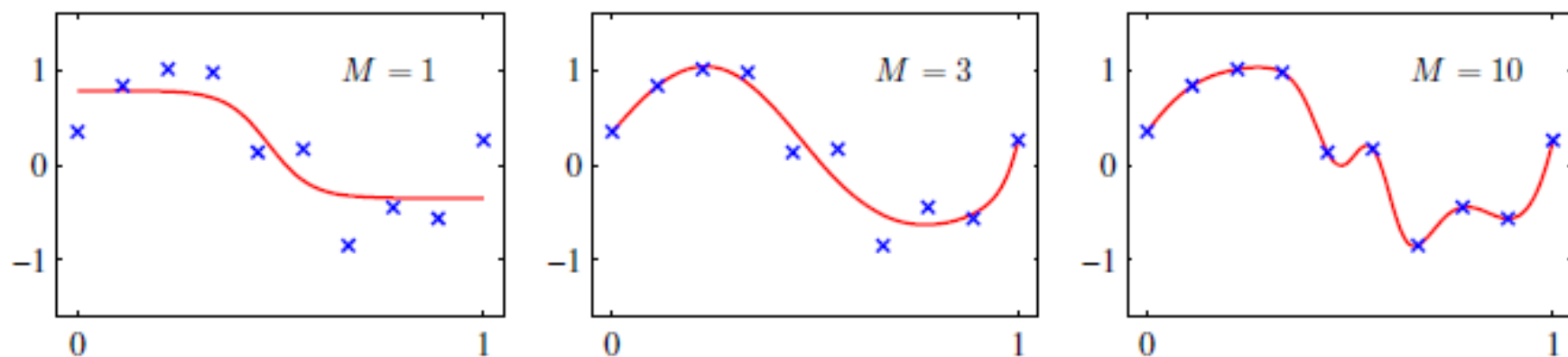


**Figure 5.5** Geometrical view of the error function  $E(\mathbf{w})$  as a surface sitting over weight space. Point  $\mathbf{w}_A$  is a local minimum and  $\mathbf{w}_B$  is the global minimum. At any point  $\mathbf{w}_C$ , the local gradient of the error surface is given by the vector  $\nabla E$ .

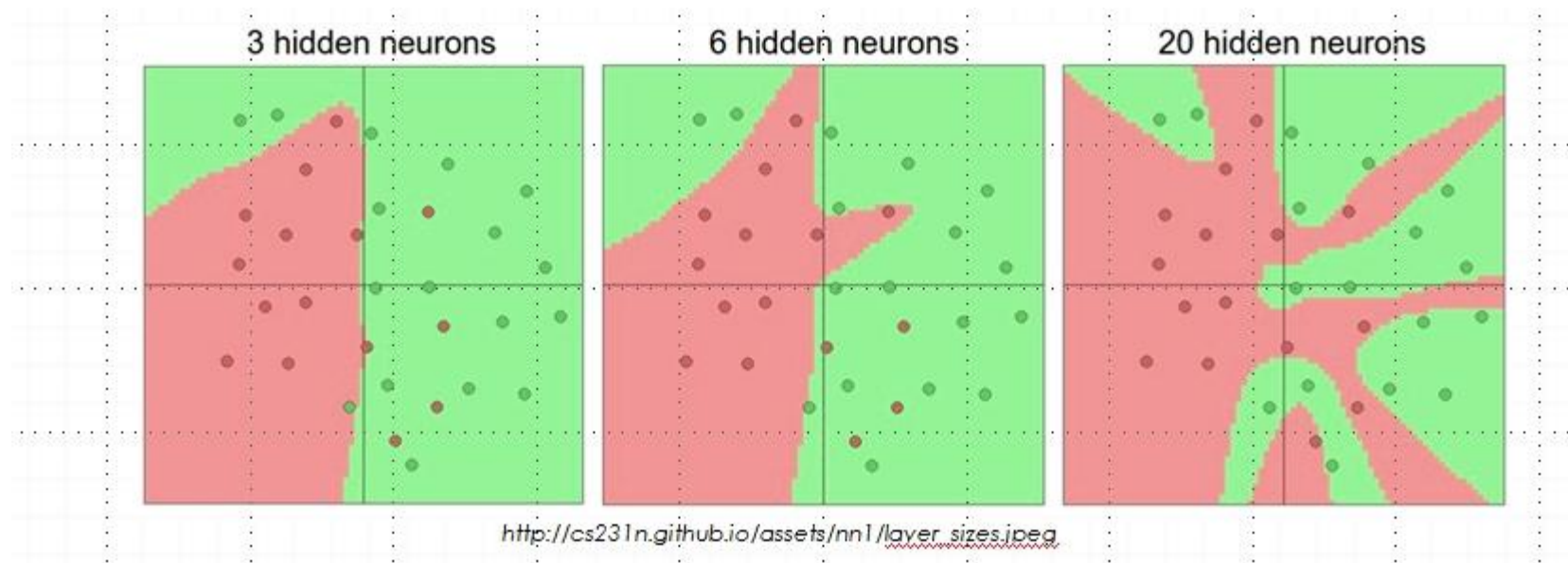


**Figure 5.6** In the neighbourhood of a minimum  $\mathbf{w}^*$ , the error function can be approximated by a quadratic. Contours of constant error are then ellipses whose axes are aligned with the eigenvectors  $\mathbf{u}_i$  of the Hessian matrix, with lengths that are inversely proportional to the square roots of the corresponding eigenvalues  $\lambda_i$ .





**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.



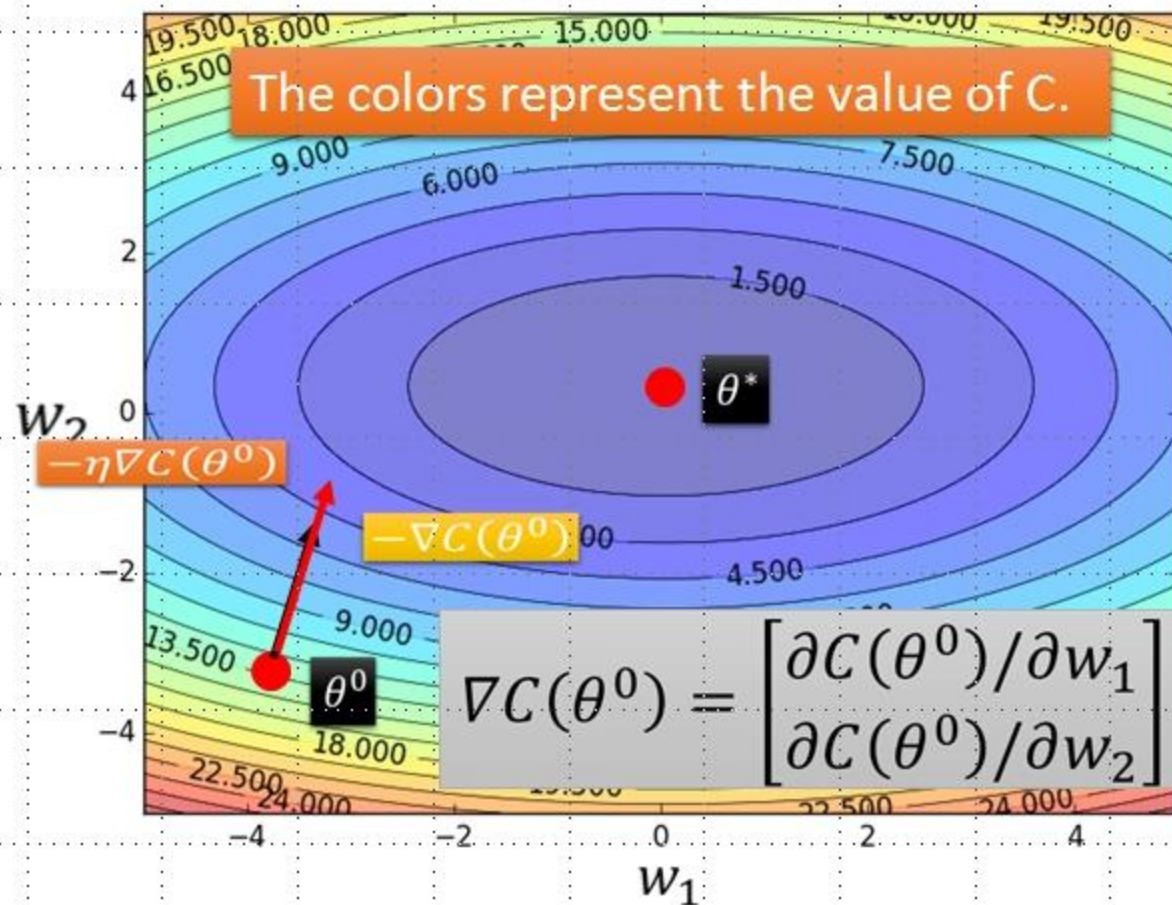


# Gradient Descent

Assume there are only two parameters  $w_1$  and  $w_2$  in a network.

$$\theta = \{w_1, w_2\}$$

Error Surface



Randomly pick a starting point  $\theta^0$

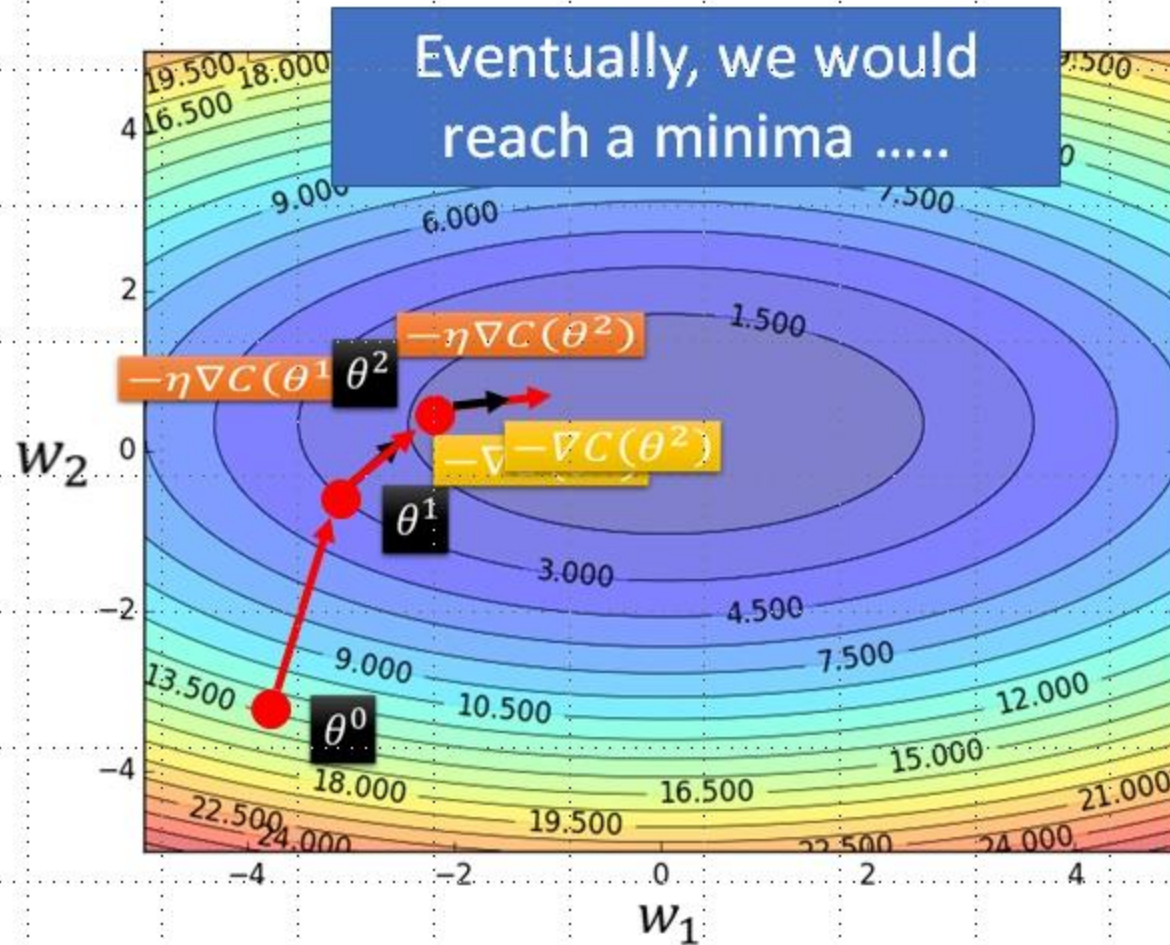
Compute the negative gradient at  $\theta^0$

→  $-\nabla C(\theta^0)$

Times the learning rate  $\eta$

→  $-\eta \nabla C(\theta^0)$

# Gradient Descent



Randomly pick a starting point  $\theta^0$

Compute the negative gradient at  $\theta^0$

→  $-\nabla C(\theta^0)$

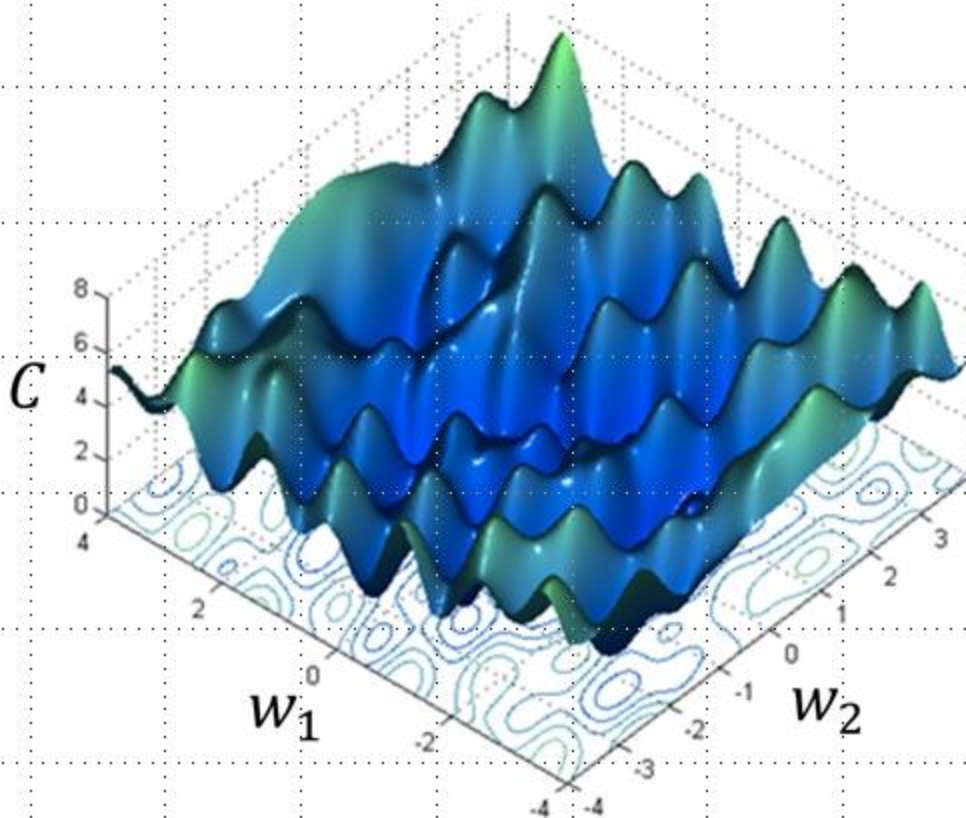
Times the learning rate  $\eta$

→  $-\eta \nabla C(\theta^0)$



# Local Minima

- Gradient descent never guarantee global minima



Different initial  
point  $\theta^0$



Reach different minima,  
so different results



# Loss functions and output

## Classification

**Training examples**

$\mathbb{R}^n \times \{\text{class\_1}, \dots, \text{class\_n}\}$   
(one-hot encoding)

**Output Layer**

Soft-max  
[map  $\mathbb{R}^n$  to a probability distribution]

$$P(y = j \mid \mathbf{x}) = \frac{e^{\mathbf{x}^T \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^T \mathbf{w}_k}}$$

**Cost (loss) function**

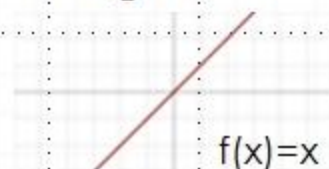
Cross-entropy

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K [y_k^{(i)} \log \hat{y}_k^{(i)} + (1 - y_k^{(i)}) \log (1 - \hat{y}_k^{(i)})]$$

## Regression

$\mathbb{R}^n \times \mathbb{R}^m$

Linear (Identity)  
or Sigmoid



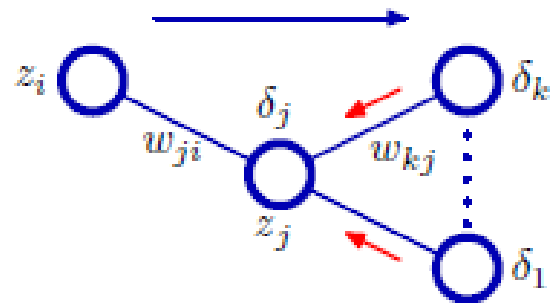
Mean Squared Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

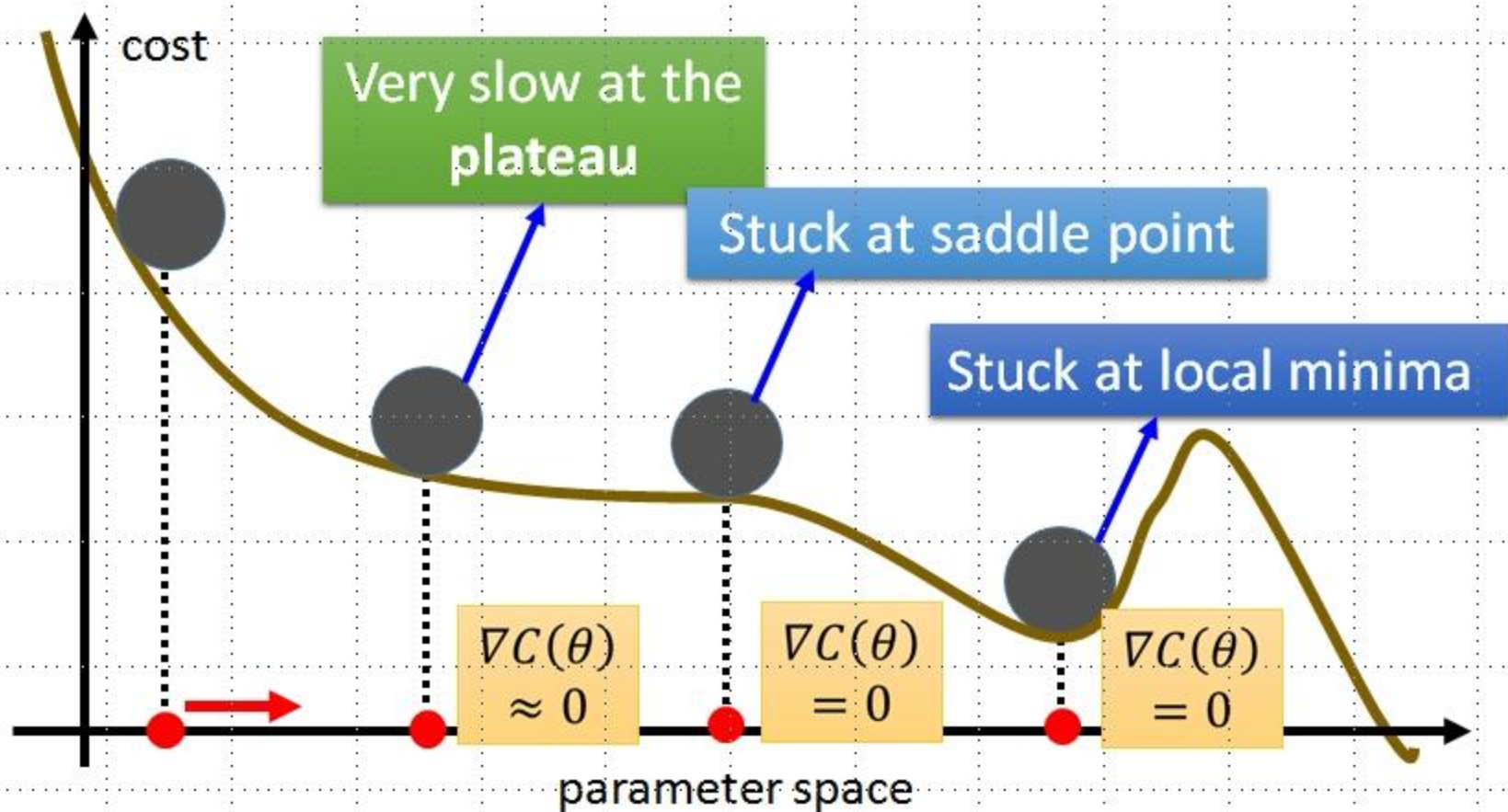
Mean Absolute Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n |y^{(i)} - \hat{y}^{(i)}|$$

**Figure 5.7** Illustration of the calculation of  $\delta_j$  for hidden unit  $j$  by backpropagation of the  $\delta$ 's from those units  $k$  to which unit  $j$  sends connections. The blue arrow denotes the direction of information flow during forward propagation, and the red arrows indicate the backward propagation of error information.



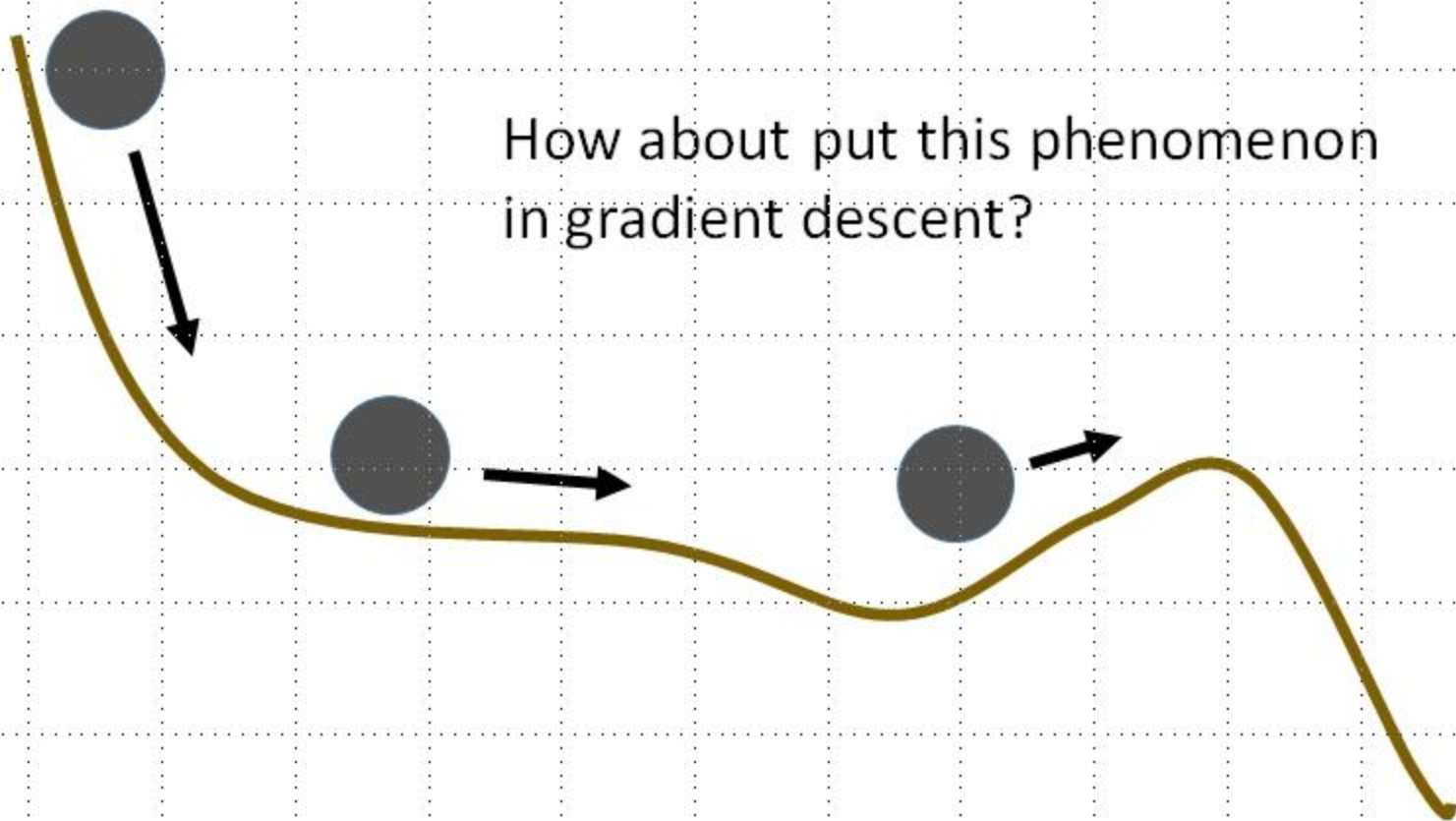
Besides local minima .....



In physical world .....

- Momentum

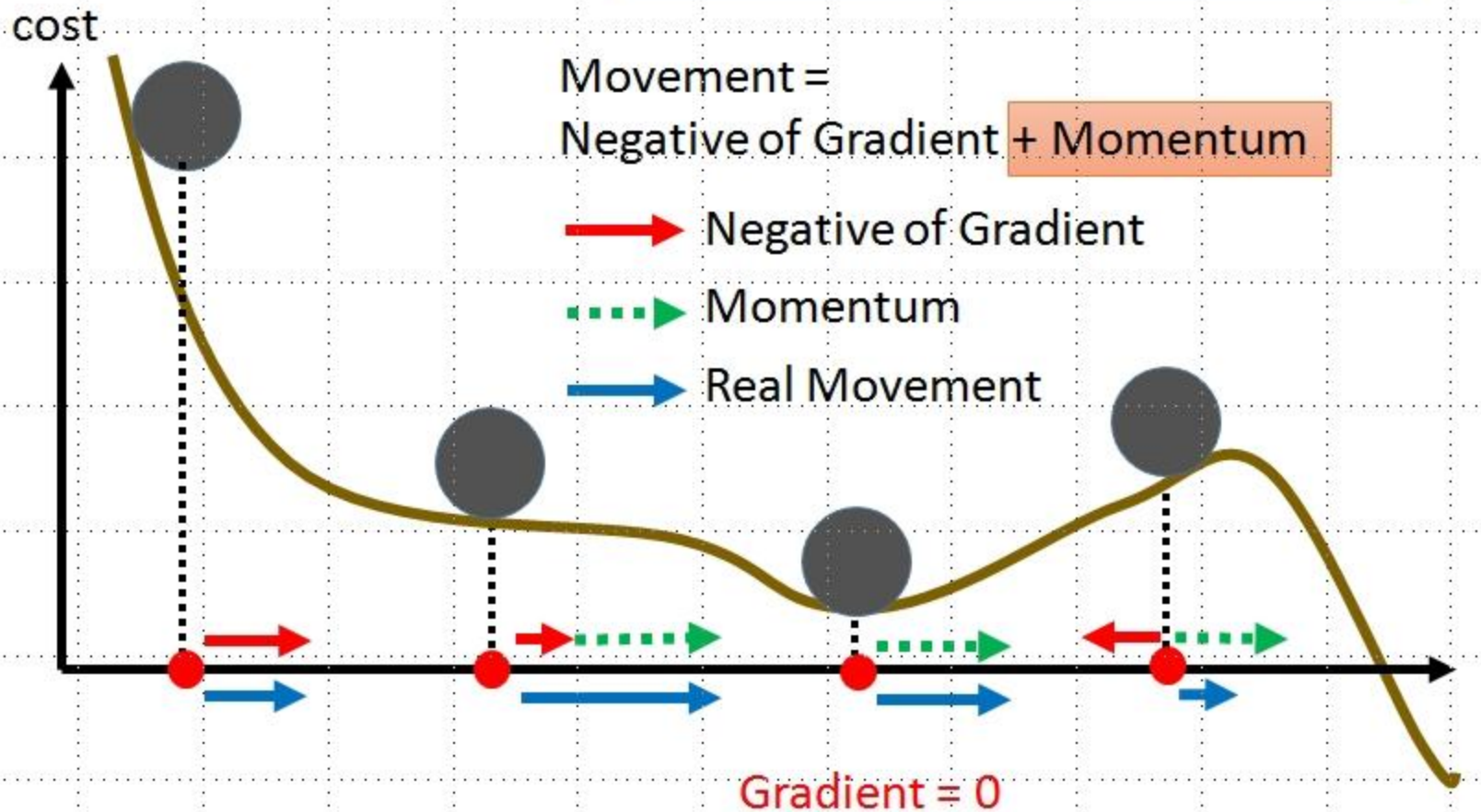
How about put this phenomenon  
in gradient descent?

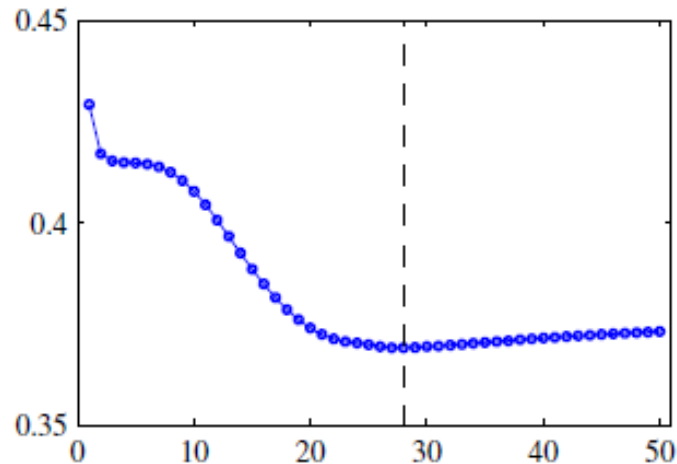
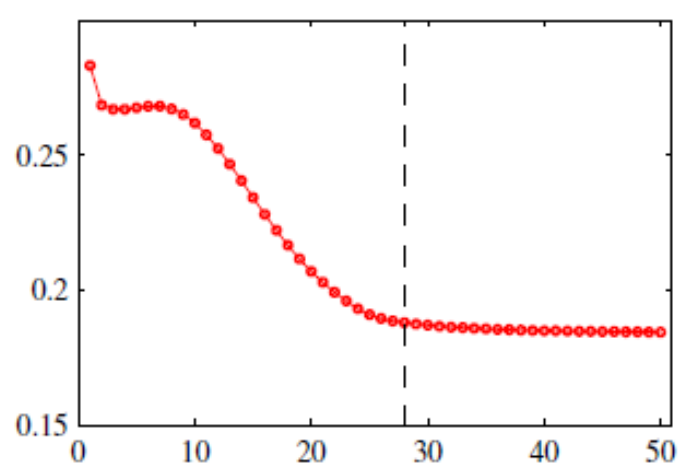




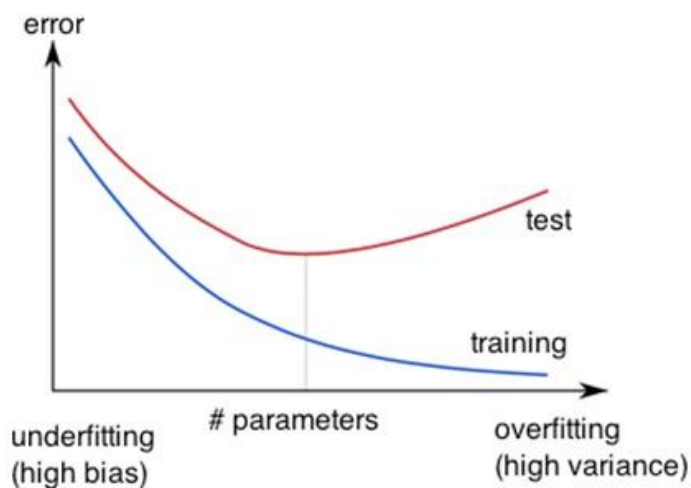
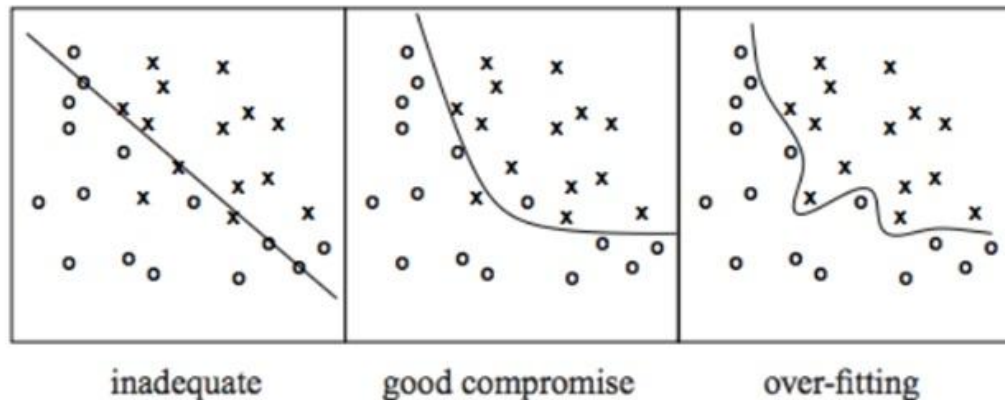
# Momentum

Still not guarantee reaching global minima, but give some hope .....

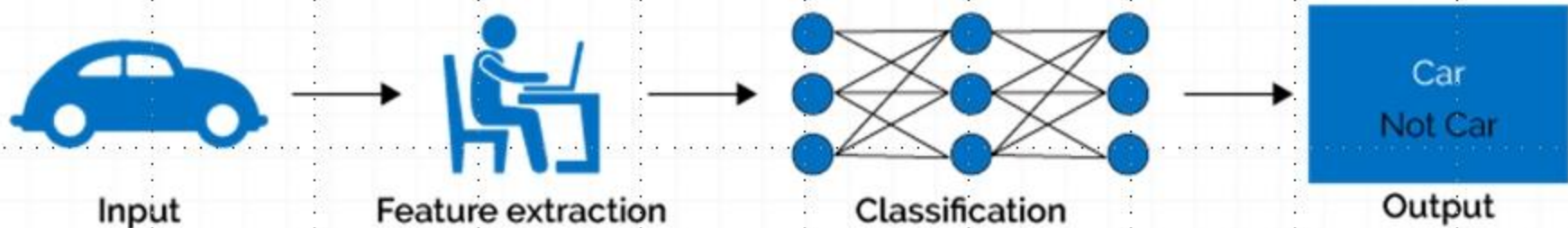




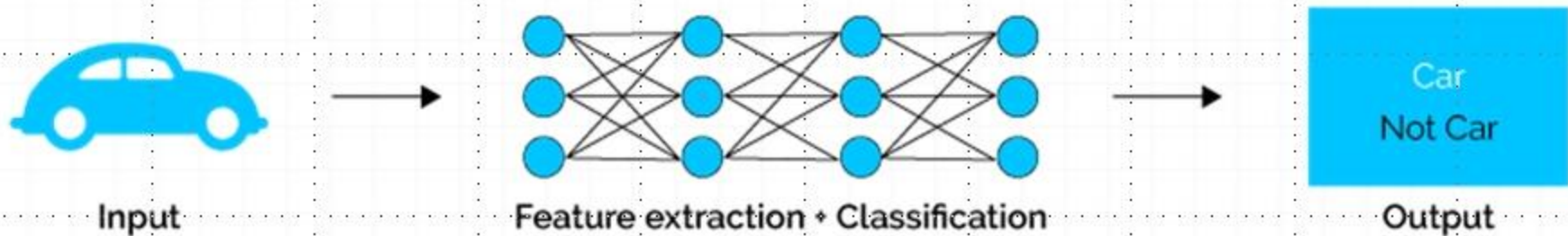
**Figure 5.12** An illustration of the behaviour of training set error (left) and validation set error (right) during a typical training session, as a function of the iteration step. The best generalization performance suggests that training dashed lines, corresponding to the minimum of the validation



## Machine Learning



## Deep Learning



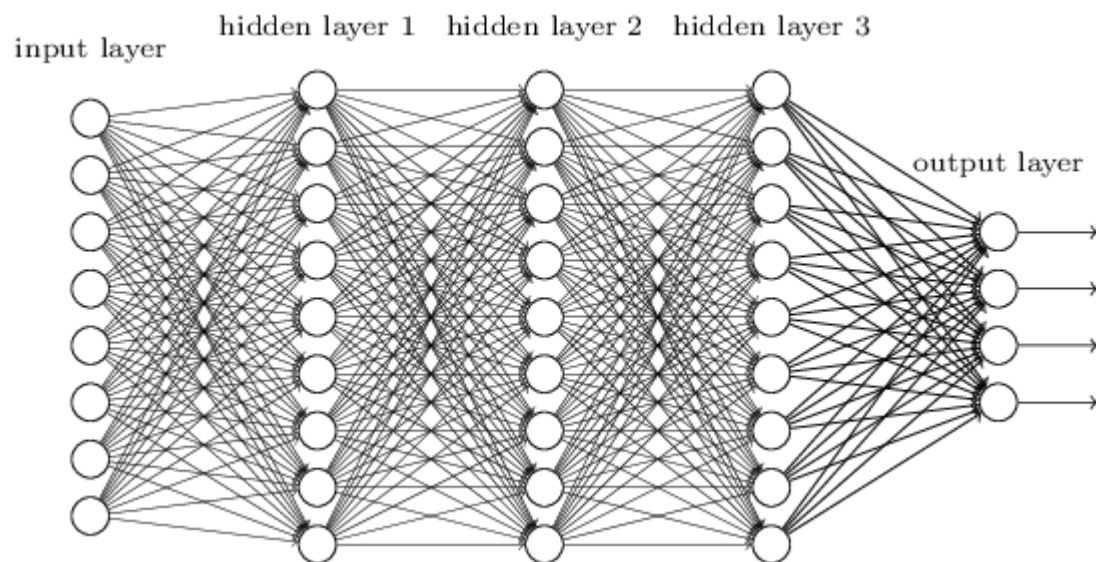
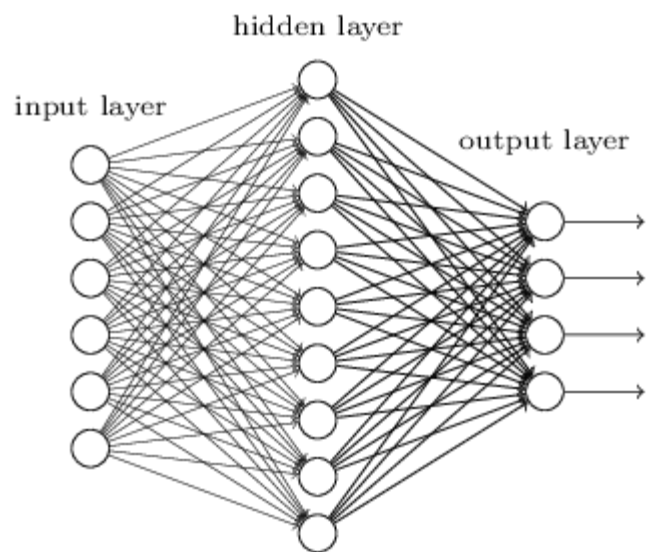






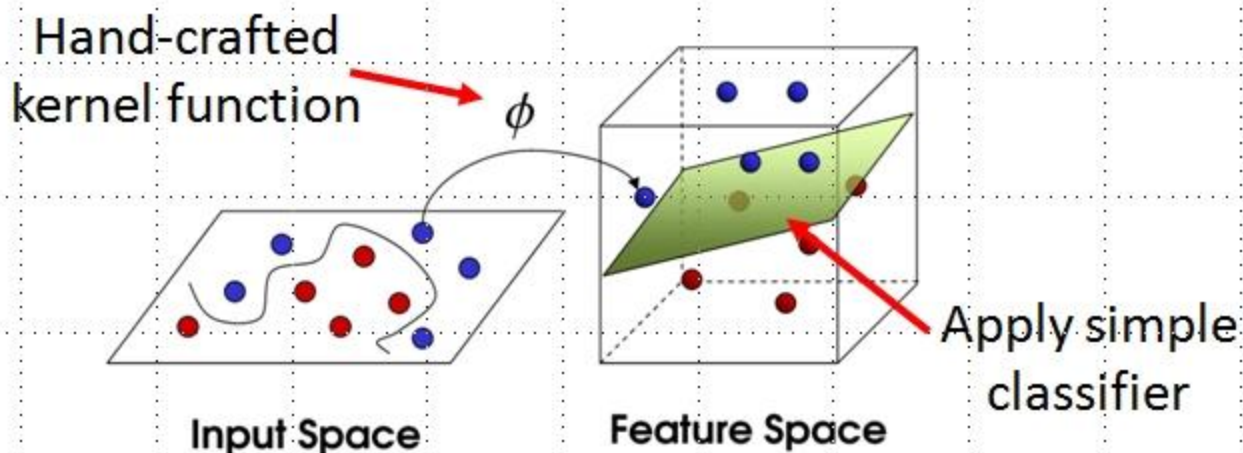






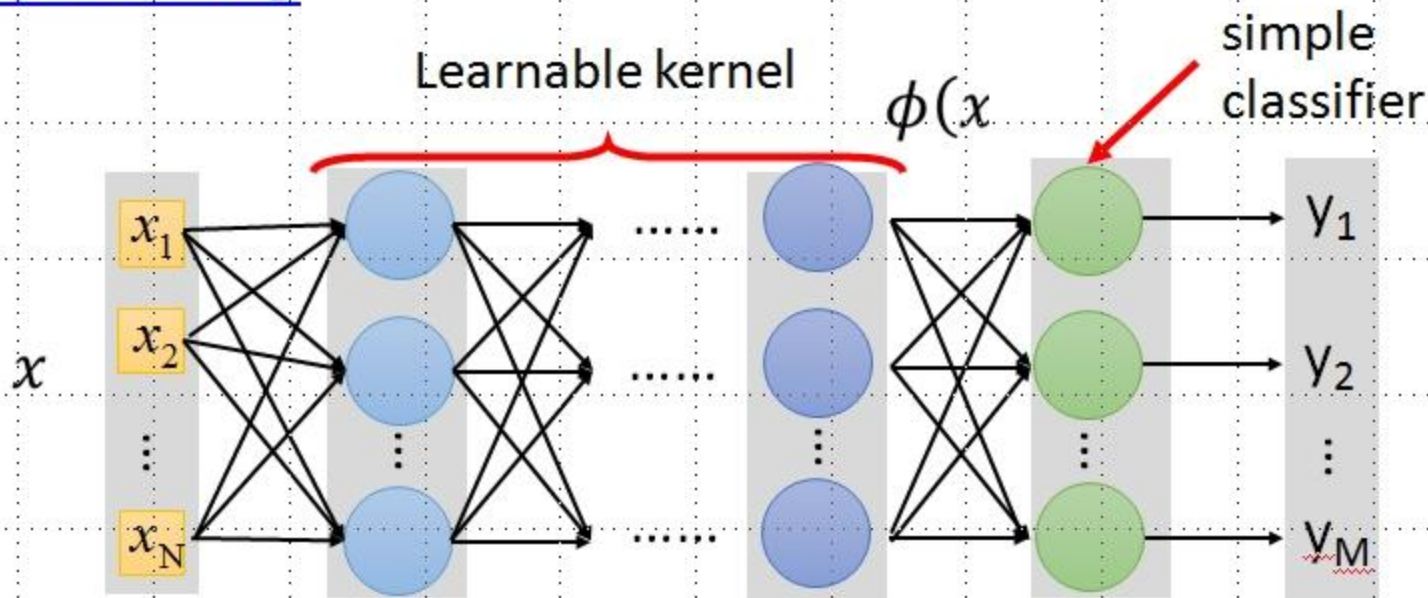


## SVM



Source of image: [http://www.gipsa-lab.grenoble-inp.fr/transfert/seminaire/455\\_Kadri2013Gipsa-lab.pdf](http://www.gipsa-lab.grenoble-inp.fr/transfert/seminaire/455_Kadri2013Gipsa-lab.pdf)

## Deep Learning







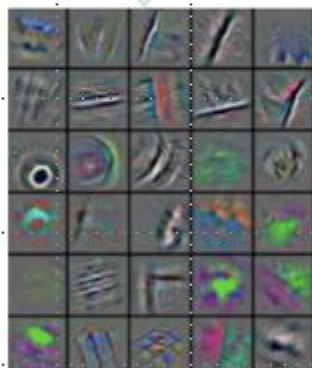
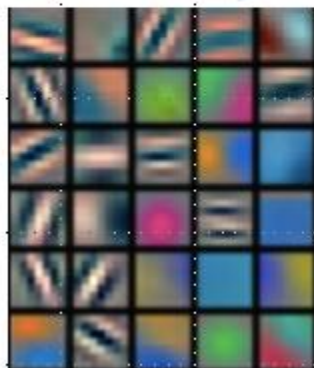
Low-level  
features

Mid-level  
features

High-level  
features

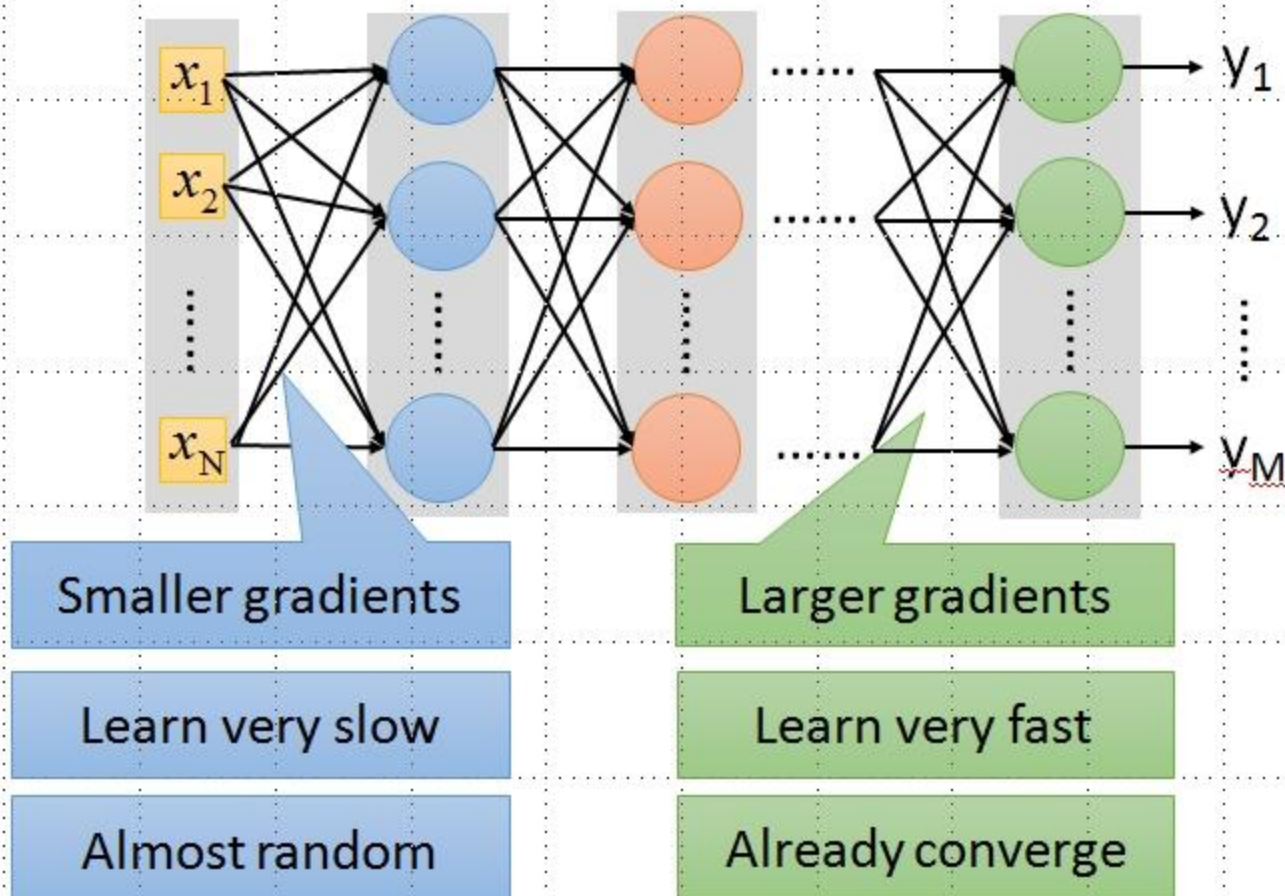
Trainable  
classifier

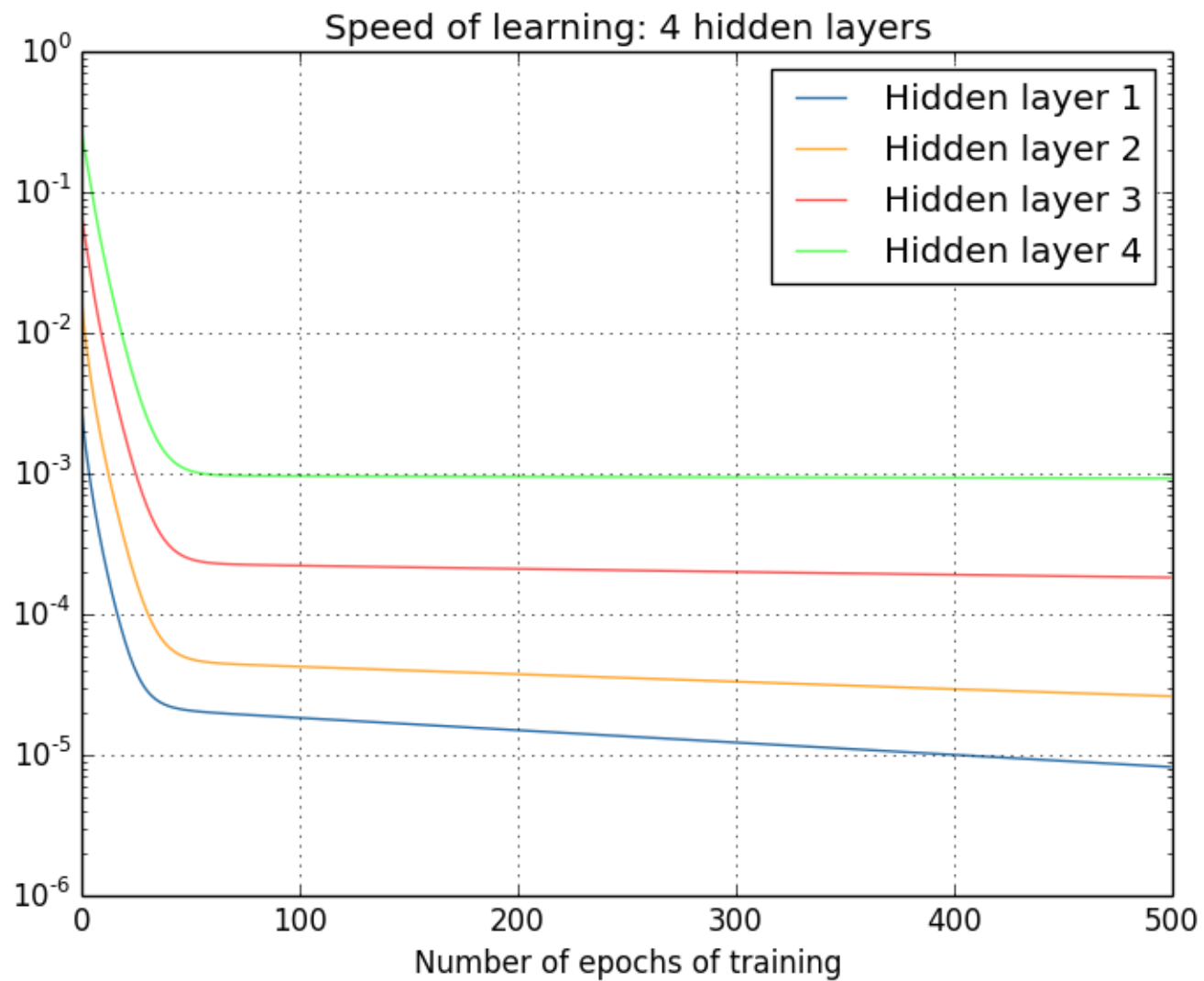
output

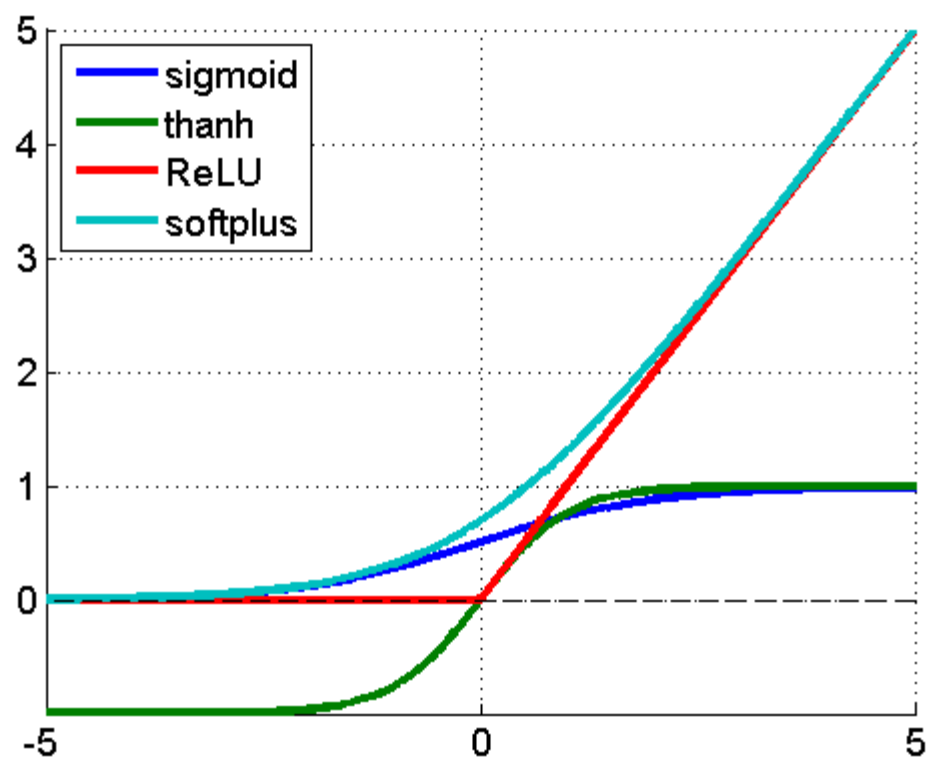
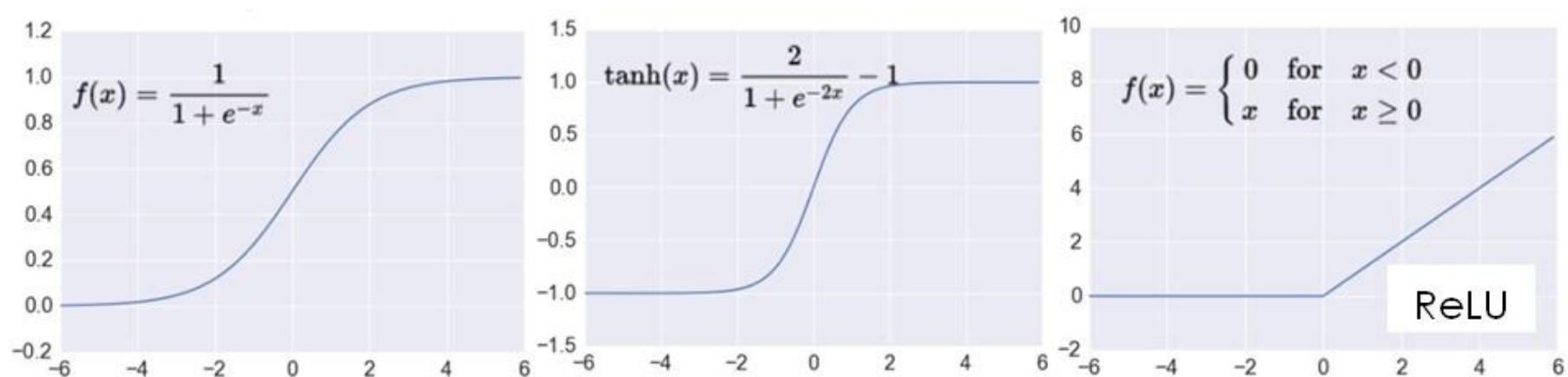


Feature visualization of convolutional net trained on ImageNet  
(Zeiler and Fergus, 2013)

# Vanishing Gradient Problem







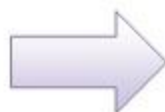


1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Input matrix

1	0	1
0	1	0
1	0	1

Convolutional  
3x3 filter



1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

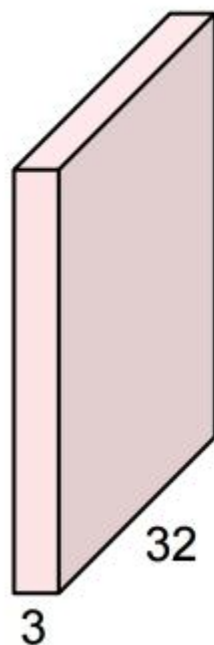
Image

4		

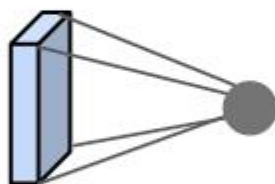
Convolved  
Feature

32x32x3 image

activation map

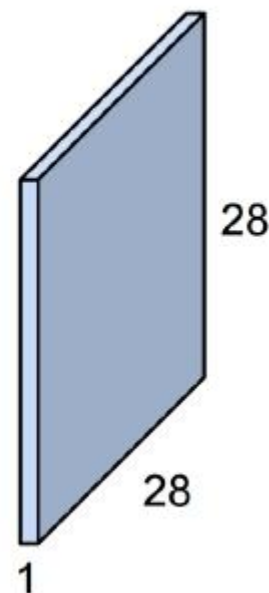
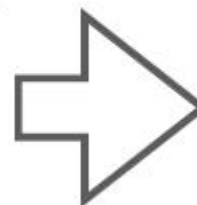


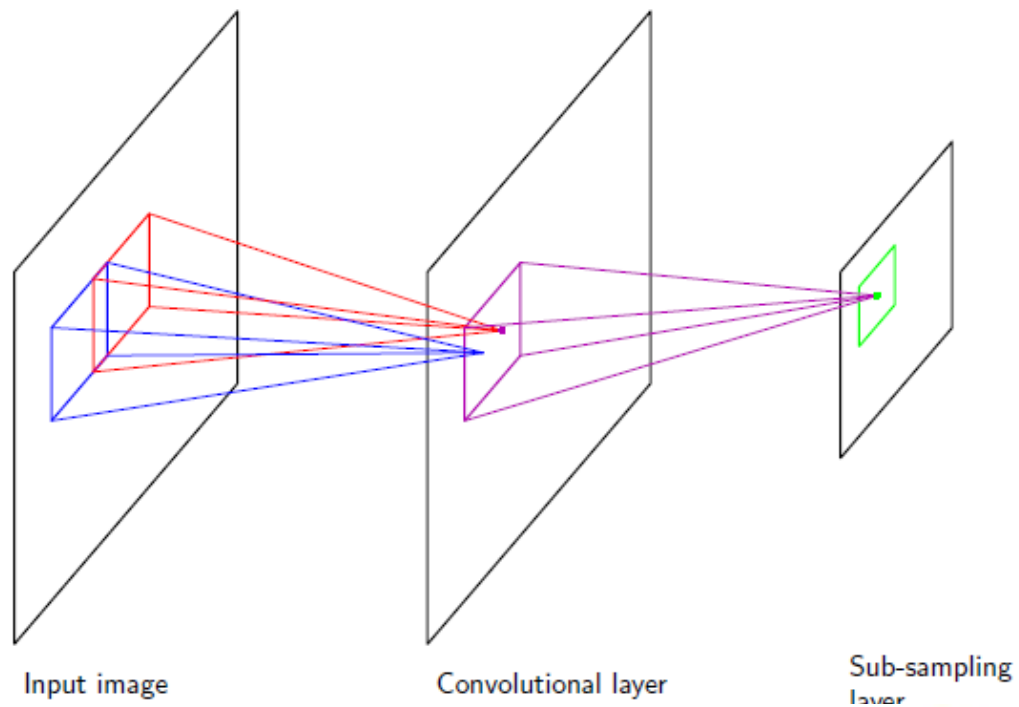
5x5x3 filter



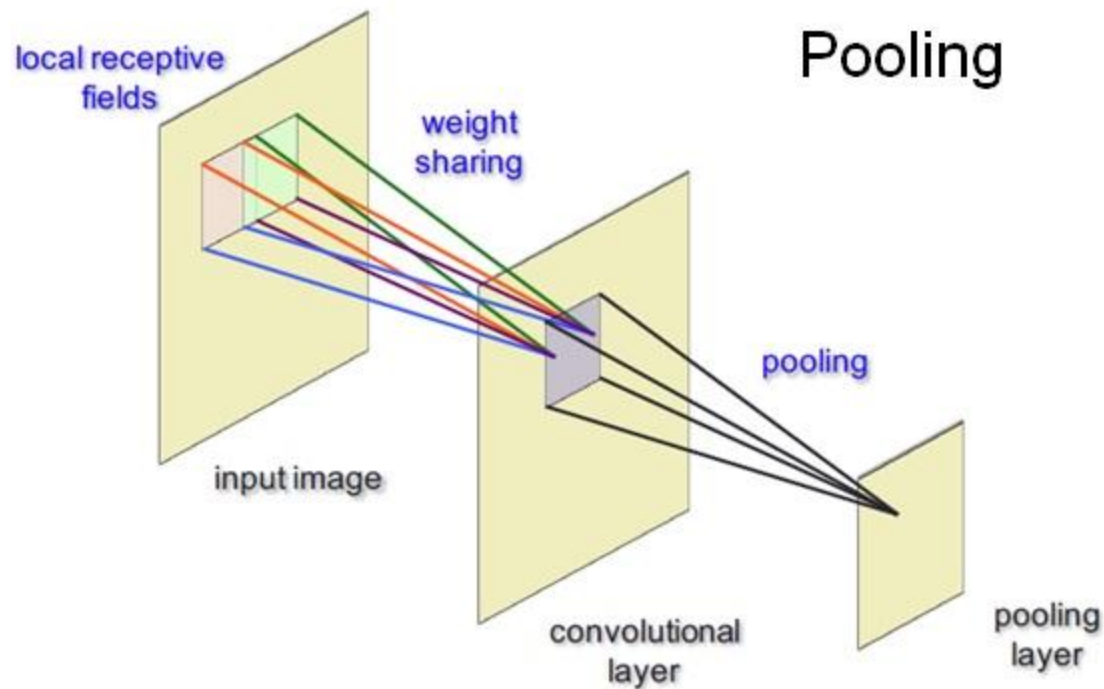
1 number to  
represent result of  
filter with small  
chunk of image

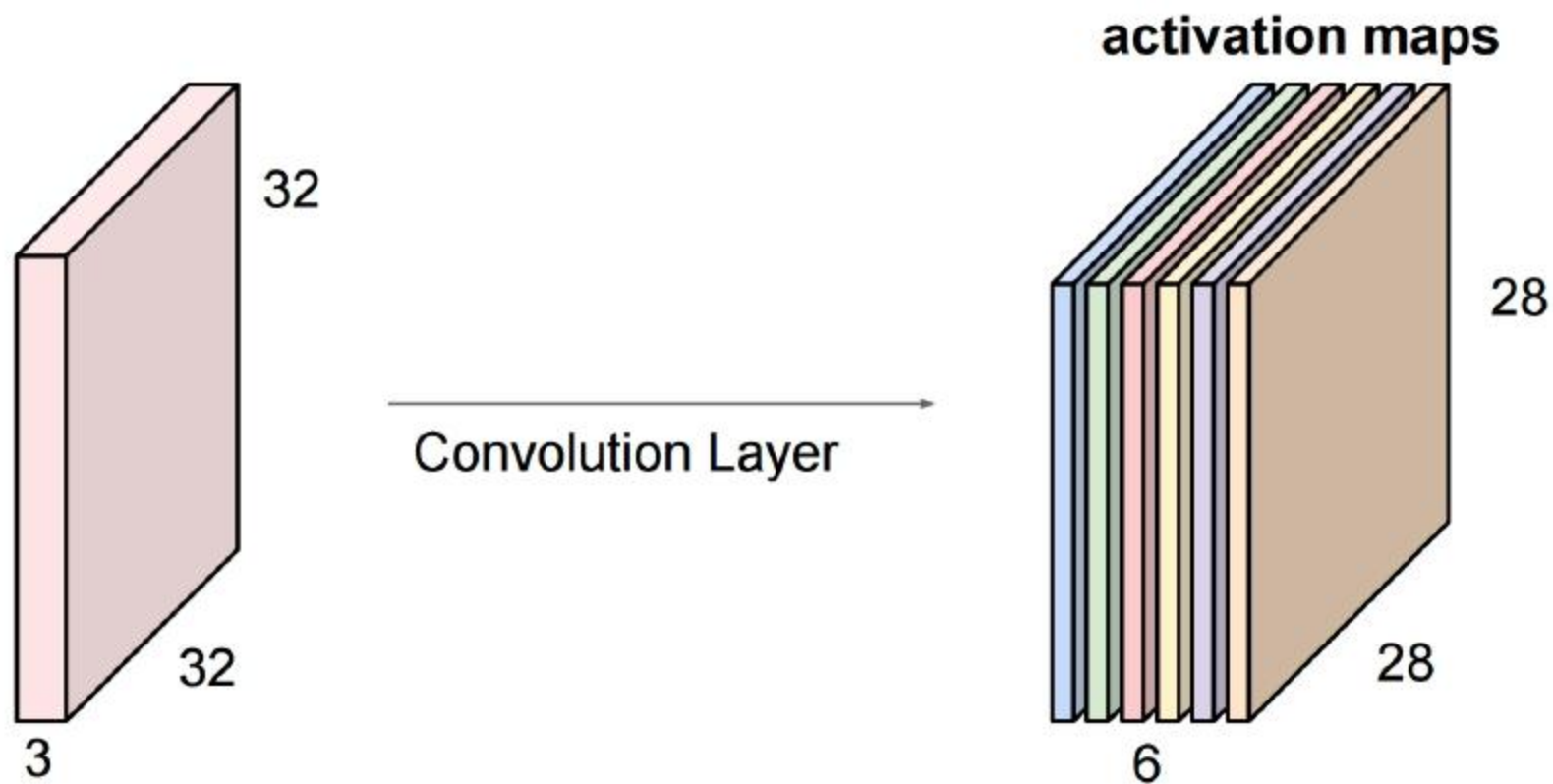
Convolve the filter with the image  
(i.e., slide over image spatially  
computing dot products)





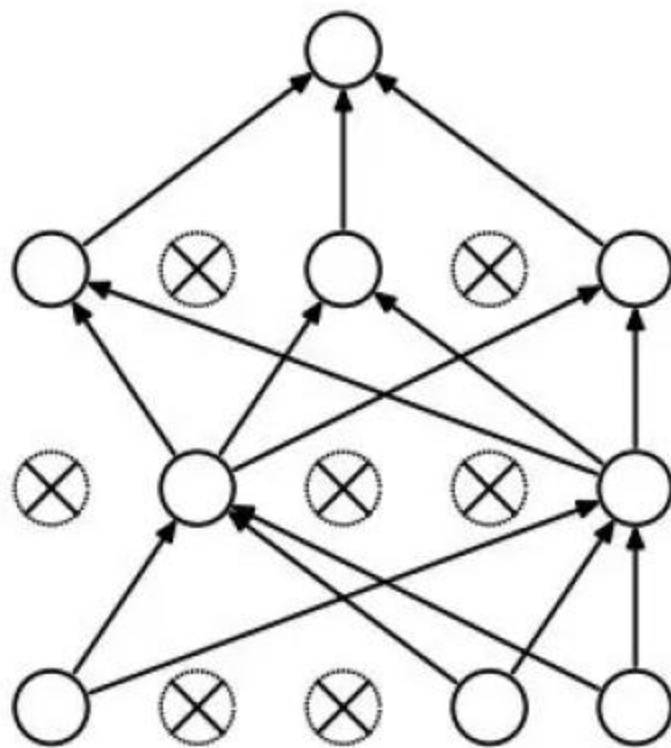
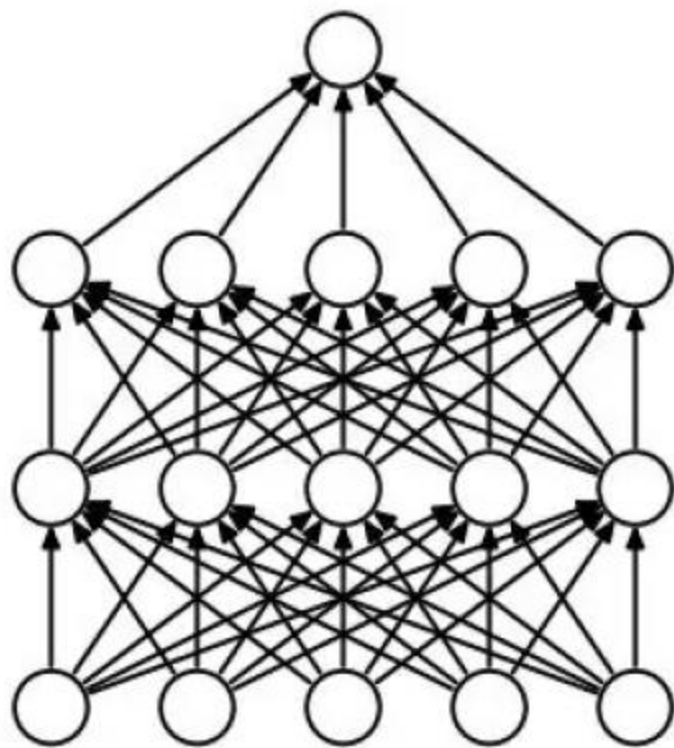
**Figure 5.17** Diagram illustrating part of a convolutional neural network. Convolutional units followed by a layer of subsampling units may be used.





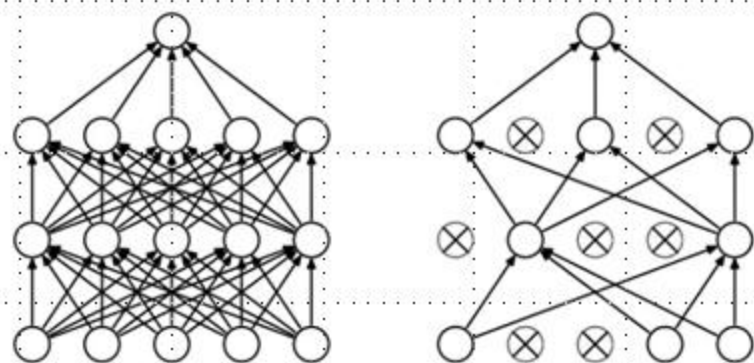
Stacking up multiple filters yields “new image”

# Dropout





# Regularization



## Dropout

- Randomly drop units (along with their connections) during training
- Each unit retained with fixed probability  $p$ , independent of other units
- Hyper-parameter  $p$  to be chosen (tuned)

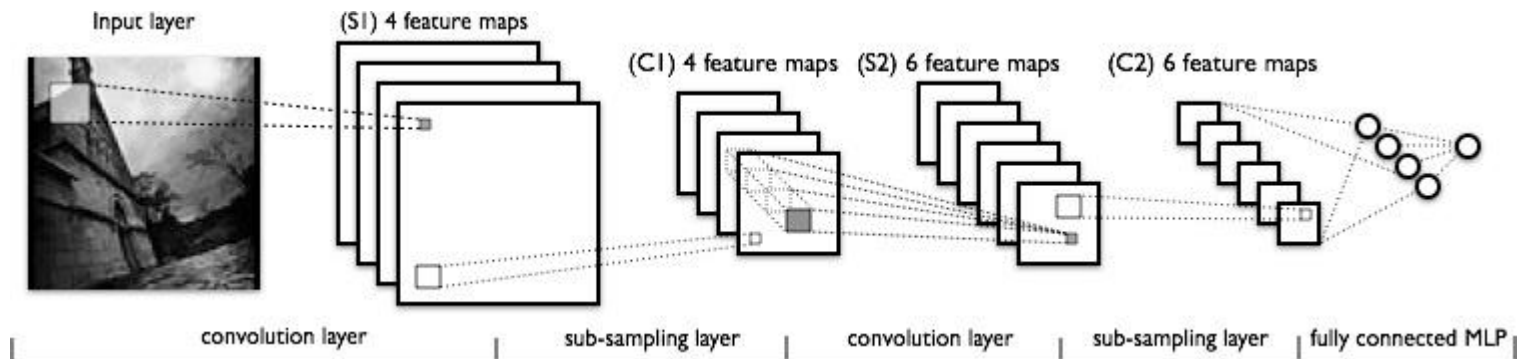
## L2 = weight decay

- Regularization term that penalizes big weights, added to the objective
- Weight decay value determines how dominant regularization is during gradient computation
- Big weight decay coefficient  $\rightarrow$  big penalty for big weights

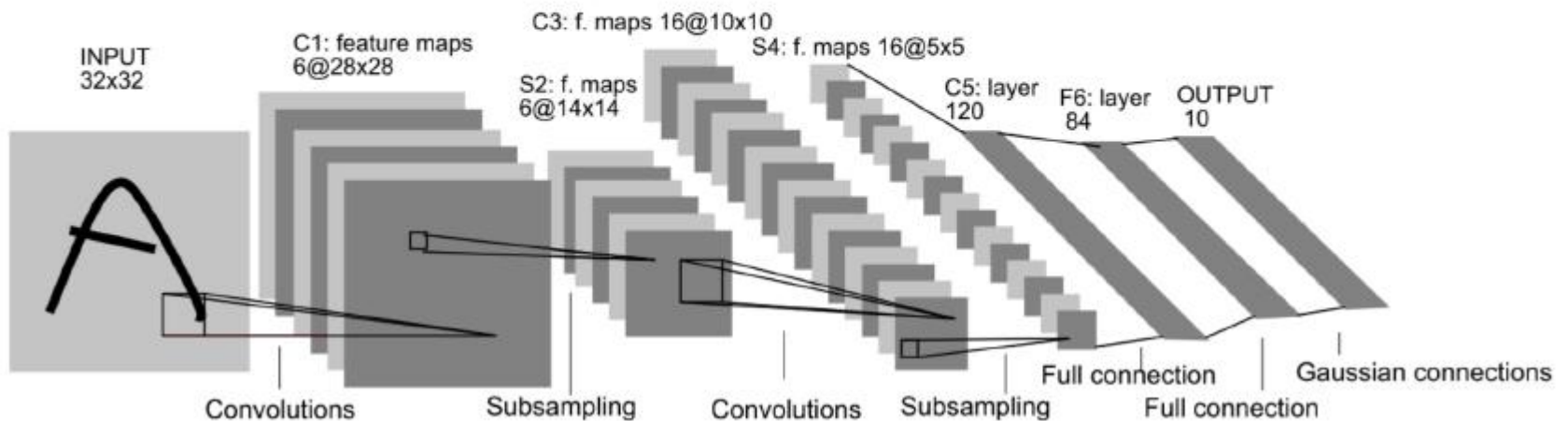
$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

## Early-stopping

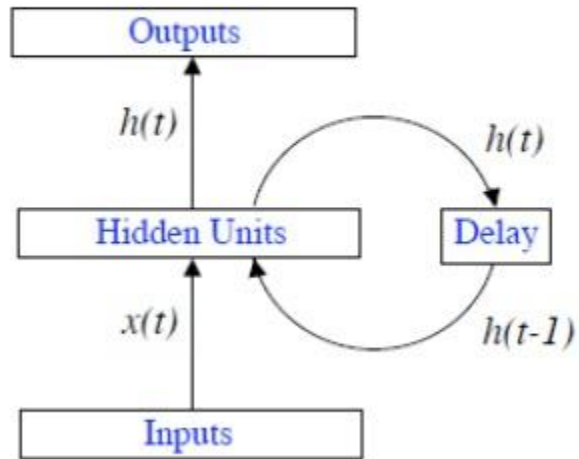
- Use validation error to decide when to stop training
- Stop when monitored quantity has not improved after  $n$  subsequent epochs



# LeNet 5 [LeCun et al., 1998]



## Recurrent neural network



## Back Propagation Through Time (BPTT)

