

PNU Industrial Data Science

Neural Network

인공신경망

번호판 인식



인공지능, 사람보다 정확하게 자동차 번호판 읽어냈다

고성능 사진 해석기능 제공, 범죄 예방, 주차 관리 등 스마트 치안 및 생활 안전에 도움

책 읽어주는 조명

**전자
신문**   

  **강소·중견기업의 빅데이터 분석, SAS와 함께 하세요!**
지금 바로 무료 웨비나 시청하고, 교촌치킨도 받으세요!  [자세히보기](#)

책 읽어주는 조명 '램프' 일주일만에 완 판...네이버 AI 대중화 선봉에

발행일 : 2020.10.27 15:18

    [URL](#)

 가

· 한국 1위 로또 명당 당첨자 또 나왔다.. 2주 연속 1등 29억!



<글로벌 램프, 사전=네이버>

Contents

산업데이터과학은 산업현장에서 수집된 데이터를 분석하는데 필요한 기초 소양을 강의합니다.

01

S-R Agent

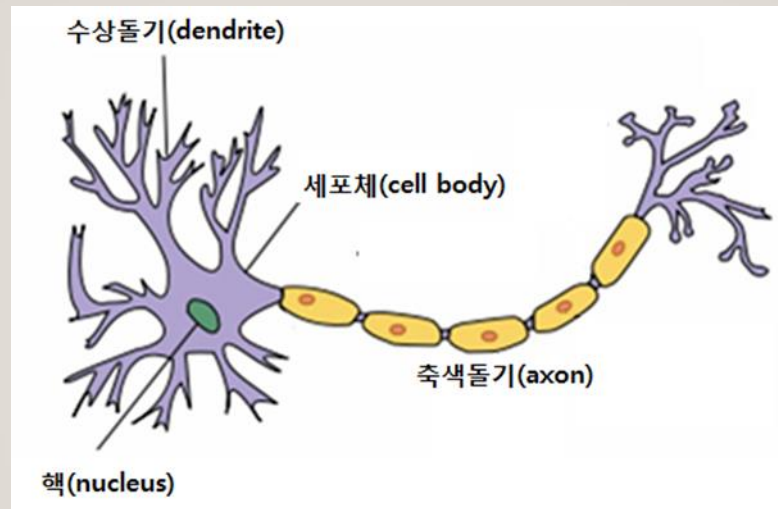
02

Perceptron

03

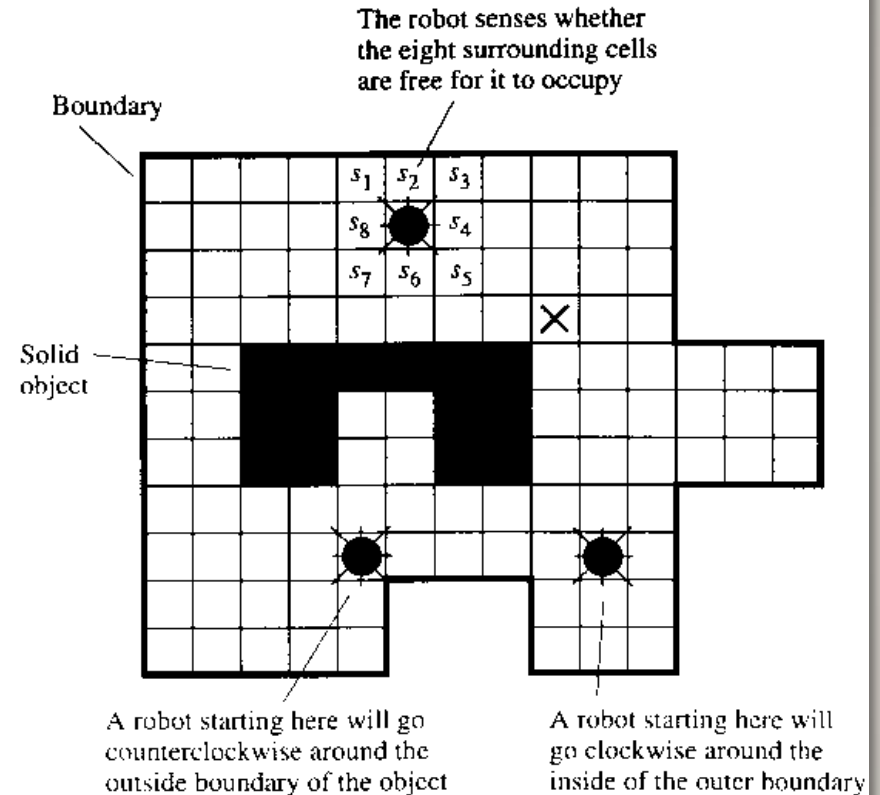
Backpropagation

Stimulus-Response Agents



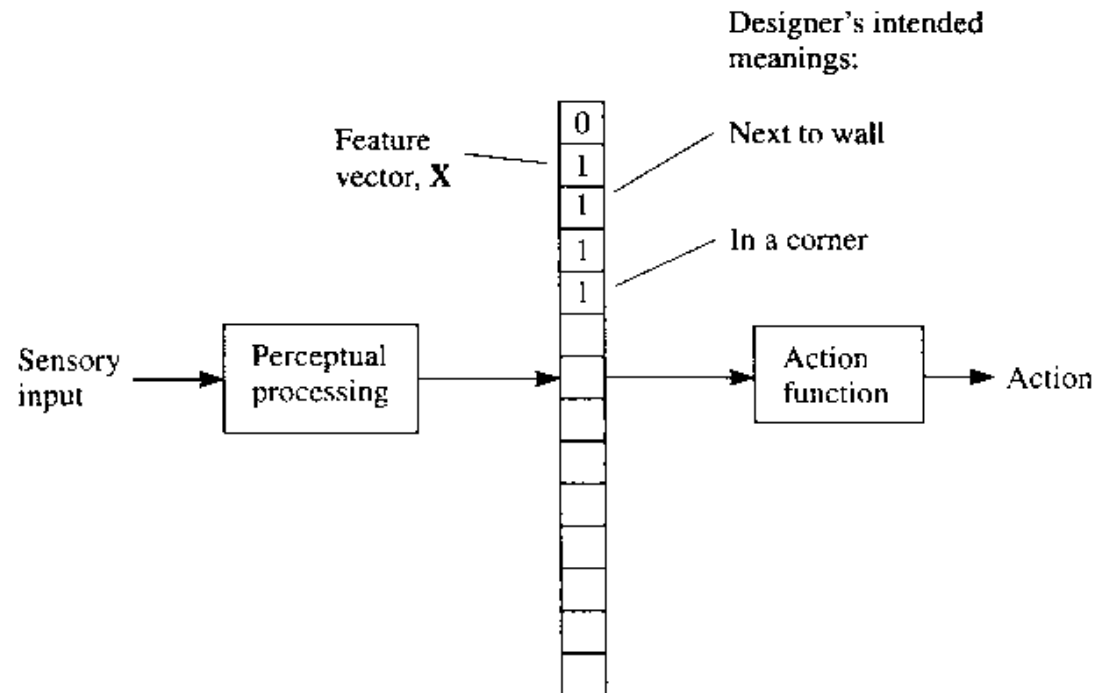
Perception and action

- S-R agents
 - Machines that have no internal state
 - Simply react to immediate stimuli
- 2-dimensional grid-space world
 - Enclosed by boundaries
 - Unmovable objects
 - Has no "tight spaces"
 - Action
 - Go to a cell adjacent to a boundary or object
 - follow that boundary along its perimeter



- Sensory inputs: $s_1, \dots s_8$
 - Have value 0 whenever the corresponding cell can be occupied by the root;
 - Otherwise, have value 1
- Robot movements
 - *north* moves the robot one cell up in the cellular grid
 - *east* moves the robot one cell to the right
 - *south* moves the robot one cell down
 - *west* moves the robot one cell to the left
- Designer's job
 - Specify a function of the sensory input
 - Selects actions appropriate for the task.
- Division of processes
 - Perception processing and action computation

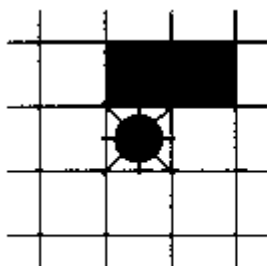
- Perceptual processing
 - produces feature vector \mathbf{X}
 - *numeric features*: real number
 - *categorical features*: categories
 - **Action computation**
 - Selects an action based on feature vector



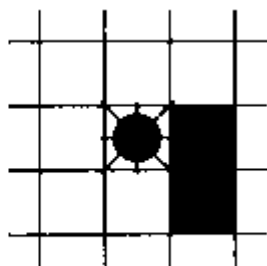
- The split between perception and action is arbitrary
 - The split is made in such a way that the same features would be used repeatedly in a variety of tasks to be performed
- The computation of features from sensory signals can be regarded as often used library routines
 - needed by many different action functions
- The next problems
 - (1) converting raw sensory data into a feature vector
 - (2) specifying an action function

Perception

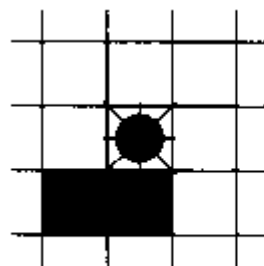
- The sensory input consists of the eight values of (s_1, \dots, s_8)
 - There are $2^8=256$ combinations of the values
- For the robot task, there are four binary-valued features of the sensory values that are useful for computing an appropriate action (x_1, \dots, x_4)
 - For example, $x_1 = 1$ if and only if $s_2=1$ or $s_3=1$
- Perceptual processing might occasionally give erroneous, ambiguous, or incomplete information about the robot's environment
 - Such errors might evoke inappropriate actions
- For robots with more complex sensors and tasks, designing appropriate perceptual processing can be challenging



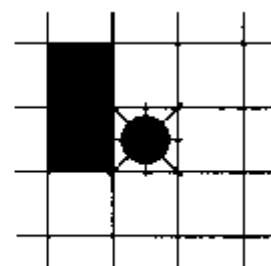
x_1



x_2



x_3



x_4

Action

- Specifying a function that selects the appropriate boundary-following action
 - If $x_1=1$ and $x_2=0$, move **east**
 - If $x_2=1$ and $x_3=0$, move **south**
 - If $x_3=1$ and $x_4=0$, move **west**
 - If $x_4=1$ and $x_1=0$, move **north**
- None of the features has value 1, the robot can move in *any* direction until it encounters a boundary
- The robot take the actions happen to be Boolean combination of the features
 - The features themselves are also Boolean combination of the sensory inputs

Representing and implementing action function: production system

- Production system comprises an ordered list of rules called *production* rules or *productions*
 - $c_i \rightarrow a_i$, where c_i is the condition part and a_i is the action part
 - Production system consists of a list of such rules
 - Condition part
 - Can be any binary-valued function of the features
 - Often a monomial
 - Action part
 - Primitive action, a call to another productive system, or a set of actions to be executed simultaneously
- Production system representation for the boundary following routine
 - An example of a durative systems-system that never ends

$$\overline{x_4 x_1} \rightarrow \textit{north}$$

$$\overline{x_3 x_4} \rightarrow \textit{west}$$

$$\overline{x_2 x_3} \rightarrow \textit{south}$$

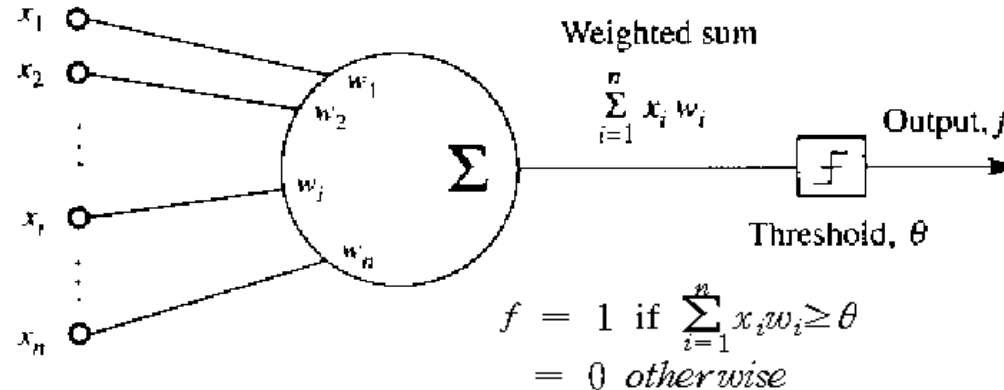
$$\overline{x_1 x_2} \rightarrow \textit{east}$$

$$1 \rightarrow \textit{north}$$

Network

- *Threshold logic unit (TLU)*

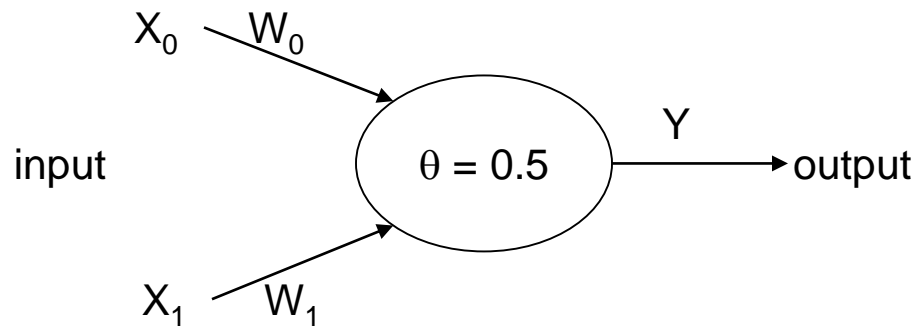
- Circuit consists of networks of threshold elements or other elements that compute a nonlinear function of a weighted sum of their inputs
- **TLU separates the space of input vectors yielding an above-threshold response from those yielding a below-threshold response by a linear space-called a *hyperplane***



- *Linearly separable functions*

- The boolean functions implementable by a TLU
- Many boolean functions are linearly separable
- Exclusive-or function of two variables is an example of not linearly separable

– AND, XOR Example



input		output		
X_0	X_1	AND	f	XOR
0	0	0	0	0
0	1	0	1	1
1	0	0	0	1
1	1	1	0	0

• AND

$$0 \times W_0 + 0 \times W_1 = 0 < 0.5$$

$$0 \times W_0 + 1 \times W_1 = W_1 < 0.5$$

$$1 \times W_0 + 0 \times W_1 = W_0 < 0.5$$

$$1 \times W_0 + 1 \times W_1 = W_0 + W_1 > 0.5$$

→ $W_0, W_1 : 0.3 \text{ or } 0.4$

- XOR

$$0 \times W_0 + 0 \times W_1 = 0 < 0.5$$

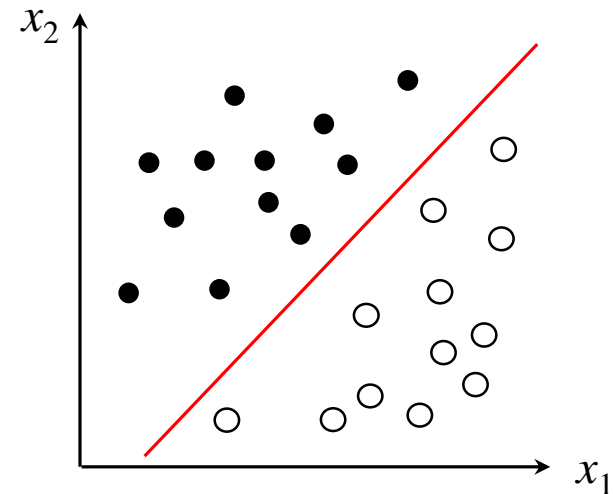
$$0 \times W_0 + 1 \times W_1 = W_1 > 0.5$$

$$1 \times W_0 + 0 \times W_1 = W_0 > 0.5$$

$$1 \times W_0 + 1 \times W_1 = W_0 + W_1 < 0.5$$

→ W_0, W_1 do not exist that satisfy above

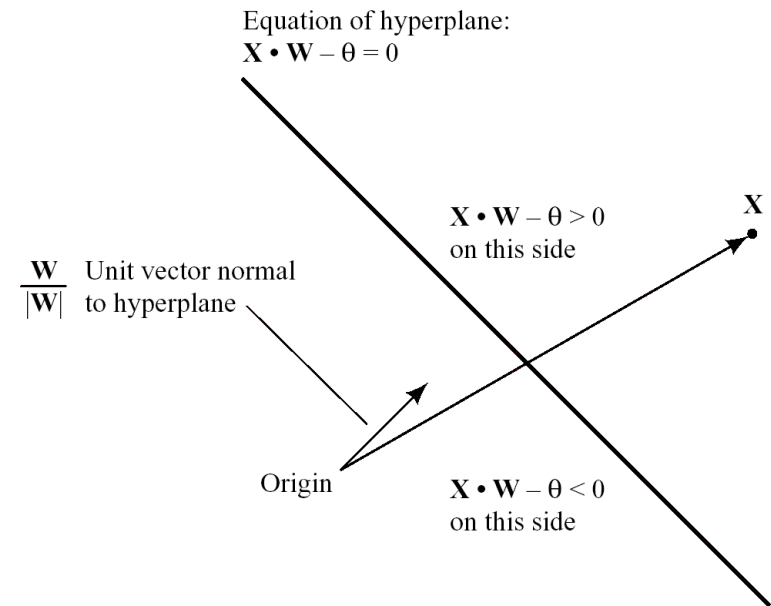
→ cannot solve XOR



Neural Network

- TLU (threshold logic unit): Basic units for neural networks
 - Based on some properties of biological neurons
- Training set
 - Input: real value, boolean value, ...
N-dimensional vector: $\mathbf{X}=(x_1, x_2, \dots x_n)$
 - Output:
 - d_i : associated actions (Label, Class ...)
- Target of training
 - Finding $f(\mathbf{X})$ corresponds "acceptably" to the members of the training set.
 - Supervised learning: Labels are given along with the input vectors.

- Training TLU: Adjusting variable weights
- A single TLU: Perceptron, Adaline (*adaptive linear element*) [Rosenblatt 1962, Widrow 1962]
- Elements of TLU
 - Weight: $\mathbf{W} = (w_1, \dots, w_n)$
 - Threshold: θ
- Output of TLU: Using weighted sum $s = \mathbf{W} \cdot \mathbf{X}$
 - 1 if $s - \theta > 0$
 - 0 if $s - \theta < 0$
- Hyperplane
 - $\mathbf{W} \cdot \mathbf{X} - \theta = 0$



Augmented vectors

- Adopting the convention that threshold is fixed to 0.
- Arbitrary thresholds: $(n + 1)$ -dimensional vector
 - $\mathbf{X} = (x_1, \dots, x_n, 1)$
 - $\mathbf{W} = (w_1, \dots, w_n, -\theta)$
- Output of TLU
 - 1 if $\mathbf{W} \cdot \mathbf{X} \geq 0$
 - 0 if $\mathbf{W} \cdot \mathbf{X} < 0$

Gradient Descent Methods

- Training TLU: minimizing the *error function* by adjusting weight values.
- Batch learning v.s. incremental learning
- Commonly used error function: squared error $\varepsilon = (d - f)^2$

– Gradient:
$$\frac{\partial \varepsilon}{\partial \mathbf{W}} \stackrel{\text{def}}{=} \left[\frac{\partial \varepsilon}{\partial w_1}, \dots, \frac{\partial \varepsilon}{\partial w_i}, \dots, \frac{\partial \varepsilon}{\partial w_{n+1}} \right]$$

– Chain rule:
$$\frac{\partial \varepsilon}{\partial \mathbf{W}} = \frac{\partial \varepsilon}{\partial s} \frac{\partial s}{\partial \mathbf{W}} = \frac{\partial \varepsilon}{\partial s} \mathbf{X} = -2(d - f) \frac{\partial f}{\partial s} \mathbf{X}$$

- f is not continuously differentiable with respect to s ($\partial f / \partial s$) :
 - Ignoring threshold function: $f = s$
 - Replacing threshold function with differentiable nonlinear function

The Widrow-Hoff Procedure

- Weight update procedure:
 - Using $f = s = \mathbf{W} \cdot \mathbf{X}$
 - Data labeled 1 \rightarrow 1, Data labeled 0 \rightarrow -1
- Gradient: if $f \neq s$,

$$\frac{\partial \varepsilon}{\partial \mathbf{W}} = -2(d - f) \frac{\partial f}{\partial s} \mathbf{X} = -2(d - f) \mathbf{X}$$

- New weight vector
- Widrow-Hoff (delta) rule $\mathbf{W} \leftarrow \mathbf{W} + c(d - f) \mathbf{X}$
 - $(d - f) > 0 \rightarrow$ increasing $s \rightarrow$ decreasing $(d - f)$
 - $(d - f) < 0 \rightarrow$ decreasing $s \rightarrow$ increasing $(d - f)$

The generalized delta procedure

- Sigmoid function (differentiable): [Rumelhart, et al. 1986]

$$f(s) = 1/(1 + e^{-s})$$

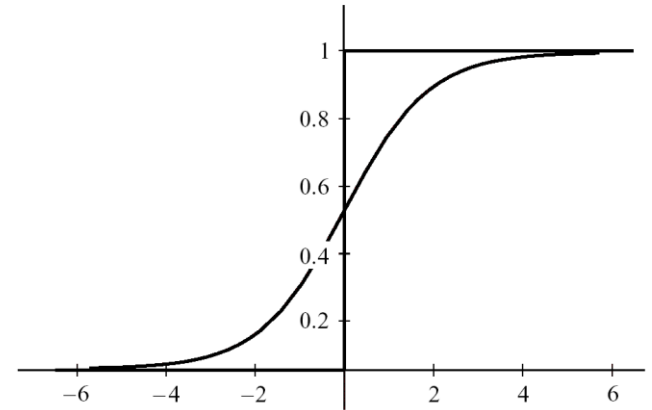
- Gradient:

$$\frac{\partial \varepsilon}{\partial \mathbf{W}} = -2(d - f) \frac{\partial f}{\partial s} \mathbf{X} = -2(d - f) f(1 - f) \mathbf{X}$$

- Generalized delta procedure:

$$\mathbf{W} \leftarrow \mathbf{W} + c(d - f)f(1 - f)\mathbf{X}$$

- Target output: 1, 0
- Output f = output of sigmoid function
- $f(1 - f) = 0$, where $f = 0$ or 1
- Weight change can occur only within 'fuzzy' region surrounding the hyperplane (near the point $f(s) = 1/2$).



The error-correction procedure

- Using threshold unit: $(d - f)$ can be either 1 or -1 .

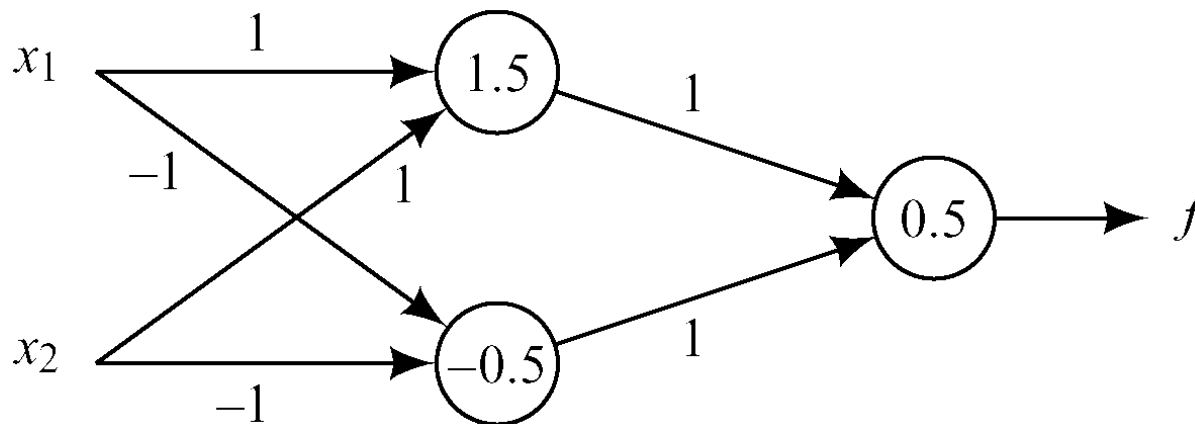
$$\mathbf{W} \leftarrow \mathbf{W} + c(d - f)\mathbf{X}$$

- In the linearly separable case, after finite iteration, \mathbf{W} will be converged to the solution.
 - If there is some weight vector, that produces a correct output for all of the input vectors, then after a finite number of input vector presentation, the procedure will find such a weight vector and thus make no more weight changes
- In the nonlinearly separable case, \mathbf{W} will never be converged.
- The Widrow-Hoff and generalized delta procedures will find minimum squared error solutions even when the minimum error is not zero.

Neural network: motivation

- Need for use of multiple TLUs
 - Feedforward network: no cycle
 - Recurrent network: cycle (treated in a later chapter)
 - Layered feedforward network
 - j_{th} layer can receive input only from $j - 1_{\text{th}}$ layer.
- Example :

$$f = x_1x_2 + \bar{x}_1\bar{x}_2$$



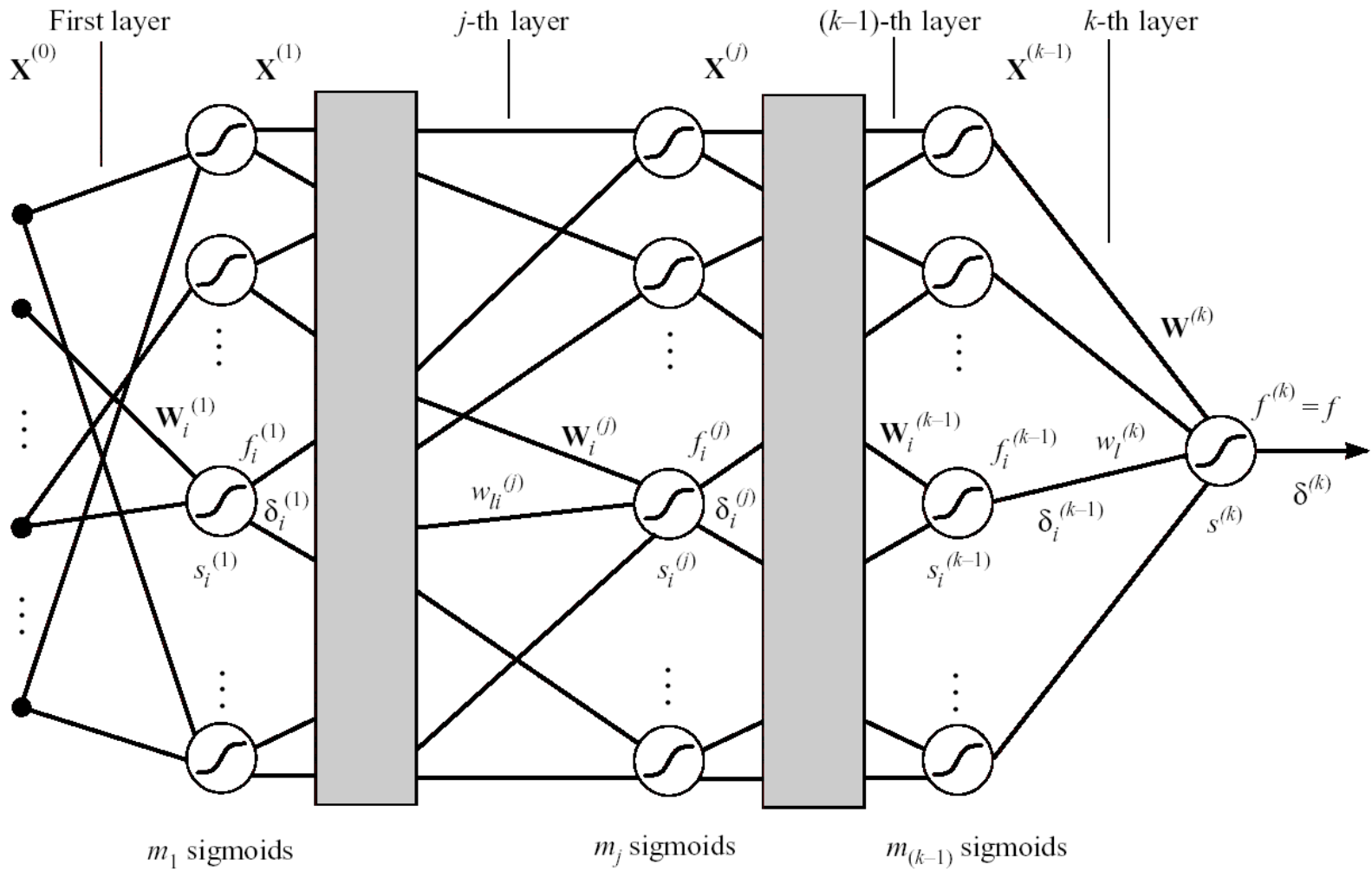
Notation

- Hidden unit: neurons in all but the last layer
- Output of j -th layer: $\mathbf{X}^{(j)} \rightarrow$ input of $(j+1)$ -th layer
- Input vector: $\mathbf{X}^{(0)}$
- Final output: f
- The weight of i -th sigmoid unit in the j -th layer: $\mathbf{W}_i^{(j)}$
- Weighted sum of i -th sigmoid unit in the j -th layer: $s_i^{(j)}$

$$s_i^{(j)} = \mathbf{X}^{(j-1)} \cdot \mathbf{W}_i^{(j)}$$

- Number of sigmoid units in j -th layer: m_j

$$\mathbf{W}_l^{(j)} = (w_{1,i}^{(j)}, \dots, w_{l,i}^{(j)}, \dots, w_{m_{(j-1)}+1,i}^{(j)})$$



The backpropagation method

- Gradient of $\mathbf{W}_i^{(j)}$:

$$s_i^{(j)} = \mathbf{X}^{(j-1)} \cdot \mathbf{W}_i^{(j)}$$

$$\begin{aligned} \frac{\partial \varepsilon}{\partial \mathbf{W}_i^{(j)}} &= \frac{\partial \varepsilon}{\partial s_i^{(j)}} \frac{\partial s_i^{(j)}}{\partial \mathbf{W}_i^{(j)}} = \frac{\partial \varepsilon}{\partial s_i^{(j)}} \mathbf{X}^{(j-1)} \\ &= -2(d - f) \frac{\partial f}{\partial s_i^{(j)}} \mathbf{X}^{(j-1)} = -2\delta_i^{(j)} \mathbf{X}^{(j-1)} \end{aligned}$$

Local gradient

$$\frac{\partial \varepsilon}{\partial s_i^{(j)}} = \frac{\partial (d - f)^2}{\partial s_i^{(j)}} = -2(d - f) \frac{\partial f}{\partial s_i^{(j)}}$$

- Weight update:

$$\mathbf{W}_i^{(j)} \leftarrow \mathbf{W}_i^{(j)} + c_i^{(j)} \delta_i^{(j)} \mathbf{X}^{(j-1)}$$

Computing weight changes in the final layer

- Local gradient:

$$\begin{aligned}\delta^{(k)} &= (d - f) \frac{\partial f}{\partial s_i^{(k)}} \\ &= (d - f) f (1 - f)\end{aligned}$$

- Weight update:

$$\mathbf{W}^{(k)} \leftarrow \mathbf{W}^{(k)} + c^{(k)} (d - f) f (1 - f) \mathbf{X}^{(k-1)}$$

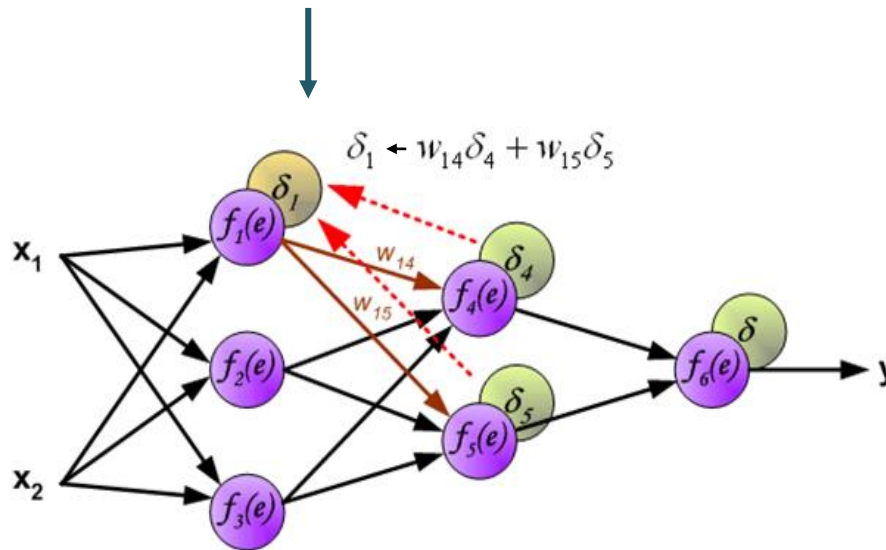
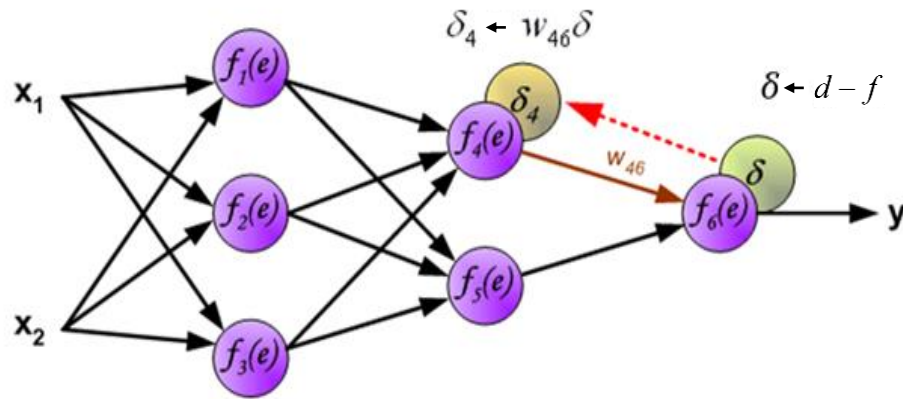
- Attention to recursive equation of local gradient!

$$\delta^{(k)} = f(1-f)(d-f)$$

$$\delta_i^{(j)} = f_i^{(j)}(1-f_i^{(j)}) \sum_{l=1}^{m_{j+1}} \delta_i^{(j+1)} w_{il}^{(j+1)}$$

- Backpropagation:
 - Error is back-propagated from the output layer to the input layer
 - Local gradient of the latter layer is used in calculating local gradient of the former layer.

$$\delta_i^{(j)} = f_i^{(j)} (1 - f_i^{(j)}) \sum_{l=1}^{m_{j+1}} \delta_l^{(j+1)} w_{il}^{(j+1)}$$



Generalization, Accuracy, and overfitting

- Generalization ability:
 - NN appropriately classifies vectors not in the training set.
 - Measurement = accuracy
- Curve fitting
 - Number of training input vectors \geq number of degrees of freedom of the network.
 - In the case of m data points, is $(m-1)$ -degree polynomial best model? **No, it can not capture any special information.**
- Overfitting
 - Extra degrees of freedom are essentially just fitting the noise.
 - Given sufficient data, the *Occam's Razor* principle dictates to choose the lowest-degree polynomial that adequately fits the data.

Validation of NN

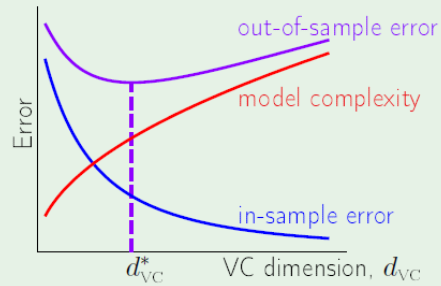
- Out-of-sample-set error rate
 - Error rate on data drawn from the same underlying distribution of training set.
- Dividing available data into a *training* set and a *validation* set
 - Usually use 2/3 for training and 1/3 for validation
- k -fold cross validation
 - k disjoint subsets (called folds).
 - Repeat training k times with the configuration: one validation set, $k-1$ (combined) training sets.
 - Take average of the error rate of each validation as the out-of-sample error.
 - Empirically 10-fold is preferred.

Error in NN

Characterizing the feasibility of learning for infinite M

Characterizing the tradeoff:

Model complexity	↑	E_{in}	↓
Model complexity	↑	$E_{out} - E_{in}$	↑



Percent error

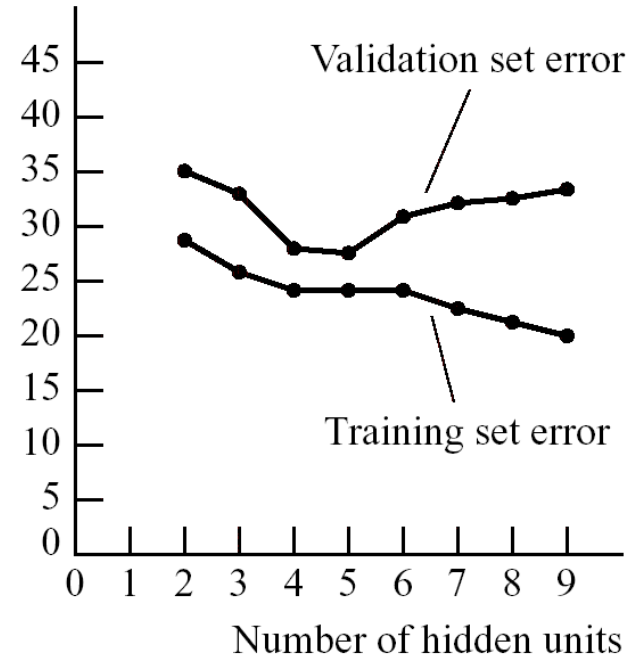


Fig 3.9 Estimate of Generalization Error Versus Number of Hidden Units

Appendix: Computing changes to the weights in intermediate layers

- Local gradient:

$$\begin{aligned}\delta_i^{(j)} &= (d - f) \frac{\partial f}{\partial s_i^{(j)}} \\ &= (d - f) \left[\frac{\partial f}{\partial s_1^{(j+1)}} \frac{\partial s_1^{(j+1)}}{\partial s_i^{(j)}} + \dots + \frac{\partial f}{\partial s_l^{(j+1)}} \frac{\partial s_l^{(j+1)}}{\partial s_i^{(j)}} + \dots + \frac{\partial f}{\partial s_{m_{j+1}}^{(j+1)}} \frac{\partial s_{m_{j+1}}^{(j+1)}}{\partial s_i^{(j)}} \right] \\ &= \sum_{l=1}^{m_{j+1}} (d - f) \frac{\partial f}{\partial s_l^{(j+1)}} \frac{\partial s_l^{(j+1)}}{\partial s_i^{(j)}} = \sum_{l=1}^{m_{j+1}} \delta_l^{(j+1)} \frac{\partial s_l^{(j+1)}}{\partial s_i^{(j)}}\end{aligned}$$

- The final output f depends on $s_l^{(j)}$ through of the summed inputs to the sigmoids in the $(j+1)$ -th layer.

- Need for computation of $\frac{\partial s_l^{(j+1)}}{\partial s_i^{(j)}}$

$$s_l^{(j+1)} = \mathbf{X}^{(j)} \cdot \mathbf{W}_l^{(j+1)} = \sum_{v=1}^{m_j+1} f_v^{(j)} w_{vl}^{(j+1)}$$

$$\begin{aligned} \frac{\partial s_l^{(j+1)}}{\partial s_i^{(j)}} &= \frac{\partial \left[\sum_{v=1}^{m_j+1} f_v^{(j)} w_{vl}^{(j+1)} \right]}{\partial s_i^{(j)}} = \sum_{v=1}^{m_j+1} w_{vl}^{(j+1)} \frac{\partial f_v^{(j)}}{\partial s_i^{(j)}} \\ &= w_{il}^{(j+1)} f_i^{(j)} (1 - f_i^{(j)}) \end{aligned}$$

– $v \neq i$:

$$\frac{\partial f_v^{(j)}}{\partial s_i^{(j)}} = 0$$

$v = i$:

$$\frac{\partial f_v^{(j)}}{\partial s_i^{(j)}} = f_v^{(j)} (1 - f_v^{(j)})$$

– Consequently,

$$\delta_i^{(j)} = f_i^{(j)} (1 - f_i^{(j)}) \sum_{l=1}^{m_{j+1}} \delta_l^{(j+1)} w_{il}^{(j+1)}$$