

# Pattern Information Processing:<sup>153</sup> Neural Networks

Masashi Sugiyama  
(Department of Computer Science)

Contact: W8E-505

[sugi@cs.titech.ac.jp](mailto:sugi@cs.titech.ac.jp)

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

# Linear/Non-Linear Models

154

- **Linear model:**  $\hat{f}(x)$  is linear with respect to  $\alpha$   
(Note: not necessarily linear with respect to  $x$ )

$$\hat{f}(x) = \sum_{i=1}^b \alpha_i \varphi_i(x)$$

- **Non-linear model:** Otherwise

# Today's Plan

155

- Neural networks
- Least-squares in neural networks:  
Error back-propagation algorithm

# Non-Linear Models

156

- A popular choice: **Hierarchical models**

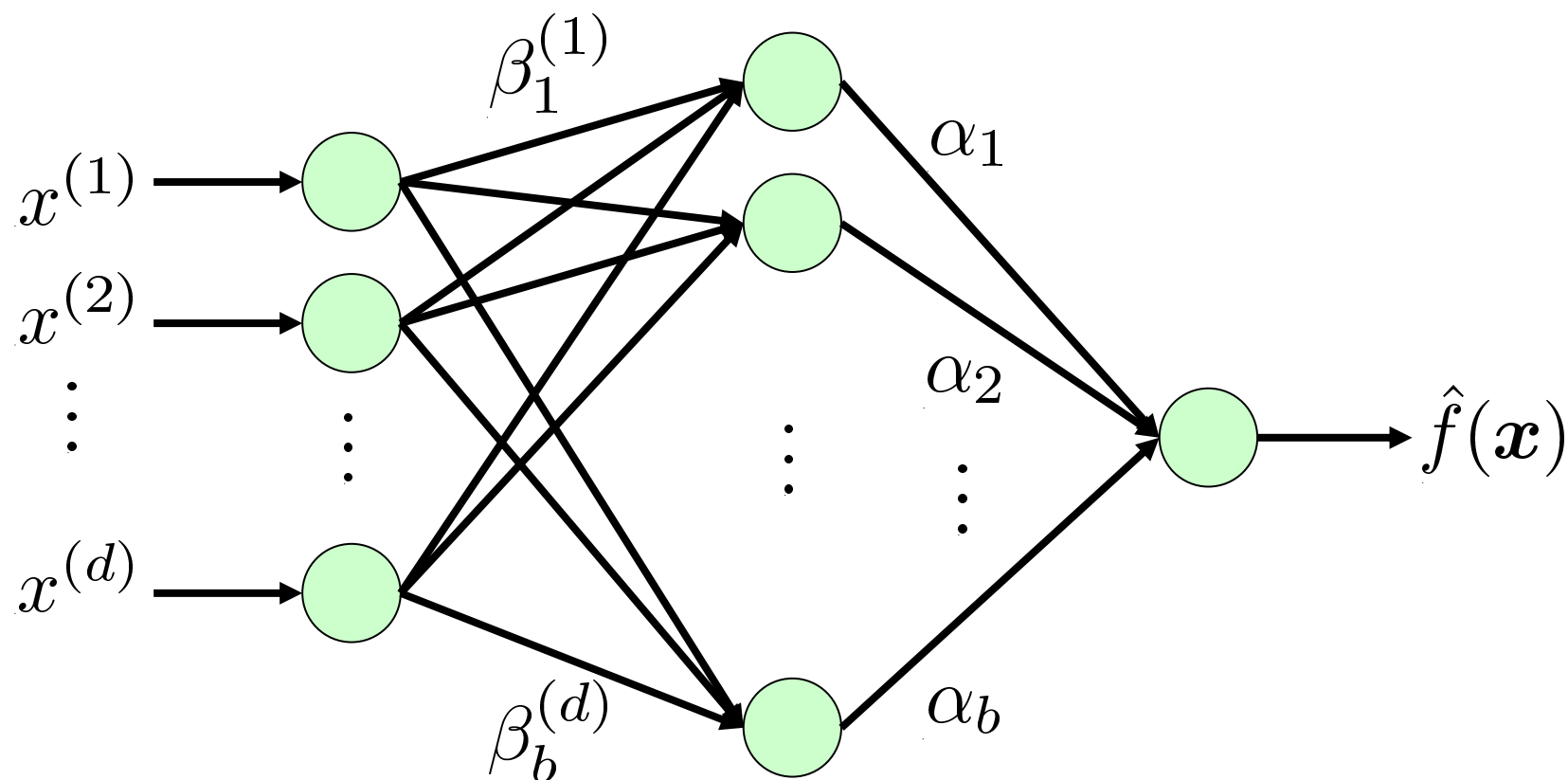
$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi(\mathbf{x}; \beta_i)$$

- Basis function is parameterized by  $\beta_i$

# Three-Layer Networks

157

- Such a hierarchical model can be represented as a **3-layer network**.

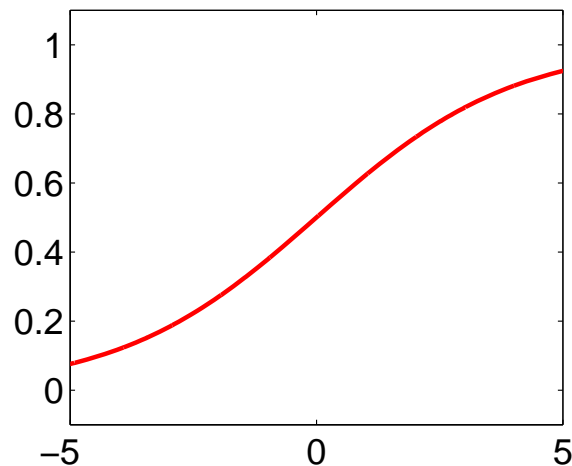


# Sigmoidal Function

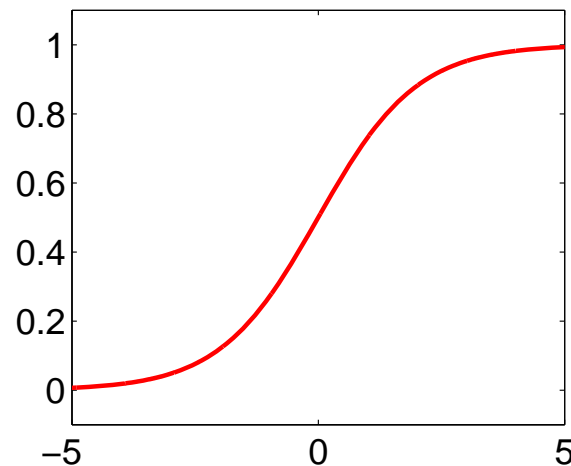
158

- A typical basis function: Sigmoidal function

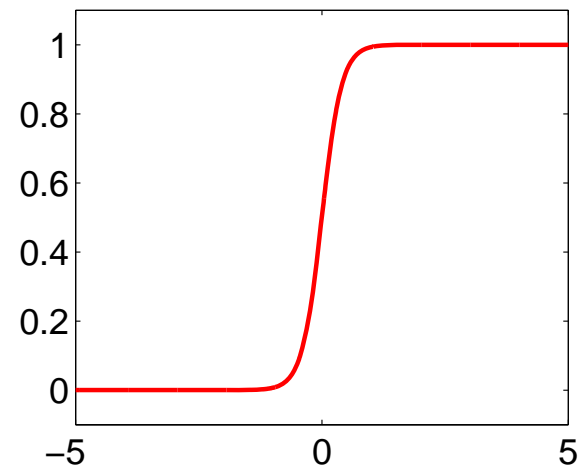
$$\varphi(\mathbf{x}; \boldsymbol{\beta}, h) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \boldsymbol{\beta} \rangle - h)}$$



$\beta$  : Small



$\beta$  : Medium



$\beta$  : Large

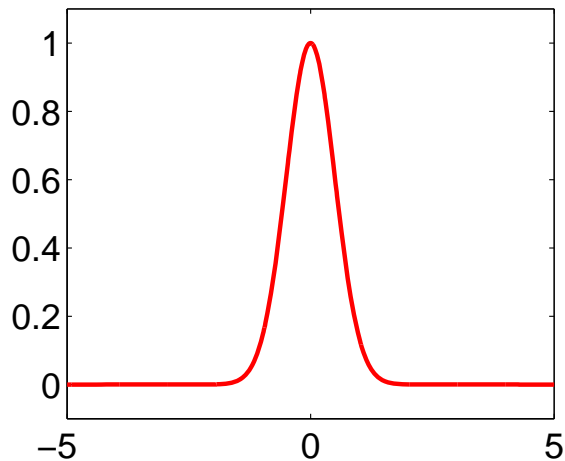
# Perceptrons

- The behavior of the sigmoidal functions is similar to the **neurons in the brain**.
- For this reason, hierarchical models with sigmoidal functions are called **artificial neural networks** or **perceptrons**.
- Mathematically, 3-layer neural networks can approximate any continuous functions with arbitrary small error (“**universal approximator**”).

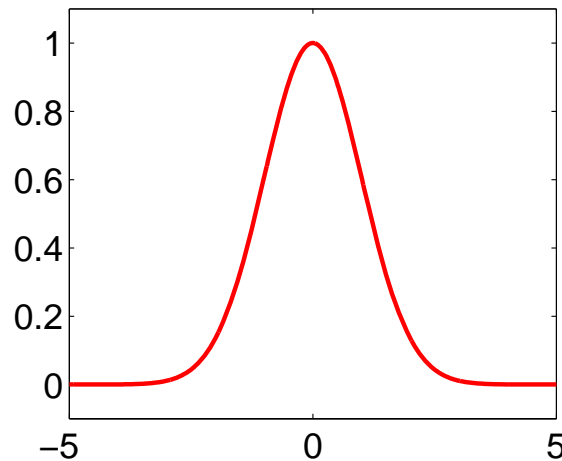
# Gaussian Radial Basis Function<sup>160</sup>

- Another popular basis function:  
Gaussian radial basis function

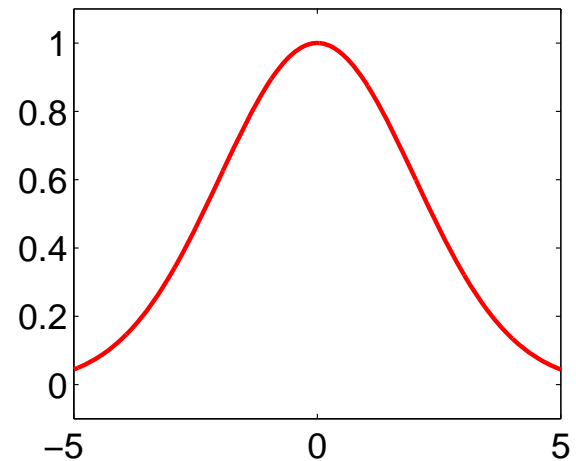
$$\varphi(\mathbf{x}; \boldsymbol{\beta}, c) = \exp \left( -\frac{\|\mathbf{x} - \boldsymbol{\beta}\|^2}{2c^2} \right)$$



$c$  : Small



$c$  : Medium



$c$  : Large



# Least-Squares Learning

161

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi(\mathbf{x}; \beta_i)$$

$$\mathbf{w} = (\alpha^\top, \beta_1^\top, \beta_2^\top, \dots, \beta_b^\top)^\top.$$

- Least-squares learning is often used for training hierarchical models.

$$\min_{\mathbf{w}} J_{LS}(\mathbf{w})$$

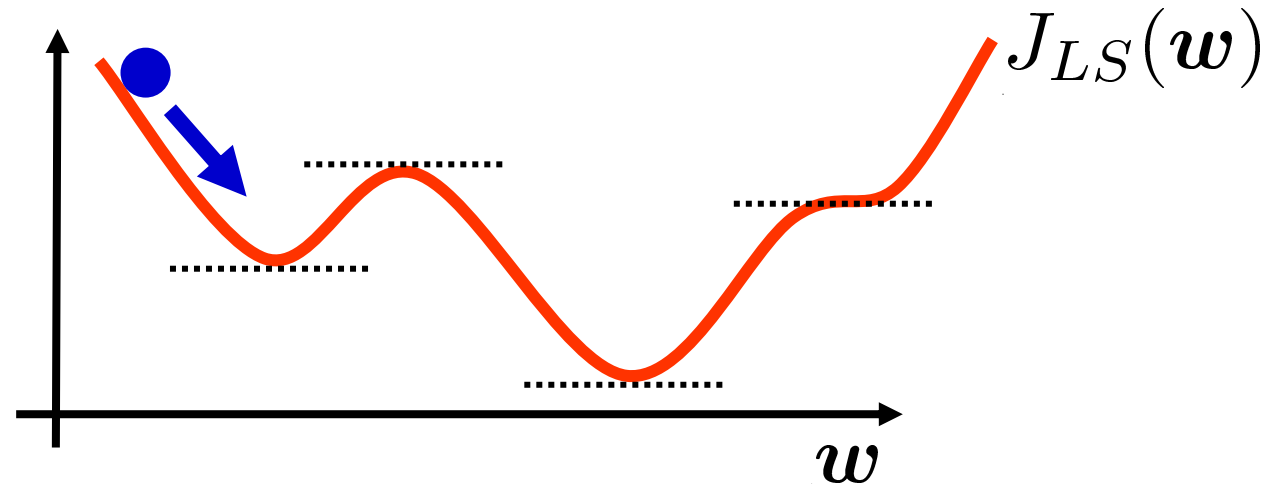
$$J_{LS}(\mathbf{w}) = \sum_{i=1}^n \left( \hat{f}(\mathbf{x}_i) - y_i \right)^2$$

# How to Obtain Solutions

162

- No analytic solution is known.
- Simple gradient search is usually used.

$$\hat{w}^{new} \leftarrow \hat{w}^{old} - \varepsilon \nabla J_{LS}(\hat{w}^{old})$$



- It converges to one of the local minima.

# Error Back-Propagation

- Efficient calculation of gradient for Sigmoidal basis functions

$$\frac{\partial J_{LS}}{\partial \alpha_j} = 2 \sum_{i=1}^n o_{i,j} \left( \hat{f}(\mathbf{x}_i) - y_i \right)$$

$$\frac{\partial J_{LS}}{\partial \beta_j^{(k)}} = 2\alpha_j \sum_{i=1}^n o_{i,j} (1 - o_{i,j}) x_i^{(k)} \left( \hat{f}(\mathbf{x}_i) - y_i \right)$$

$$\frac{\partial J_{LS}}{\partial h_j} = 2\alpha_j \sum_{i=1}^n o_{i,j} (1 - o_{i,j}) \left( \hat{f}(\mathbf{x}_i) - y_i \right)$$

$$o_{i,j} = \varphi(\mathbf{x}_i; \boldsymbol{\beta}_j, h_j)$$

# Error Back-Propagation (cont.)<sup>164</sup>

- When the output values of the network are calculated, the input points are propagated following the forward path.
- On the other hand, when the gradients are calculated, the output error  $(\hat{f}(x_i) - y_i)$  is propagated backward.
- For this reason, this algorithm is called the error back-propagation.
- However, it is gradient descent so global convergence is not guaranteed.

# Stochastic Gradient Descent<sup>165</sup>

- In the usual gradient method, all training examples are used **at the same time**.
- In practice, the following **stochastic method** would be computationally advantageous.

- Randomly choose one of the training examples (say,  $(\mathbf{x}_i, y_i)$ )

- Update the parameter vector by

$$\hat{\mathbf{w}}^{new} \leftarrow \hat{\mathbf{w}}^{old} - \varepsilon \nabla J_i(\hat{\mathbf{w}}^{old})$$

$$J_i(\mathbf{w}) = (\hat{f}(\mathbf{x}_i) - y_i)^2$$

- Repeat this procedure until convergence.

# Avoiding Overfitting

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi(\mathbf{x}; \beta_i)$$

$$\mathbf{w} = (\alpha^\top, \beta_1^\top, \beta_2^\top, \dots, \beta_b^\top)^\top.$$

■ LS overfits to noisy samples.

- **Regularization:**

$$\min_{\mathbf{w}} [J_{LS}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2]$$

$$J_{LS}(\mathbf{w}) = \sum_{i=1}^n \left( \hat{f}(\mathbf{x}_i) - y_i \right)^2 \quad \lambda > 0$$

- **Early stopping:** Stop gradient descent before it converges

# Special Structure of Neural Networks

- Parameters and functions are not one-to-one mapping.

$$\hat{f}(x) = \sum_{i=1}^b \alpha_i \varphi(x; \beta_i)$$

- If  $\alpha_i = 0$ , any  $\beta_i$  gives the same function.
- Any permutation gives the same function.
- This **non-identifiability** causes great difficulty in the mathematical analysis of neural networks.

# Homework

1. Prove that the gradients for Sigmoidal basis functions are given as

$$\frac{\partial J_{LS}}{\partial \alpha_j} = 2 \sum_{i=1}^n o_{i,j} \left( \hat{f}(\mathbf{x}_i) - y_i \right) \quad o_{i,j} = \varphi(\mathbf{x}_i; \boldsymbol{\beta}_j, h_j)$$

$$\frac{\partial J_{LS}}{\partial \beta_j^{(k)}} = 2\alpha_j \sum_{i=1}^n o_{i,j} (1 - o_{i,j}) x_i^{(k)} \left( \hat{f}(\mathbf{x}_i) - y_i \right)$$

$$\frac{\partial J_{LS}}{\partial h_j} = 2\alpha_j \sum_{i=1}^n o_{i,j} (1 - o_{i,j}) \left( \hat{f}(\mathbf{x}_i) - y_i \right)$$



# Homework (cont.)

169

2. For your own toy 1-dimensional data, perform simulations using
  - 3-layer neural networks with sigmoid functions
  - Error back-propagationand analyze the results, e.g., by changing
  - Target functions
  - Number of training samples
  - Noise level
  - Number of neurons in the second layer ( $b$ )

# Notification of Final Assignment

1. Apply supervised learning techniques to your data set and analyze it.
  2. Write your opinion about this course
- Final report deadline: Aug 6<sup>th</sup> (Fri.)
  - E-mail submission is also accepted!  
*sugi@cs.titech.ac.jp*

# Mini-Workshop on Data Mining<sup>171</sup>

- On July 20<sup>th</sup> and 27<sup>th</sup>, we have a **mini-workshop on data mining**, instead of regular lectures.
- We have 16 speakers in total.
- Presentation: **10 minutes per person (including Q&A)**.
  - Specification of your dataset
  - Employed methods
  - Outcome

# Schedule

- July 13<sup>th</sup> : Preparation for workshop  
(no lecture)
- July 20<sup>th</sup> : Mini-workshop Day 1
- July 27<sup>th</sup> : Mini-workshop Day 2

# Mini-Workshop: Program

July 20<sup>th</sup>, 2010: 10:40-12:10

- Wittawat Jitkrittum
- Nakamasa Inoue
- Yukihiro Fukunaga
- Tatsuya Shigemura
- Ahmet Cetinkaya
- Michael Hamarneh
- Yonatan Andy Fajar Nugraha
- Youhei Namiki

July 27<sup>th</sup>, 2010: 10:40-12:10

- Johan Rohdin
- Hilman F. Pardede
- Sangeeta Biswas
- Rubaiya Rahtin Khan
- Mori Syogo
- M. Rasyid Aqmar
- Nozomi Kurihara
- Yuan Liang

Talk: 9 mins. Q&A: 2 mins.