

# NEURAL NETWORK ESSENTIAL

Ke Chen

Department of Computer Science, The University of Manchester

*Ke.Chen@manchester.ac.uk*

## INTRODUCTION

History and application

## NEURON MODEL

Neuron, artificial neuron model, activation (transfer) functions

## NEURAL ARCHITECTURE

Feedforward versus recurrent, fully connected versus partially connected, homogeneous components versus heterogeneous components

## LEARNING

Loss functions, stochastic gradient optimisation, back-propagation (BP) algorithm, practical issues

# INTRODUCTION

- In **1943**, **Neural networks (NNs)** were originated by **W. McCulloch & W. Pitts** who created computational models to simulate neurons in human brain.
- In **1957**, **F. Rosenblatt** proposed the first ever learning algorithm named **perceptron** to train artificial NNs for binary classification tasks. (**1st wave of neural networks**)
- In **1969**, **M. Minsky & S. Pappert** wrote a book entitled **perceptrons** that pointed out the limitation due to a lack of learning algorithm for multilayered NNs.
- In **1986**, **Parallel Distributed Processing (PDP)** Project led to several seminal works in neural computation where the most **influential** one is **back-propagation** learning developed by **D. Rumelhart, G. Hinton & R. Williams**. (**2nd wave of neural networks**)
- From **end of 1990s** to **2006**, difficulties in training **deep NNs** diverted ML research to simple yet theoretically justified learning models, e.g. SVM and Adaboost.
- In **2006**, **G. Hinton & his students** proposed a **new learning strategy** to train deep NNs and further coined the term **deep learning** to replace **neural networks**.
- Since **2006**, ML has shifted its focus to deep learning, which lifts **AI** to a new era. In **2018**, **G. Hinton, Y. Bengio & Y. LeCun** received **ACM Alan Turing Award** for their contributions in deep learning. (**3rd wave of neural networks**)

As an **underpinning** technology, **deep learning** has been applied to many AI domains.

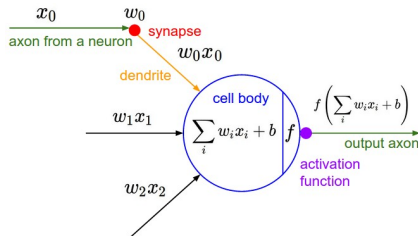
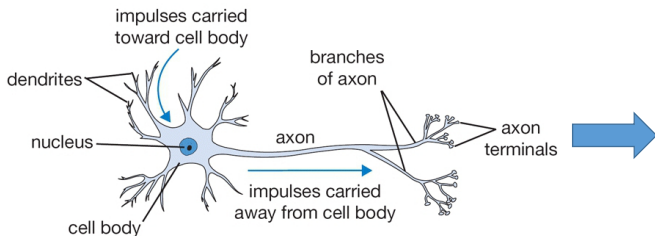
- **Computer vision**: deep learning has become a **pre-dominated** techniques and led to **super-human** performance in some visual recognition tasks.
- **Speech recognition**: as one of the most important **technical components**, deep learning dramatically improves **recognition accuracy**, e.g. Google voice recogniser.
- **Natural language processing**: as one of the most important **technical components**, deep learning has substantially improve **information retrieval** and **machine translation**, e.g., Google translate.
- **Game agent**: as one of the most important **technical components**, deep learning has created game agents outperforming human beings, e.g., Alpha Go and Atari agents.
- **Miscellaneous**: deep learning has played a crucial role in many real applications such as **industry 4.0 manufacture automation** and **medical diagnosis and treatment**.

# NEURON MODEL

- **Biological neuron** has been well studied in **biology** and **neuroscience**.
- **Computational neuron model** may be **biologically plausible** or **artificial**.
- **Biologically plausible model**: modelling all biological mechanisms and functions via differential equation system, e.g., Hodgkin–Huxley model for spiking neuron
- **Artificial model**: modelling main functions abstractly by ignoring biological meaning

$$a = \mathbf{w}^T \mathbf{x} + b = \sum_i w_i x_i + b, \quad \text{output} = f(a),$$

$w_i$ : weights,  $b$ : bias and  $x_i$ : input,  $a$ : action potential,  $f(\cdot)$ : activation function



## Activation Function

- Step (perceptron) function

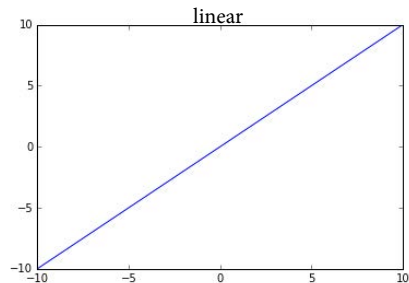
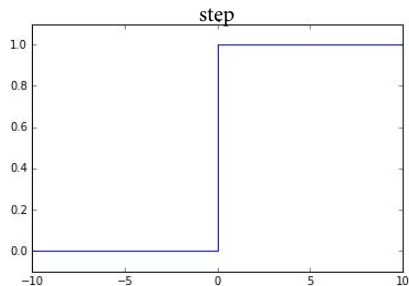
$$f(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

Step function is **discontinuous** hence has no gradient.

- Linear (identity) function

$$f(x) = x$$

$$\frac{df(x)}{dx} = 1.$$



## Activation Function

- Sigmoid (logistic) function

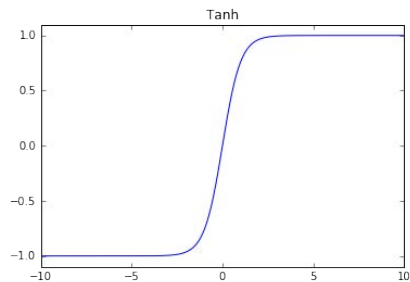
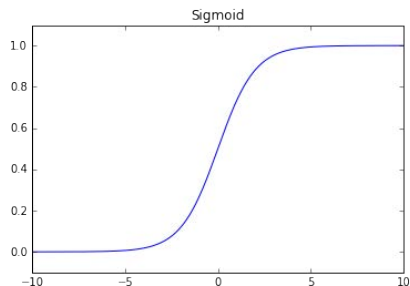
$$f(x) = \frac{1}{1 + e^{-x}}.$$

$$\frac{df(x)}{dx} = f(x)(1 - f(x)).$$

- Hyperbolic tangent function

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}.$$

$$\frac{df(x)}{dx} = 1 - f^2(x).$$



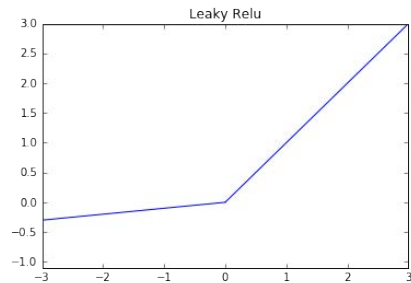
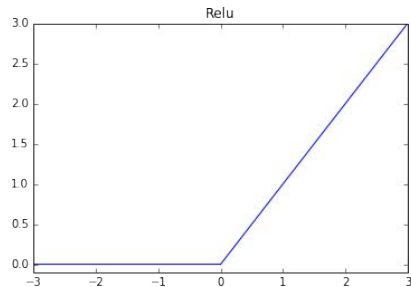
## Activation Function

- Rectified linear unit (ReLU) function

$$f(x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases}$$
$$\frac{df(x)}{dx} = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

- Leaky ReLU function

$$f(x) = \begin{cases} 0.1x & x < 0 \\ x & x \geq 0 \end{cases}$$
$$\frac{df(x)}{dx} = \begin{cases} 0.1 & x < 0 \\ 1 & x \geq 0 \end{cases}$$





## Activation Function

- Scaled-exponential linear unit (SeLU) function

$$f(\alpha, x) = \begin{cases} \alpha (e^x - 1) & x < 0 \\ x & x \geq 0 \end{cases}, \quad \frac{df(x)}{dx} = \begin{cases} f(\alpha, x) + \alpha & x < 0 \\ 1 & x \geq 0 \end{cases}$$

- Maxout function

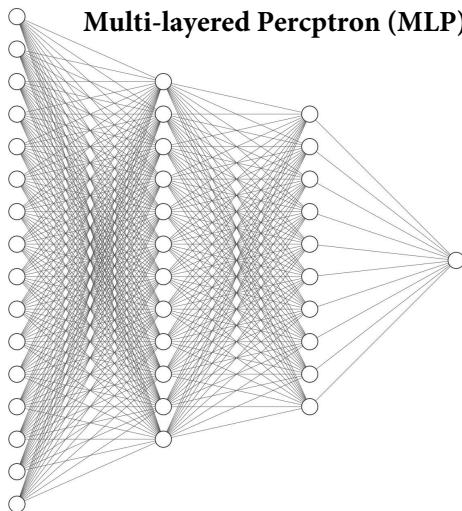
$$f(\mathbf{x}) = \max_{i=1}^{|\mathbf{x}|} x_i, \quad \frac{df(\mathbf{x})}{dx_j} = \begin{cases} 1 & j = \operatorname{argmax}_{i=1}^{|\mathbf{x}|} x_i \\ 0 & j \neq \operatorname{argmax}_{i=1}^{|\mathbf{x}|} x_i \end{cases}$$

- Softmax function

$$f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{|\mathbf{x}|} e^{x_j}}, \quad \frac{df(x_i)}{dx_j} = f(x_i) (\delta(i, j) - f(x_i)) \quad i = 1, 2, \dots, |\mathbf{x}|.$$

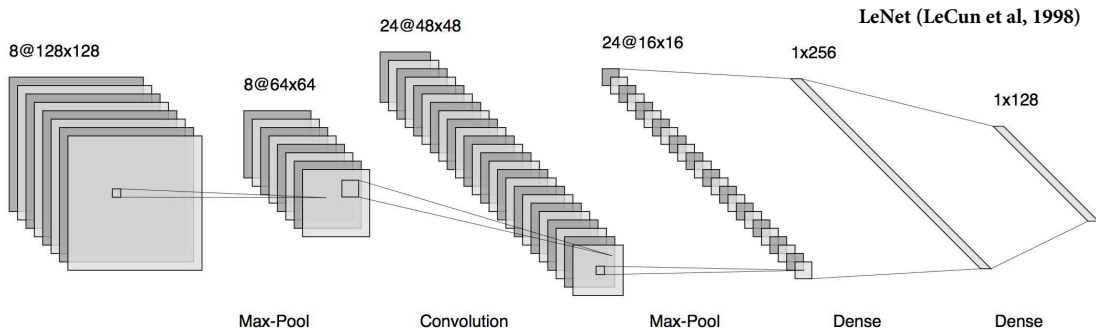
## Feedforward Neural Networks

- **Fully-connected**: all neurons in one layer connected to all in its successive layers
- **Homogeneous**: all neurons apart from those in input/output layer are the same.



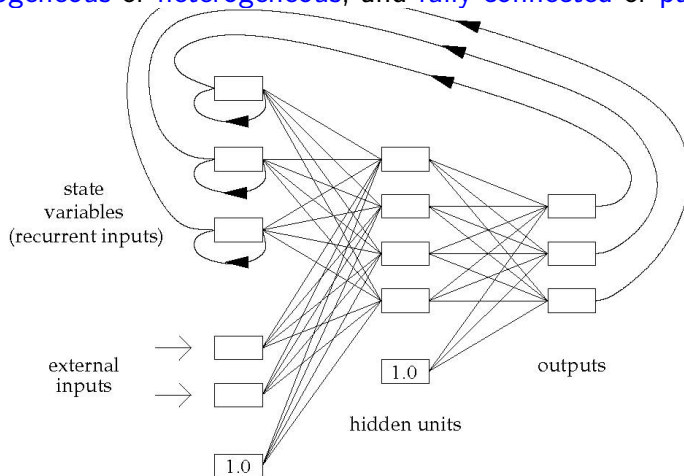
## Feedforward Neural Networks

- **Partially-connected**: neurons in a layer only connected to some in its successive layer
- **Heterogeneous**: neurons in different layers are various for different purposes
- For example, **Convolutional Neural Networks (CNNs)**



## Recurrent Neural Networks

- **Recurrent-connected**: neurons with feedback connections to neurons in previous and/or the same layers (lateral connection)
- May be **homogeneous** or **heterogeneous**, and **fully-connected** or **partially-connected**



## Loss Function

Given a training dataset,  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{|\mathcal{D}|}$ , **loss functions** are used to train NNs with parameters (weights and bias),  $\Theta$  (a collective notation of all parameters).

- **Mean squared error (MSE) loss** for regression

$$\mathcal{L}(\Theta; \mathcal{D}) = \frac{1}{2|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|^2.$$

$\hat{\mathbf{y}}_i$ : output of NNs for input  $\mathbf{x}_i$  and **linear activation function** used in **output layer**

- **Cross-entropy loss** for binary classification

$$\mathcal{L}(\Theta; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \left( y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \right), \quad y_i \in \{0, 1\}.$$

$\hat{y}_i$ : output of NNs for input  $\mathbf{x}_i$  and **sigmoid activation function** used in **output layer**

## Loss Function

Given a training dataset,  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{|\mathcal{D}|}$ , loss functions are used to train NNs with parameters (weights and bias),  $\Theta$  (a collective notation of all parameters).

- **Categorical cross-entropy loss** for  $C$ -class classification ( $C > 2$ )

$$\mathcal{L}(\Theta; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{j=1}^C y_{ij} \log \hat{y}_{ij}, \quad \mathbf{y}_i = \{y_{i1}, y_{i2}, \dots, y_{iC}\}, \quad \hat{\mathbf{y}}_i = \{\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iC}\}.$$

$\hat{\mathbf{y}}_i$ : output of NNs for input  $\mathbf{x}_i$  and **softmax activation function** used in **output layer**

- **Regularised loss** for generalisation and exclusion of ill-posed solution

$$\mathcal{L}_R(\Theta; \mathcal{D}) = \mathcal{L}(\Theta; \mathcal{D}) + \lambda \mathcal{R}(\Theta),$$

where  $\mathcal{R}(\Theta)$  is a **regularisation penalty** and  $\lambda$  is a **trade-off coefficient**. For instance, **weight decay** is often used for regularisation,  $\mathcal{R}(\Theta) = \frac{1}{2} \sum_{\text{all } \mathbf{w}} \|\mathbf{w}\|^2$ .

## • Stochastic Gradient Optimisation

- ① Randomly initialise all the parameters,  $\Theta_0$ .
- ② In each epoch, randomly split a training dataset,  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{|\mathcal{D}|}$ , into  $K$  mini-batches,  $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K$ , with equal size  $|\mathcal{B}|$ , e.g.,  $|\mathcal{B}| = 16$ .
- ③ Update all parameters via gradient descent on a mini-batch basis

$$\Theta_{k+1} \leftarrow \Theta_k - \frac{\eta}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{k+1}} \nabla_{\Theta} \mathcal{L}(\Theta; \mathbf{x}, \mathbf{y})|_{\Theta=\Theta_k}, \quad k = 0, 1, \dots, K$$

where  $0 < \eta < 1$  is a learning rate.  $\Theta_0$  refers to the one after last epoch.

- ④ Repeat steps 2 and 3 until a stopping condition is satisfied.
- Add momentum, e.g.,  $\Delta\Theta_k = \Theta_k - \Theta_{k-1}$ , to update rule to speed up training

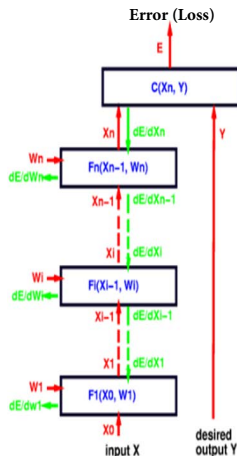
$$\Theta_{k+1} \leftarrow \Theta_k - \frac{\eta_1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{k+1}} \nabla_{\Theta} \mathcal{L}(\Theta; \mathbf{x}, \mathbf{y})|_{\Theta=\Theta_k} + \eta_2 \Delta\Theta_k, \quad k = 0, 1, \dots, K$$

where  $0 < \eta_1 < 1$  and  $0 < \eta_2 < 1$  are different learning rates.

**Fact:** Optimisation package, Adam, can compute gradient for any loss of DNNs automatically.

## Back-propagation Procedure

- Main stages: **forward propagation**, **backward gradient propagation**, **parameter update**
- Compute **gradients** of a **cost (loss)** function with respect to parameters, **weights and biases**, associated with different layers by applying **chain rule** recursively



LeCun et al (1998)

$$\frac{\partial E}{\partial X_n} = \frac{\partial C(X_n, Y)}{\partial X_n}$$

$$\frac{\partial E}{\partial X_{n-1}} = \frac{\partial E}{\partial X_n} \frac{\partial F_n(X_{n-1}, W_n)}{\partial X_{n-1}}$$

$$\frac{\partial E}{\partial W_n} = \frac{\partial E}{\partial X_n} \frac{\partial F_n(X_{n-1}, W_n)}{\partial W_n}$$

$$\frac{\partial E}{\partial X_{n-2}} = \frac{\partial E}{\partial X_{n-1}} \frac{\partial F_{n-1}(X_{n-2}, W_{n-1})}{\partial X_{n-2}}$$

$$\frac{\partial E}{\partial W_{n-1}} = \frac{\partial E}{\partial X_{n-1}} \frac{\partial F_{n-1}(X_{n-2}, W_{n-1})}{\partial W_{n-1}}$$

....etc, until we reach the first module.

we now have all the  $\frac{\partial E}{\partial W_i}$  for  $i \in [1, n]$ .



## Back-propagation Algorithm

### • Forward propagation

An MLP has  $L$  hidden layers with neurons of activation function,  $f(\cdot)$  and output layer of activation function,  $g(\cdot)$ . For training example,  $(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{k+1}$  (mini-batch),

- ① Input:  $\mathbf{h}^{(0)} \leftarrow \mathbf{x}$ ; (input layer viewed as layer 0)
- ② Compute activation of neurons in hidden layers: for  $l = 1, 2, \dots, L$ , compute

$$\mathbf{a}^{(l)}(\mathbf{x}) = W_k^{(l)} \mathbf{h}^{(l-1)}(\mathbf{x}) + \mathbf{b}_k^{(l)}, \quad \mathbf{h}^{(l)} \leftarrow \mathbf{f}(\mathbf{a}^{(l)}(\mathbf{x})).$$

- ③ Compute activation of output units: for output layer viewed as layer  $L + 1$ ,

$$\mathbf{a}^{(L+1)}(\mathbf{x}) = W_k^{(L+1)} \mathbf{h}^{(L)}(\mathbf{x}) + \mathbf{b}_k^{(L+1)}, \quad \hat{\mathbf{y}} = \mathbf{h}^{(L+1)} \leftarrow \mathbf{g}(\mathbf{a}^{(L+1)}(\mathbf{x})).$$

- ④ Compute loss of this training example:  $\mathcal{L}(\Theta_k; \mathbf{x}, \mathbf{y})$ , where  $\Theta_k = \{W_k^{(l)}, \mathbf{b}_k^{(l)}\}_{l=1}^L$ .  
(Step 4 is optional but required by an early stop.)

## Back-propagation Algorithm

- **Backward gradient propagation**

- ① Compute gradient at output layer:

$$\delta^{(L+1)}(\mathbf{x}, \mathbf{y}) \leftarrow \frac{\partial \mathcal{L}(\Theta_k; \mathbf{x}, \mathbf{y})}{\partial \mathbf{a}^{(L+1)}(\mathbf{x})} = \frac{\partial \mathcal{L}(\Theta_k; \mathbf{x}, \mathbf{y})}{\partial \mathbf{h}^{(L+1)}(\mathbf{x})} \odot \mathbf{g}'(\mathbf{a}^{(L+1)}(\mathbf{x})).$$

- ② Compute gradient at different hidden layers: for  $l = L, L-1, \dots, 1$ , compute

$$\frac{\partial \mathcal{L}(\Theta_k; \mathbf{x}, \mathbf{y})}{\partial \mathbf{h}^{(l)}(\mathbf{x})} \leftarrow \left( W_k^{(l+1)} \right)^T \delta^{(l+1)}(\mathbf{x}, \mathbf{y}), \quad \delta^{(l)}(\mathbf{x}, \mathbf{y}) \leftarrow \frac{\partial \mathcal{L}(\Theta_k; \mathbf{x}, \mathbf{y})}{\partial \mathbf{h}^{(l)}(\mathbf{x})} \odot \mathbf{f}'(\mathbf{a}^{(l)}(\mathbf{x})).$$

- **Update parameters on mini-batch:** for  $l = L, L-1, \dots, 0$ ,

$$\mathbf{b}_{k+1}^{(l+1)}(\mathbf{x}) \leftarrow \mathbf{b}_k^{(l+1)}(\mathbf{x}) - \frac{\eta}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{k+1}} \delta^{(l+1)}(\mathbf{x}, \mathbf{y}),$$

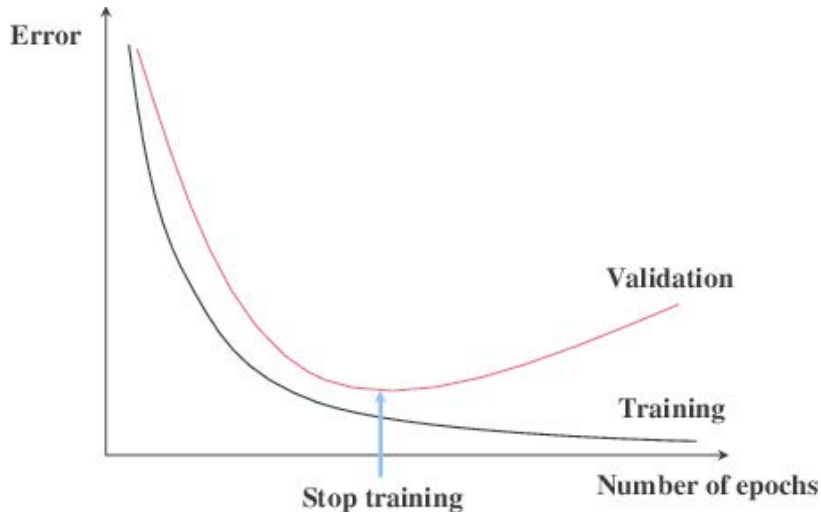
$$W_{k+1}^{(l+1)} \leftarrow W_k^{(l+1)} - \frac{\eta}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{k+1}} \delta^{(l+1)}(\mathbf{x}, \mathbf{y}) \left( \mathbf{h}^{(l)}(\mathbf{x}) \right)^T.$$

## Practical Issues

- **Initialisation**: randomly initialise weights with small values, e.g.,  $U[-\frac{1}{\sqrt{|\mathbf{x}|}}, \frac{1}{\sqrt{|\mathbf{x}|}}]$
- **Hyper-parameter issues**: many hyper-parameters to be set properly before training
  - **Architectural/structural**: number of hidden layers, number of hidden neurons/layer, other parameters in chosen activation functions
  - **learning-related**: learning rate(s), trade-off coefficient for regularisation, mini-batch size and so on
- **Hyper-parameter tuning**: **grid search** or **random search** from a range of values
- **Model selection and evaluation**
  - Main method: **held-out validation** and  **$K$ -fold cross validation**
  - During learning, use the performance on **validation sets** to find out **optimal hyper-parameter values** and decide the **early stopping** to avoid **overfitting**

## Practical Issues

- **Early stopping**: effective measure to avoid **over-fitting**



If you want to deepen your understanding and learn something beyond this lecture, you can self-study the optional references below.

[Goodfellow et al., 2016] Goodfellow I., Bengio Y., and Courville A. (2016): *Deep Learning*, MIT Press. (Chapter 6 & Sections 11.1-11.5)

[Schmidhuber, 2015] Schmidhuber J. (2015): Deep learning in neural networks: An overview. *Neural Networks*, Vol. 61, pp. 85-117.