

다층퍼셉트론 신경회로망

Multilayer Perceptron Neural Network

# 선형 회귀분석 vs 비선형 회귀분석

- 선형 회귀분석

$$\hat{y} = \Theta + \sum_{i=1}^p w_i x_i$$

- “선형”: 모델 파라미터인 beta 또는 w 들간의 관계
- x 대신  $x^2$ ,  $x^3$ ,  $e^x$ ,  $\log x$ ,  $\sin x$  가 대체 되어도 모두 선형 회귀분석

# 선형 회귀분석 vs 비선형 회귀분석

- 비선형 회귀분석
  - 모델 파라미터 간의 관계가 비 선형
  - 다양한 모델 존재
  - 그 가운데 가장 인기 좋은 신경회로망
  - 신경 회로망 가운데 가장 많이 사용되는 multi-layer perceptron 다층 퍼셉트론

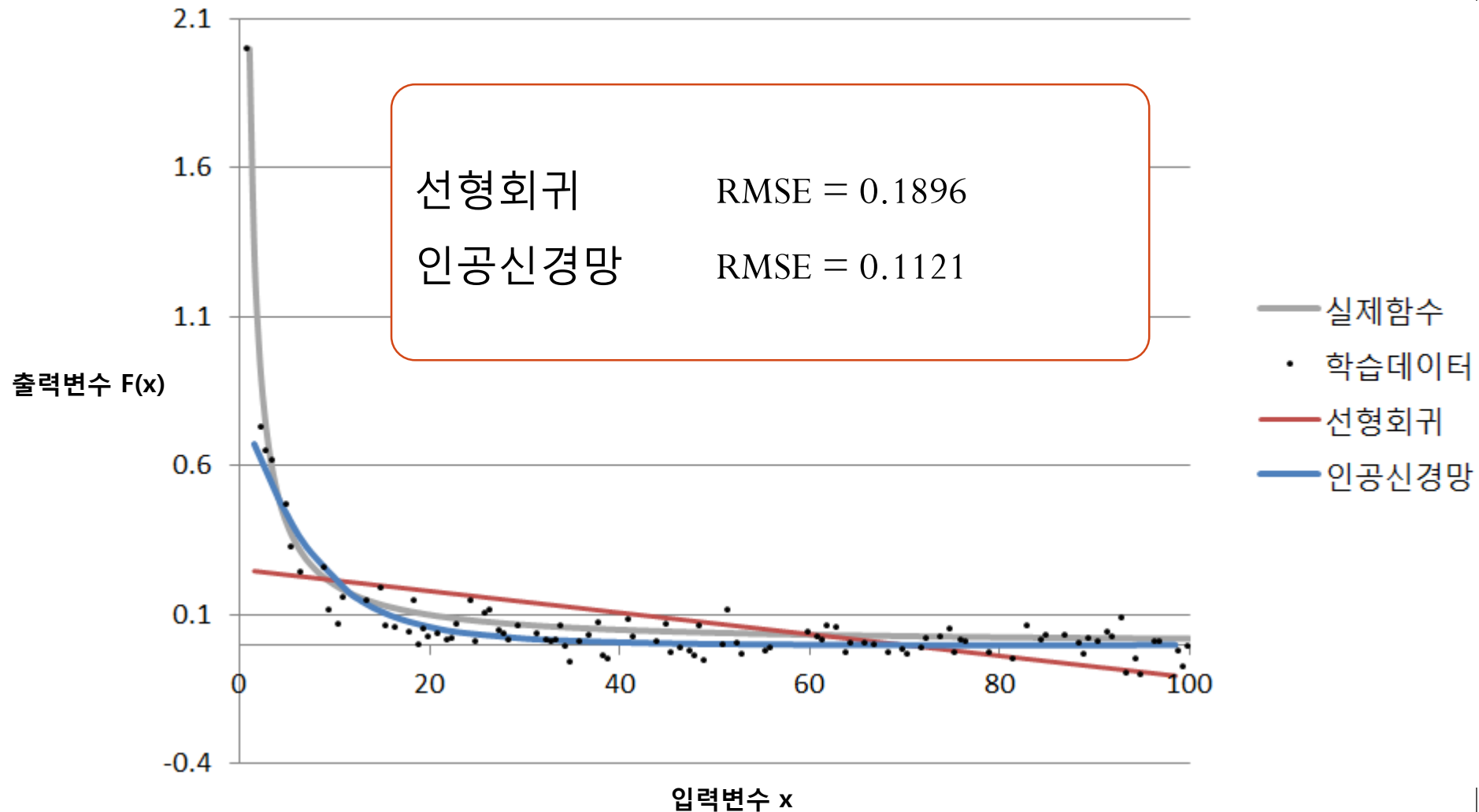
# 선형 회귀분석 vs 비선형 회귀분석

- 비선형 회귀분석
  - 다양한 모델 존재
  - 그 가운데 가장 인기 좋은 신경회로망
  - 신경 회로망 가운데 가장 많이 사용되는 multi-layer perceptron 다층 퍼셉트론
- 선형과 비선형의 비교
  - 선형은 직선 fit ( $x^2$  이나  $e^x$  없는 경우)
  - 비선형은 곡선 fit

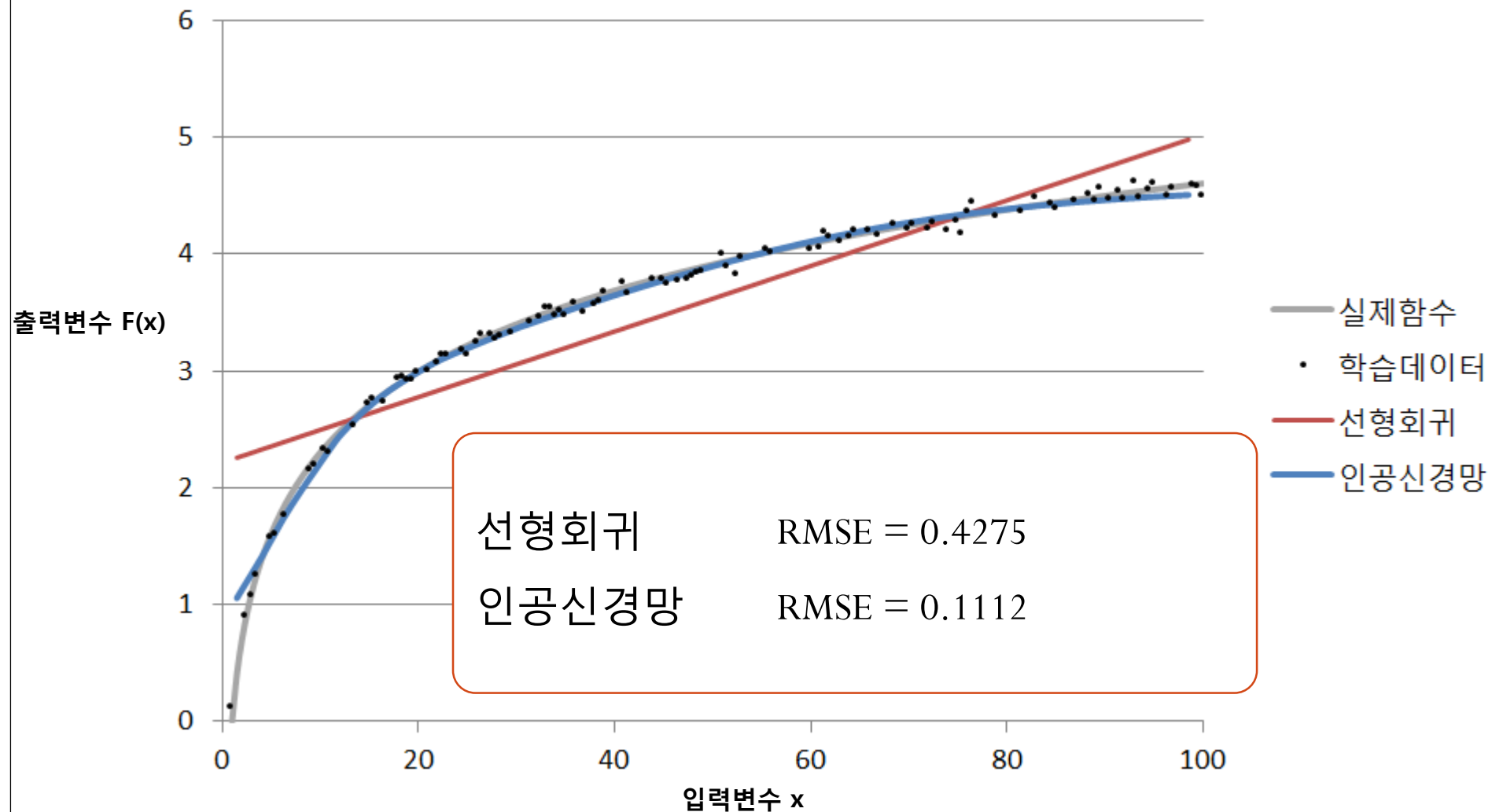
# 함수 fit

- 주어진 함수  $f$ 로부터 데이터  $D$  생성
- 생성된  $D$ 를 바탕으로  $y = f'(X)$  구축
- 모델  $f'$ 와  $f$  비교 ( $f'$ 이  $f$ 와 비슷한가?)
- 주어진 함수  $f$ 
  - $y = 2/x$
  - $y = \log_2 x$
  - $y = \exp(-0.2 * x)$
  - $y = \sin(x)$

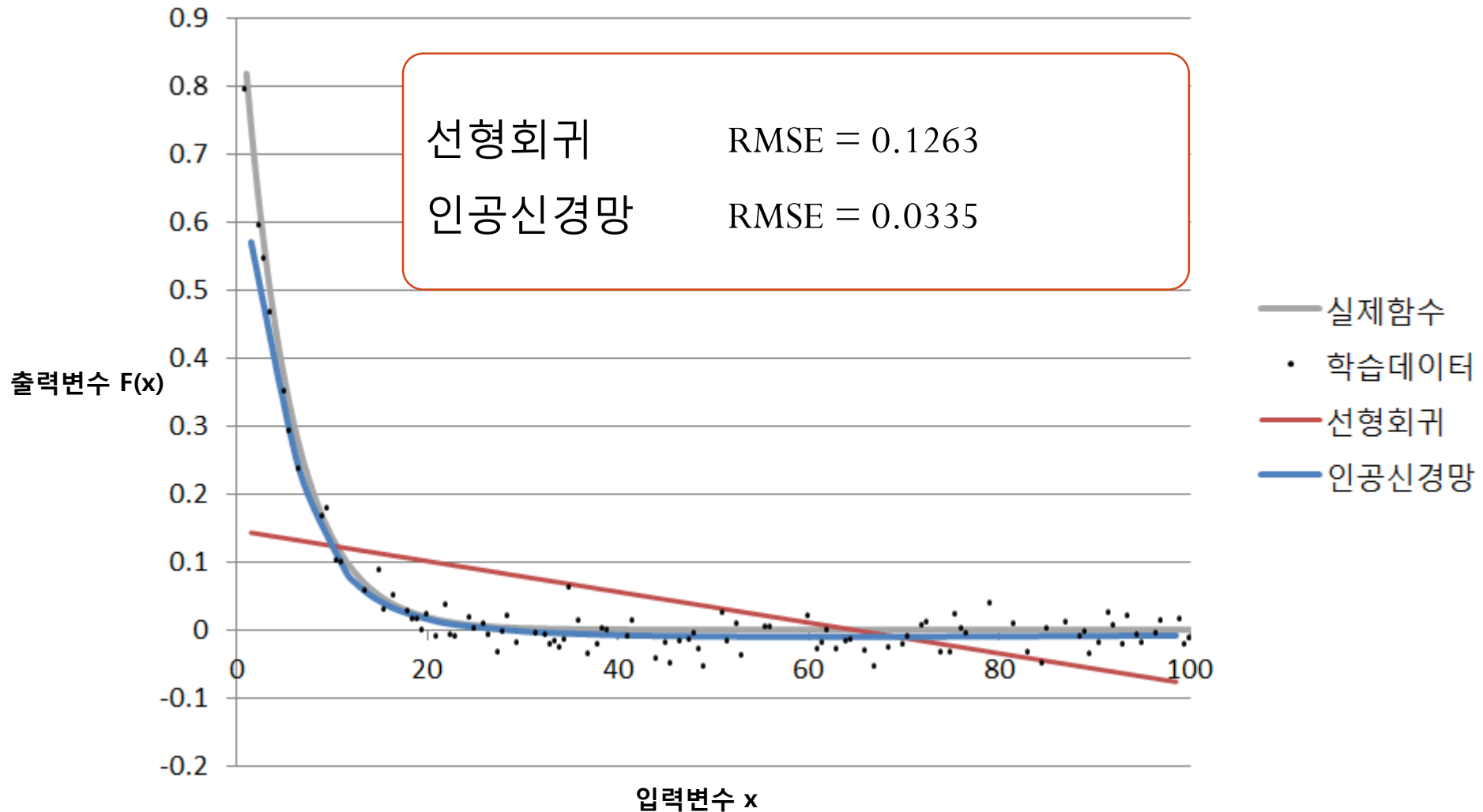
$$F(x) = 2/x \quad (1 \leq x \leq 100)$$



$$F(x) = \log(x) \quad (1 \leq x \leq 100)$$

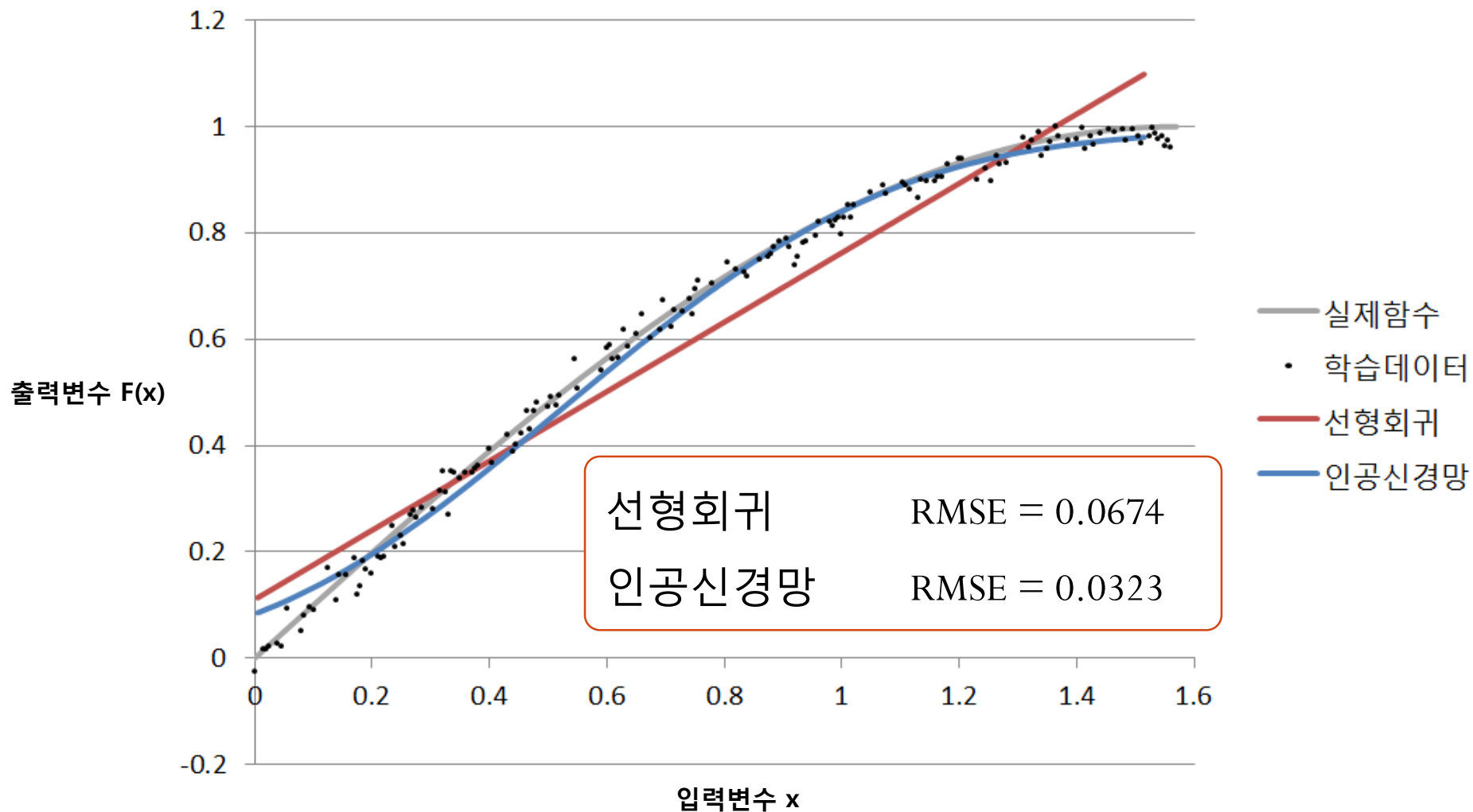


$$F(x) = \exp(-0.2 \cdot x) \quad (1 \leq x \leq 100)$$





$$F(x) = \sin(x) \quad (0 \leq x \leq \pi/2)$$



# 선형으로도 fit 가능하지만...

- 선형 회귀분석으로도  $2/x$ ,  $\log x$ ,  $e(-0.2x)$ ,  $\sin x$  항을 넣으면 위 함수들을 정확히 fit 할 수 있음
- 그러나 현실적으로 데이터 세트  $D$  만 주어졌을 때에, 어떤 “비선형 항”을 넣어야 하는지 판단 불가
- 따라서 신경망과 같은 general nonlinear model 이 사용성 측면에서 뛰어남

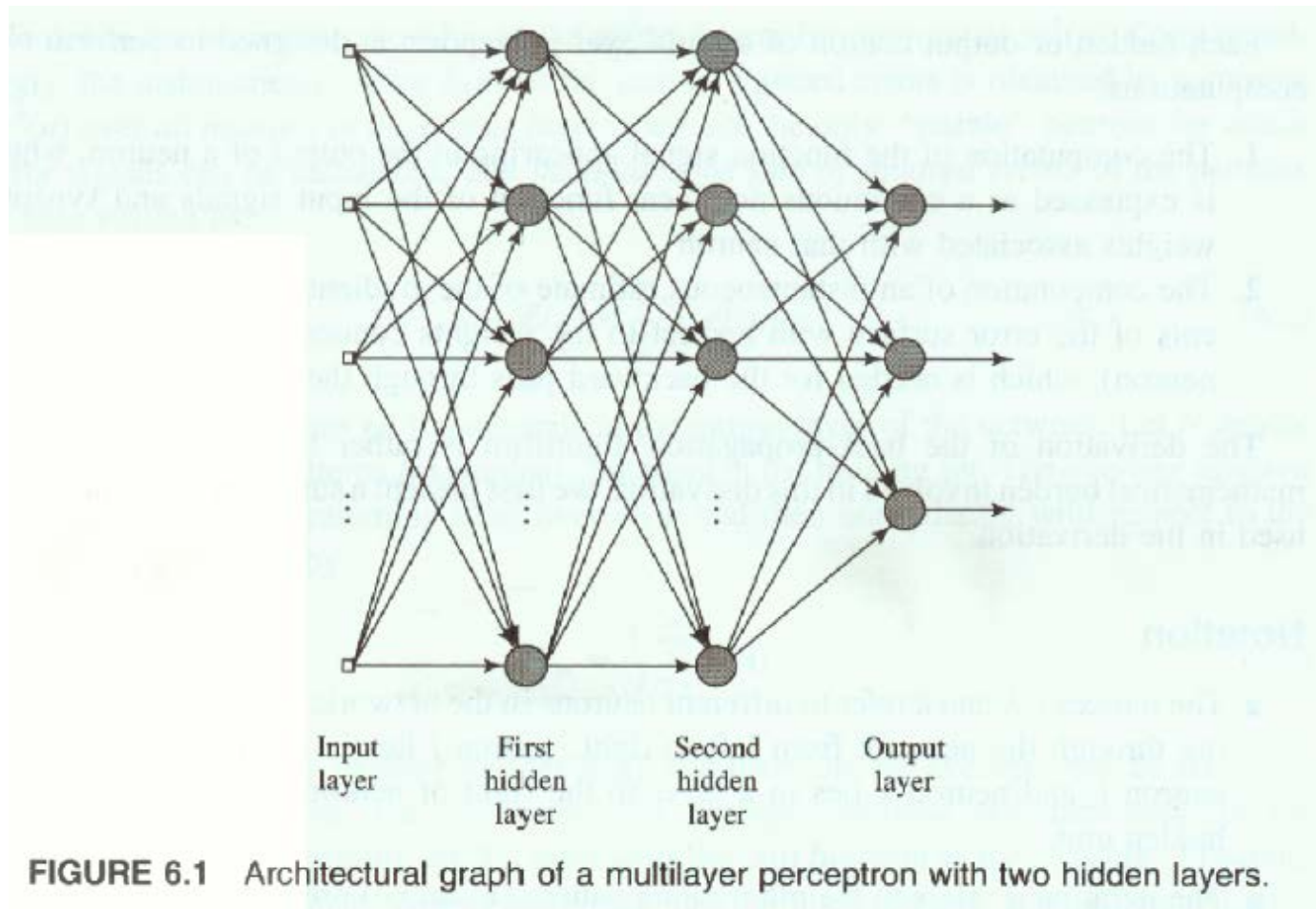
# Neural networks

- 신경회로망



- 인간: ~1천억 개 뉴론 들이 10조 개의 시냅스를 통해 연결됨

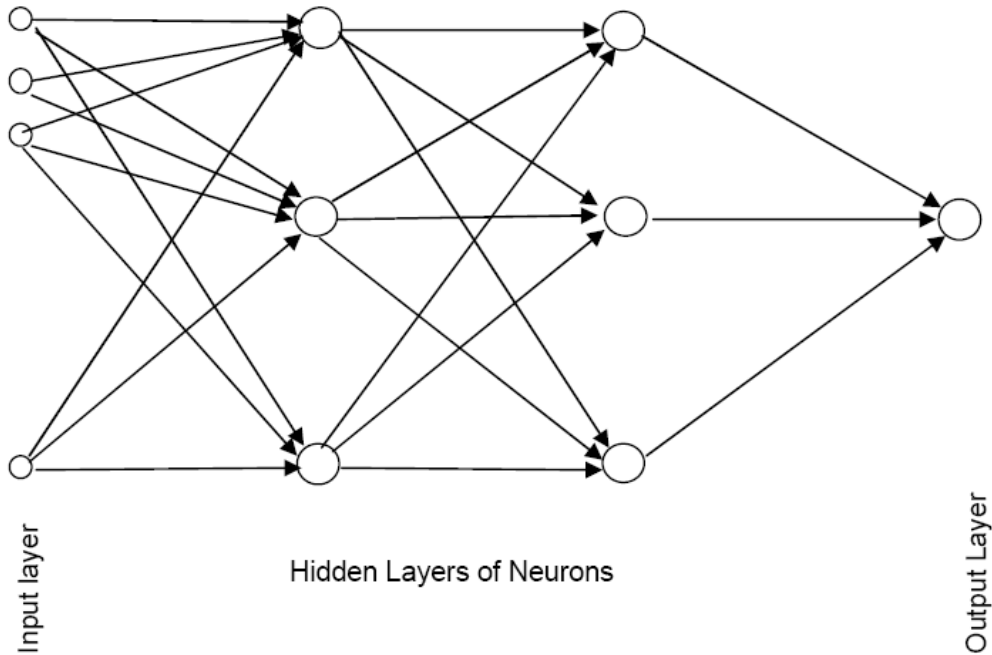
# 다층 퍼셉트론 multi-layer perceptron



# 망 구조

- 노드, 뉴런 (회귀식 변수)
- 노드 층
  - 입력층 input layer
  - 은닉층 hidden layer
  - 출력층 output layer
- 에지, 시냅스 (회귀식 계수)

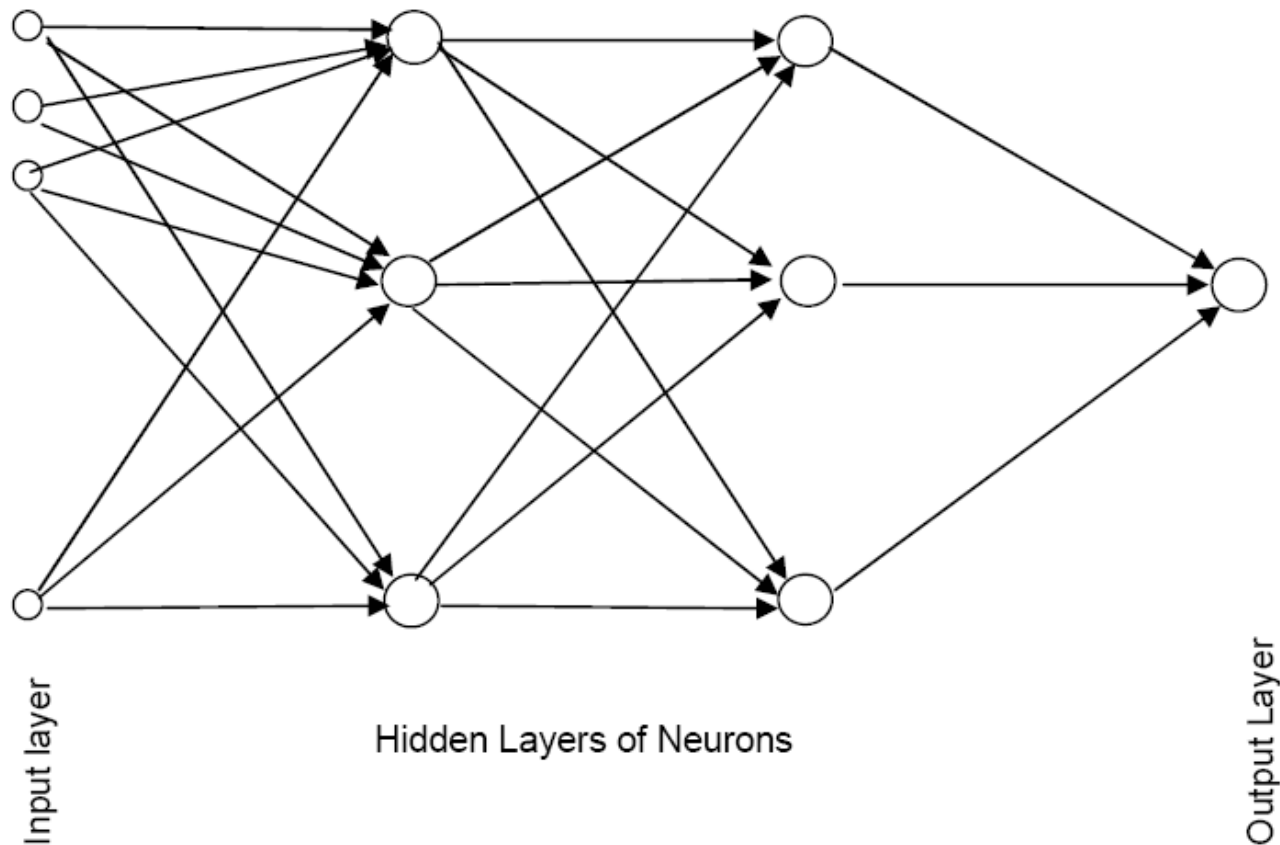
# 3층 퍼셉트론 구조



\* 영국식 층 계산 법

# 3층 퍼셉트론 구조

- input layer: input nodes = input or **independent variables**  $x$
- output layer: output node = output or **dependent variable**  $y$
- hidden layer: hidden nodes = ?  $h$



# 각 노드에서 하는 계산

- Action potential, nonlinearity, threshold, synapse, other neuron's

$$output_j = g(\theta_j + \sum_{i=1}^p w_{ij} x_i)$$



# 1층 퍼셉트론 구조는?

- P 개의 입력 노드와 1개의 출력노드를 가진...
- $P=3$
- 시냅스 수는?

# 선형 회귀 모델 Linear Regression!!

출력 노드가 하나이고 중간층이 없는 망은, 여기서  $g$ 는 항등함수, 선형 회귀분석 모델과 같은 형태를 취한다.

$$\hat{y} = \Theta + \sum_{i=1}^p w_i x_i$$

## 2층 퍼셉트론 구조는?

- P 개의 입력 노드, H 개의 은닉 노드, 1 개의 출력 노드
- $P=3, H=4$
- 시냅스 수는?

# 비선형 2층 퍼셉트론 모델

- Nonlinear regression 인 경우,
  - 히든 노드의 g 함수는 sigmoid 이고,
  - 출력 노드의 g 함수는 identity (or linear) 를 사용

- 수식으로 표현하면

$$y = \Theta_0 + \sum_{j=1}^H w_j \{g(\Theta_j + \sum_{i=1}^p w_{ij} x_i)\}$$

- Logistic Regression 몇 개?

# Example – Using fat & salt content to predict consumer acceptance of cheese

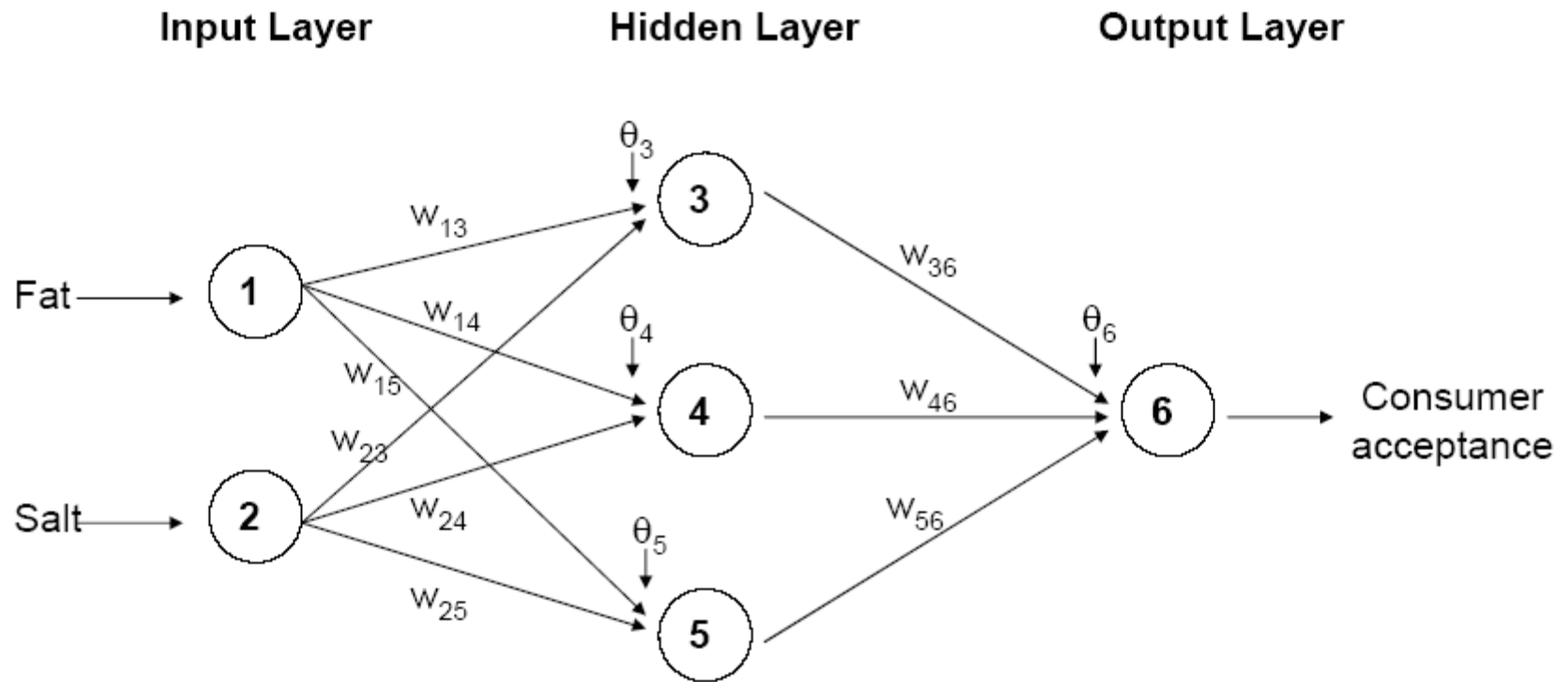


Figure 11.2: Neural network for the tiny example. Circles represent nodes,  $w_{i,j}$  on arrows are weights, and  $\theta_j$  are node bias values.

# Example - Data

<i>Obs.</i>	<i>Fat Score</i>	<i>Salt Score</i>	<i>Acceptance</i>
1	0.2	0.9	1
2	0.1	0.1	0
3	0.2	0.4	0
4	0.2	0.5	0
5	0.4	0.5	1
6	0.3	0.8	1

모델 작동

# 입력층

입력층에서, 입력 = 출력

- E.g., record #1에서:

지방 입력 = 출력 = 0.2

염분 입력 = 출력 = 0.9

입력층의 출력 = 은닉층으로 입력



# 은닉층

이 예에서, 은닉층은 3개의 노드를 가짐

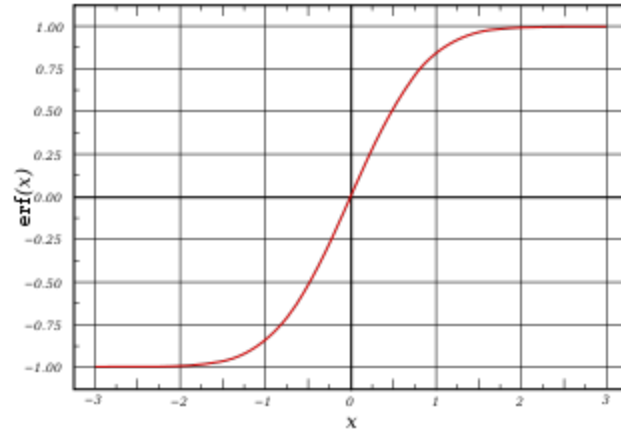
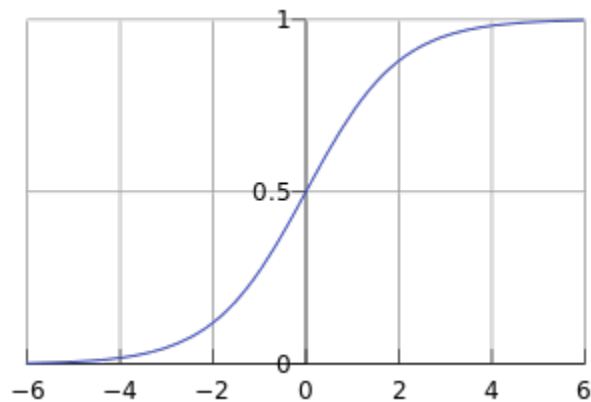
각 노드는 전체 입력 노드의 출력을 입력함

각 은닉층의 출력은 입력 가중치 합에 함수

$$output_j = g(\theta_j + \sum_{i=1}^p w_{ij} x_i)$$

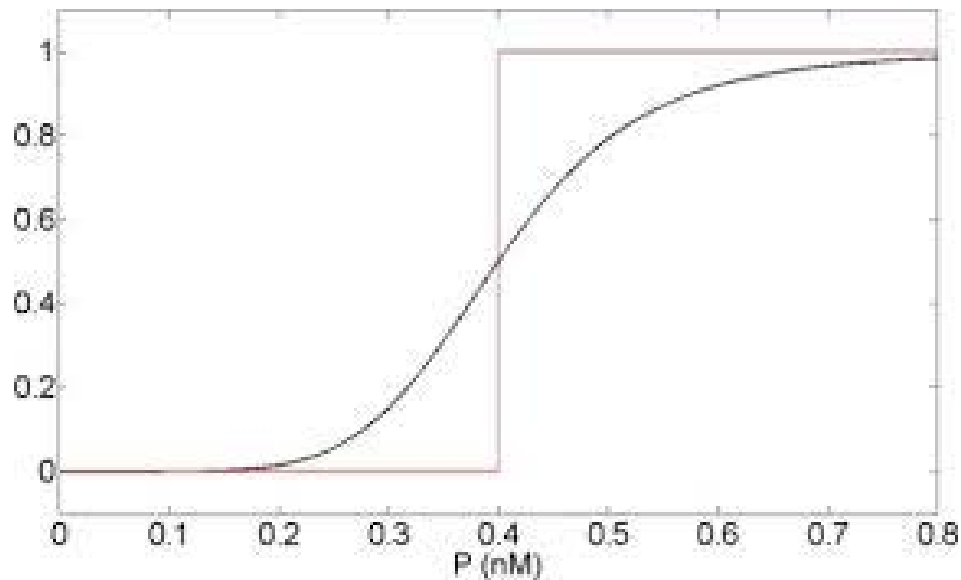
# Function g?

- $g(x) = 1/(1+\exp(-x))$
- 시그모이드, 로지스틱
- 뉴론의 활성화 함수 또는 학습 함수



# Function g?

- $g(x) = 1/(1+\exp(-k*x))$
- $k$  값이 아주 크면, 시그모이드, 로지스틱 함수는 어떤 모양이 되는가?



# 선형 분리가능성 linear separability

- OR, AND 문제
  - “Decision Boundary”
  - 1층 perceptron 으로 분리 가능
- XOR 문제
  - 1층 perceptron 으로 분리 불가능
  - 1969 “Perceptron” by Minsky
  - 여러 개의 1층 perceptron 으로는 분리 가능!
  - 어떻게? Stacking!

# Synaptic weight

$\theta$ (theta)와  $w$ 는 전형적으로 -0.05에서 +0.05 범위의 랜덤 값으로 초기화됨

이러한 초기 연결강도는 학습의 처음 단계에 사용됨

노드 3의 출력: 문제가 예측이면  $g$ 가 identity 함수이고, 분류이면  $g$ 가 로지스틱

$$output_j = g(\theta_j + \sum_{i=1}^p w_{ij} x_i)$$

$$output_3 = \frac{1}{1 + e^{-[-0.3 + (0.05)(0.2) + (0.01)(0.9)]}} = 0.43$$

# 신경망의 초기 통과

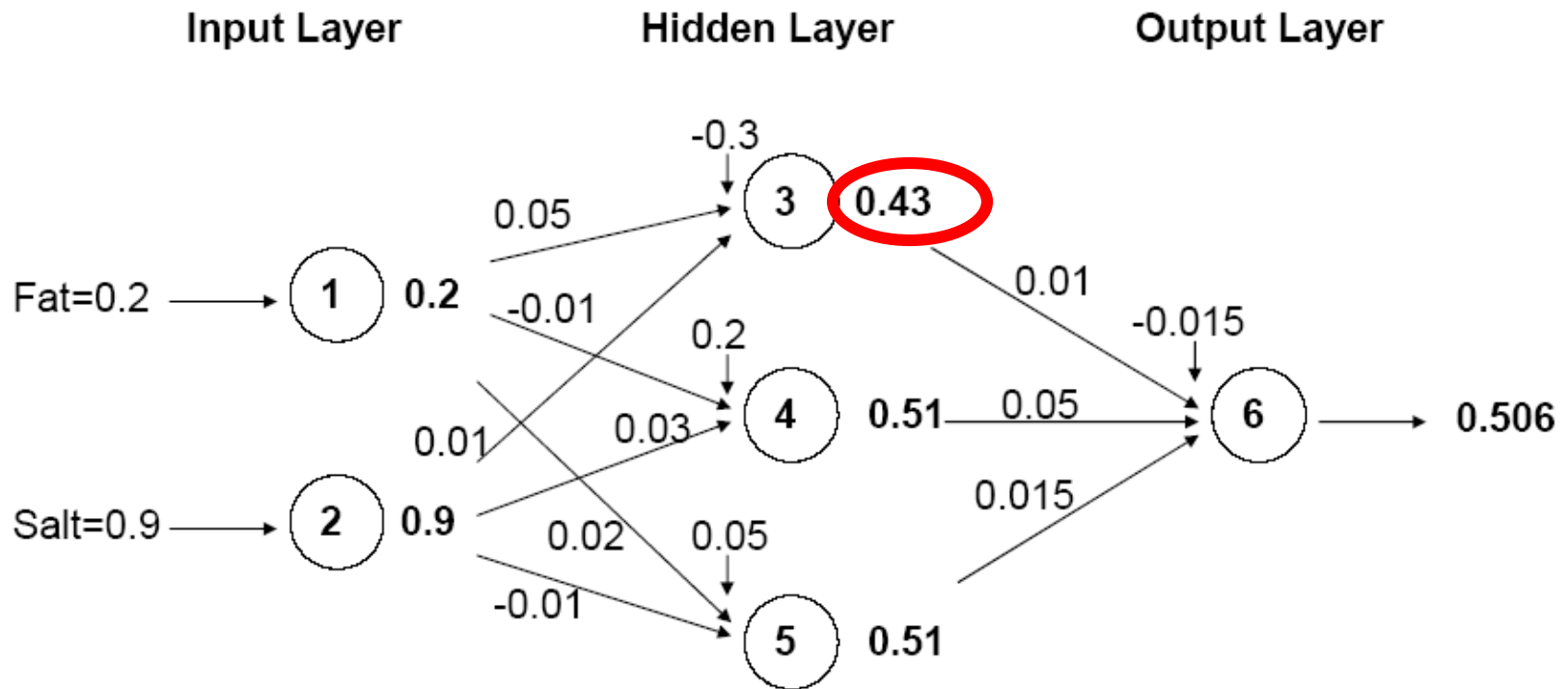


Figure 11.3: Computing node outputs (in boldface type) using the first observation in the tiny example and a logistic function.

# 출력층

마지막 중간층의 출력이 출력층의 입력이 됨

위와 같은 함수 사용, i.e. 가중평균의  $g$  함수

$$output_6 = \frac{1}{1 + e^{-[-0.015 + (0.01)(0.43) + (0.05)(0.507) + (0.015)(0.511)]}} = 0.506$$



# 출력 노드

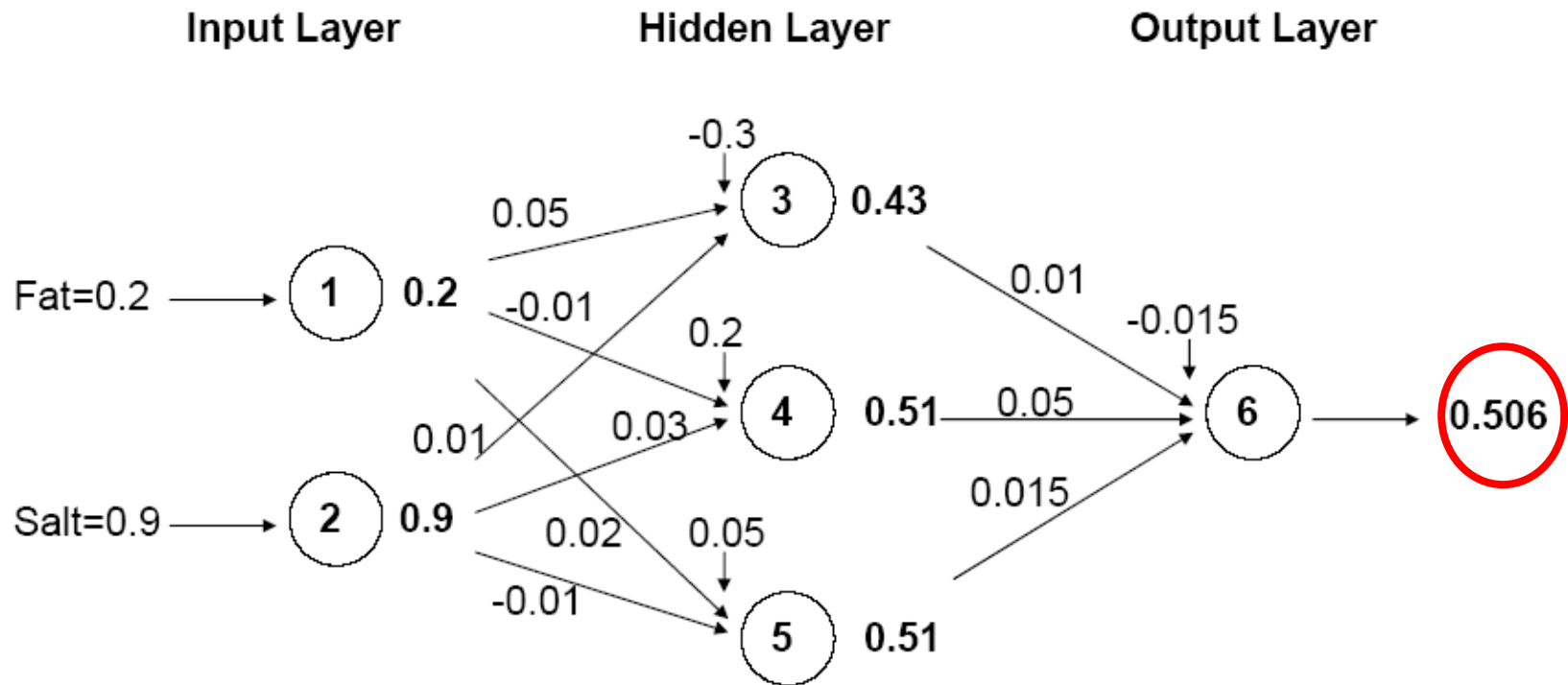


Figure 11.3: Computing node outputs (in boldface type) using the first observation in the tiny example and a logistic function.

# 모델 학습

# 전처리 단계

- 0-1로 변수 조정
- 범주형 변수
- Dummy variable 생성
- 비대칭 변수를 변환 (e.g., log)

# 학습

- 에러 함수를 최소로 하는  $w$  를 구하자
- 에러 함수 Error function
  - 에러 제곱 합: 에러 = “desired” output 과 “actual” output 의 차이

$$J(w) = \frac{1}{2} (y - \hat{y})^2$$

- $y$ : data,  $\hat{y}$ : model output
- 에러 함수는  $w$  에 대한 비선형 함수

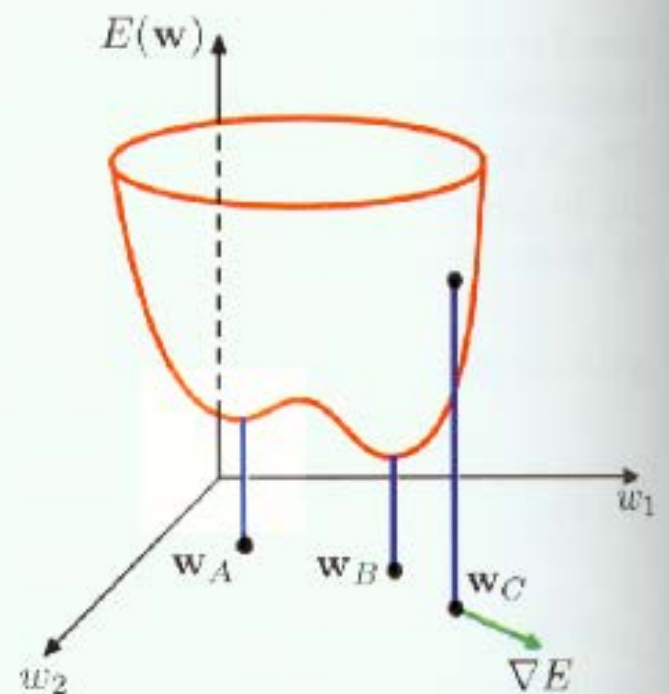
# 학습

- Steepest (Gradient) Descent = backpropagation
  - 1차 도함수를 구하고
  - $w$  를 도함수의 반대 방향으로 이동
  - $w$  가 더 이상 변화하지 않을 때까지 (즉, 도함수=0) 반복

$$\Delta w = -\eta \frac{\partial J}{\partial w}$$

# Training

**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  $\nabla E$ .



# 신경망 학습 back-propagation

목표: 에러제곱합을 최소화 하는 synaptic weight 찾기

- 학습 데이터  $x$  를 신경망에 입력
- 입력층에서 은닉층을 거쳐 출력층으로 시그널을 forward propagate 시킴
- 출력 노드 출력값 (=모델 출력값) 과 타겟값 차이가 출력 노드의 오류값으로 계산
- 이를 은닉층쪽으로 역전파 (back propagated) 하여 모든 은닉 노드들의 “오류”값으로 분배하고, 이를 연결강도 수정하는데 사용
- 이 과정을 모든 학습 데이터에 반복 적용

# Iteration 을 멈추는 조건들

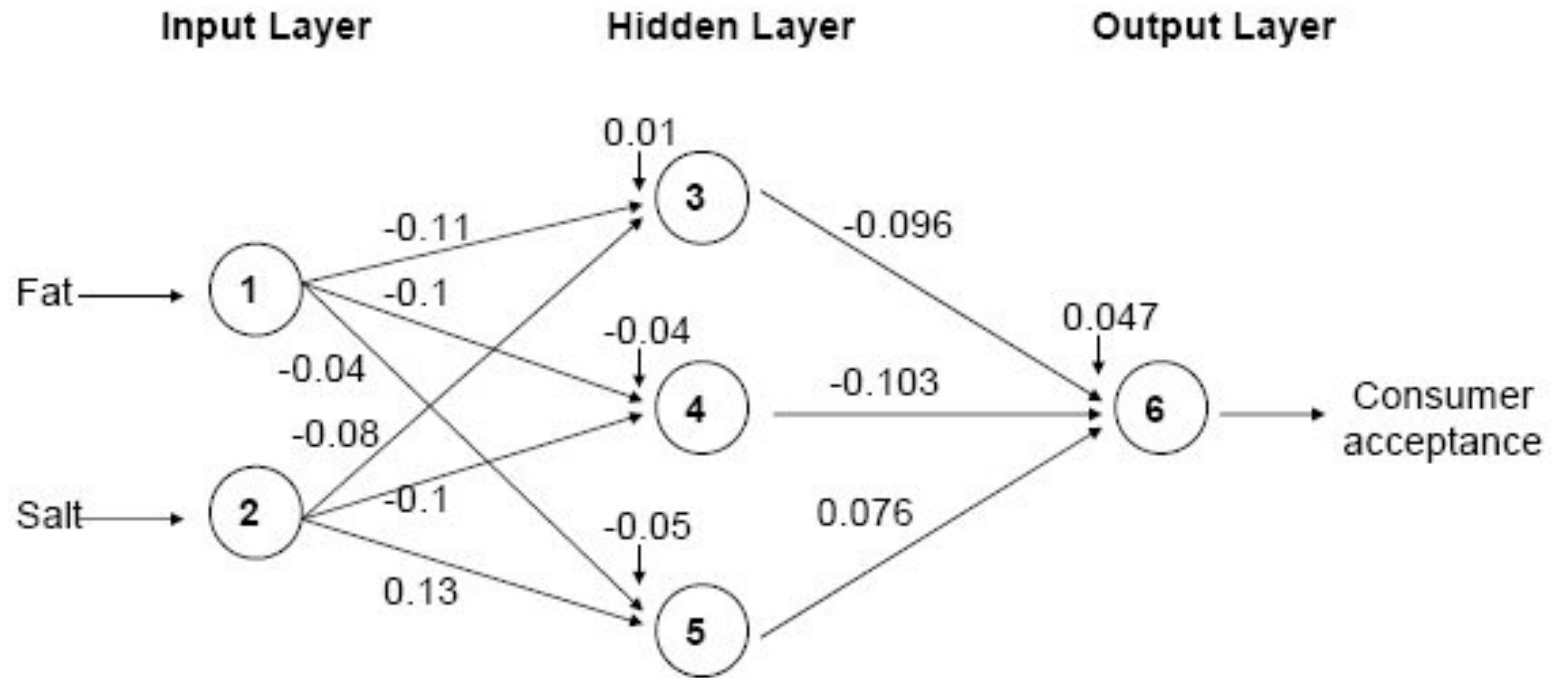
- 연결강도가 한 반복과 다음 반복 사이에 거의 변화하지 않을 때
- 오분류율이 요구된 목표치에 도달했을 때
- 실행횟수의 한계에 도달했을 때



# 다른 Training 방법

- Steepest Descent (=Back-propagation) 는 비선형 최적화 방법 가운데 가장 단순, biological plausibility
- 수 많은 다른 방법들 제안 및 시도 Conjugate Gradient 등
- 가장 인기 Levenberg-Marquadt (Newton's approximation)
  - 1,2차 도함수 모두 이용
  - Steepest Descent 비해 10~100배 빠름
  - 그러나 생물학적 의미는 전혀 없음

# 지방/염분 예: 마지막 연결강도



# 과적합 Overfitting

- With sufficient hidden nodes and training iterations, neural net can easily overfit the data
- 학습 데이터와 검증/테스트 데이터가 공유하는 내재 함수 성질만 학습하는 것이 아니라,
- 학습 데이터만이 가지고 있는 특이성(noise)도 학습

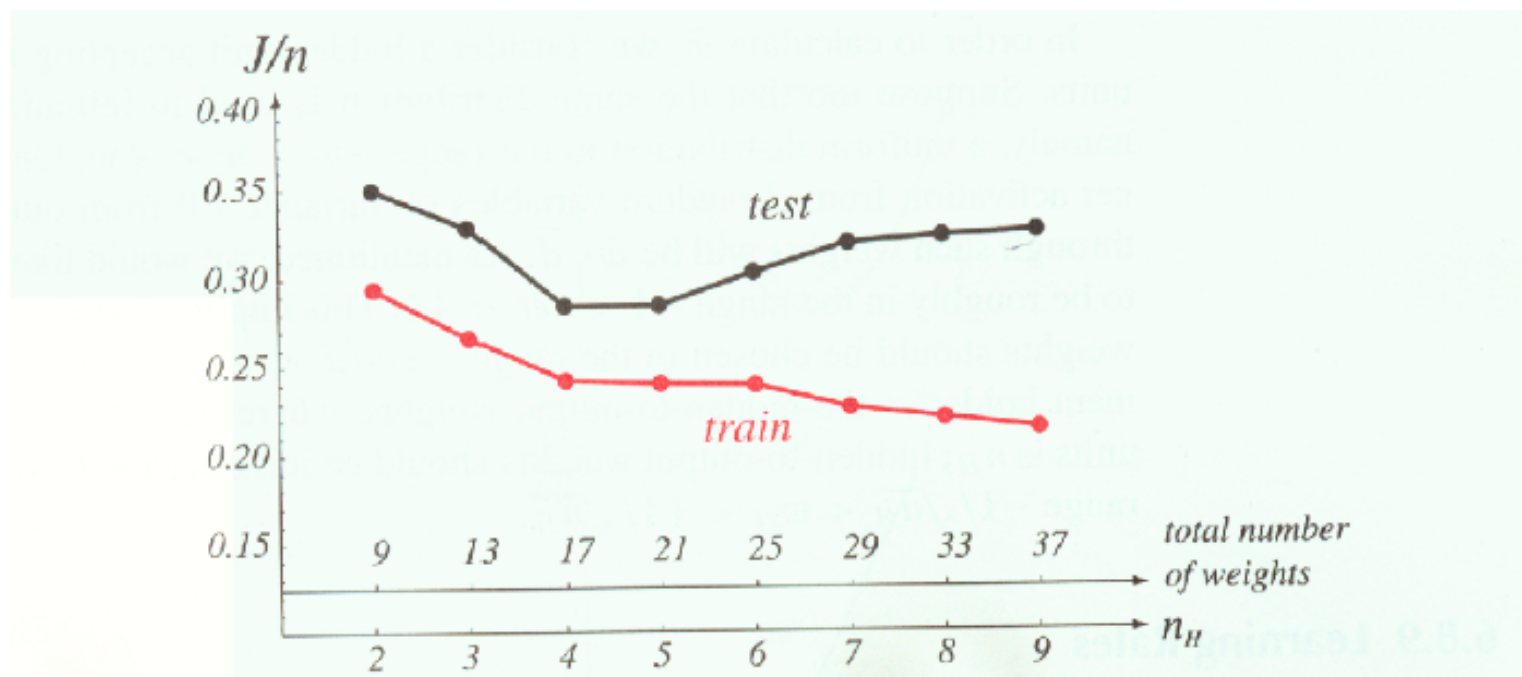
# Avoiding Overfitting 원인 제거

1. Too many hidden nodes (모델이 너무 복잡)

Limit complexity of network

다양한 hidden node 수를 검증 validation

# “최적”의 히든 노드 수

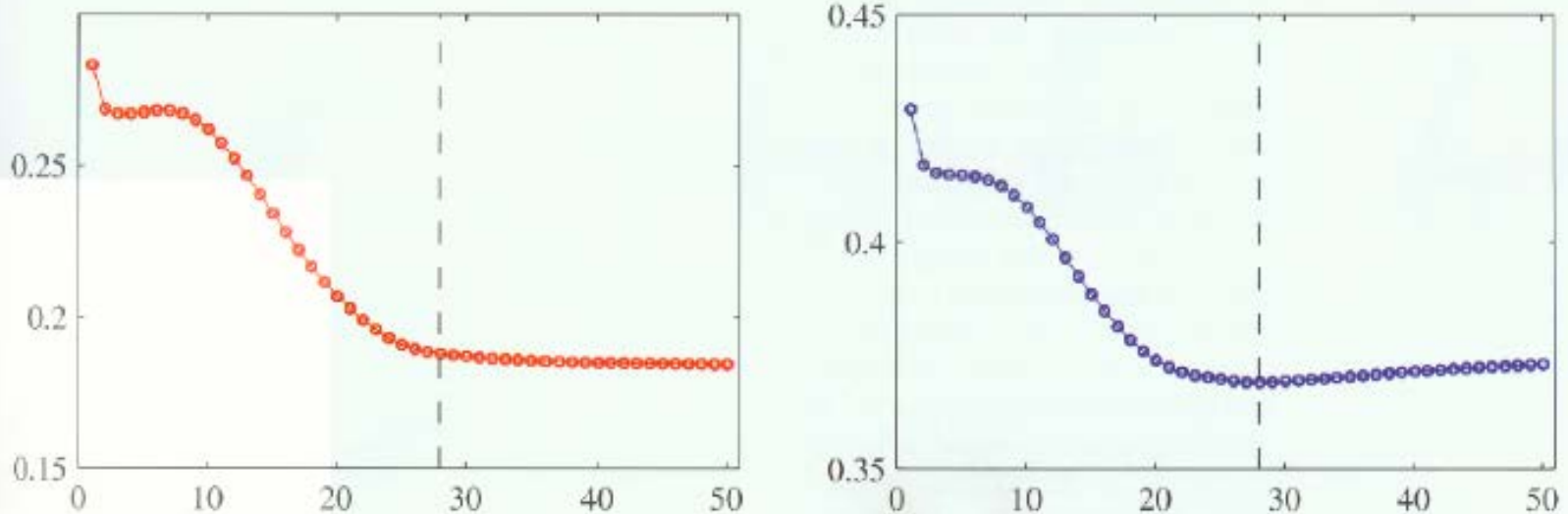


**FIGURE 6.15.** The error per pattern for networks fully trained but differing in the numbers of hidden units,  $n_H$ . Each  $2 - n_H - 1$  network with bias was trained with 90 two-dimensional patterns from each of two categories, sampled from a mixture of three Gaussians, and thus  $n = 180$ . The minimum of the test error occurs for networks in the range  $4 \leq n_H \leq 5$ , i.e., the range of weights 17 to 21. This illustrates the rule of thumb that choosing networks with roughly  $n/10$  weights often gives low test error.

# Avoiding Overfitting

1. Too long training (모델의 복잡도를 100% 사용)  
학습 도중에 계속적으로 validation error 측정  
도중에 학습 중지 early stopping

# Early stopping



**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, for the sinusoidal data set. The goal of achieving the best generalization performance suggests that training should be stopped at the point shown by the vertical dashed lines, corresponding to the minimum of the validation set error.

# 모델링, 구조 결정



# 망 구조 결정

## 중간층의 개수

- 가장 일반적인 경우 - 하나의 은닉층 (cf Deep Belief Network)

## 중간층에서 노드의 개수

- 노드가 많을수록 복잡성을 잡아내기 쉽지만, 과적합의 가능성도 높아진다.

## 출력노드의 개수

- 분류에서, 클래스마다 하나의 노드(또한 이항 응답에서 하나를 사용할 수 있다)
- 수치형 예측에서 하나를 사용

# 망 구조 (계속)

## “학습률” (l)

- 낮은 값은 각 반복에서의 오류로부터 새로운 정보의 “반영도를 줄인다”.
- 이는 학습을 늦추지만, 국부구조에서 과적합 경향을 축소시킨다.

## “관성”

- 높은 값은 이전 반복에서와 같은 방향으로 연결강도를 계속 변화시킨다.
- 이 또한, 국부구조에서 과적합을 피하도록 돕지만, 또한 학습을 저하시킨다.

# 자동화

- 몇몇 소프트웨어는 입력 파라미터의 최적선택을 자동화한다.
- XLMiner는 중간층의 개수와 노드의 개수를 자동화한다.

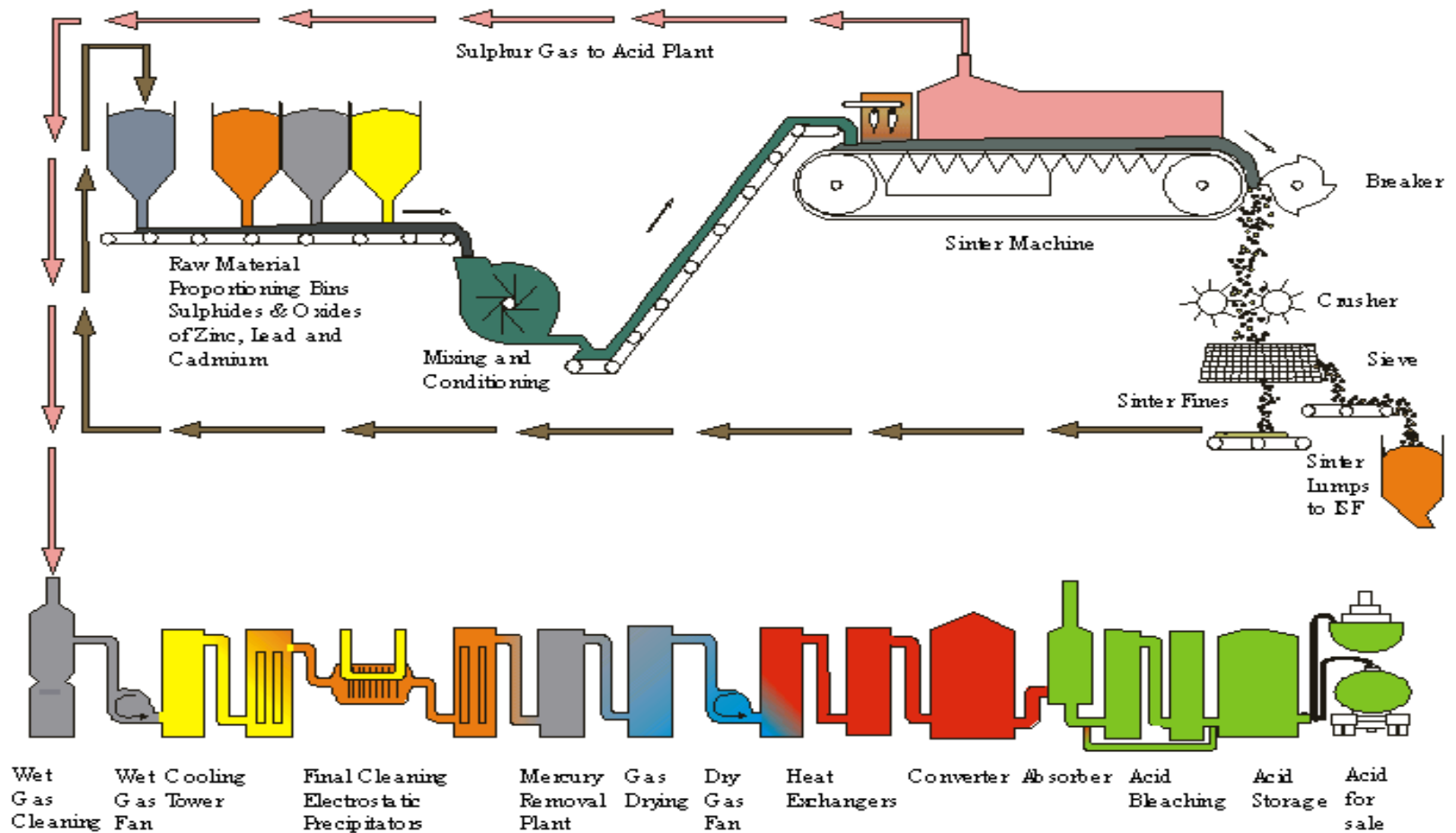
사례

# 소결 공정

- 분말 형태의 원광석은 고로에 넣어서 생산할 수 없음
- 이를 덩어리 형태로 구워내는 공정



# 소결 공정의 펠릿 속도 제어 자동화



# 소결 공정의 펄릿 속도 제어 자동화

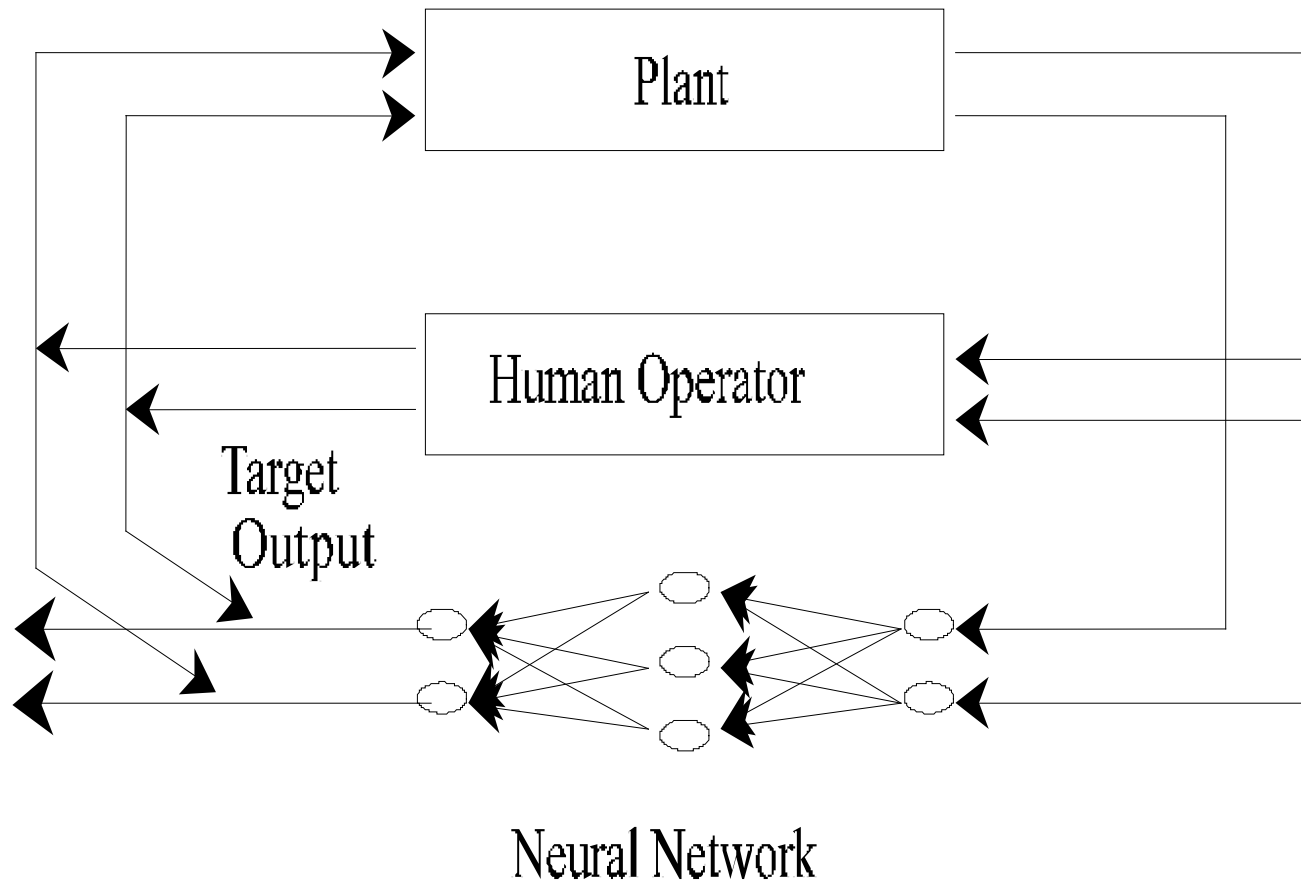
- 펄릿의 속도 제어가 가장 중요한 태스크
  - 너무 느리면: 품질 저하, 그리고 먼지, CO 및 CO2 증가
  - 너무 빠르면: 생산량 저하
  - 최적의 속도로 운영하는 것이 생산량, 품질, 비용, 환경 모든면에 중요함

# 소결 공정의 팰릿 속도 제어 자동화

- 경력 20~30년의 조업자가 매뉴얼로 제어하고 있음
  - 공정 유지 수준의 일관성 확보 어려움
  - 조업자 은퇴 후, 후계자가 없음
- 자동화 모델 필요하여 조업 데이터로부터 조업자와 같은 제어를 할 수 있는 제어 모델 구축 (copy control, 다음 페이지 그림 참조)
- 속도 =  $f$  (원자재 관련 변수 120 개 => 5개) 로 모델 구축 (다음 페이지 그림 참조)



# 소결 공정 Copy Controller



# 소결공정: 중요 변수 선택

Significance level (%)	Factor	Nermal Range	Update freque-ncy	NN con-trol-ler
19.9	Hot zone properties (height)	180±5cm	1 min	
12.9	Windbox temperature profile	400±20°C	1 min	O
12.4	Bed height	550-10mm	constant	
8.7	Quick lime ratio	1.5±0.2%	10 days	O
8.2	Main Blower air pressure	1525±25 mmH <sub>2</sub> O	1 min	O
6.3	Main Blower power	2.9MW±0.1MW	1 min	O
6.1	Coke size	1.40mm±0.2mm	10 days	
5.9	SI	91±0.5%	4 hours	
5.2	Return fine ratio	13+2 %	10 days	
3.7	Screening ratio	17.5±2.5%	10 days	
3.5	Blending material moisture	6.0±0.3 %	1 min	O
2.3	Coke ratio	3.7±0.1%	10 days	
2.3	Feeding density	1.955±0.05ton/m <sup>3</sup>	2 ~ 3 days	
1.2	RDI	33.2+2.0%	8 hours	

# 냉간 압연



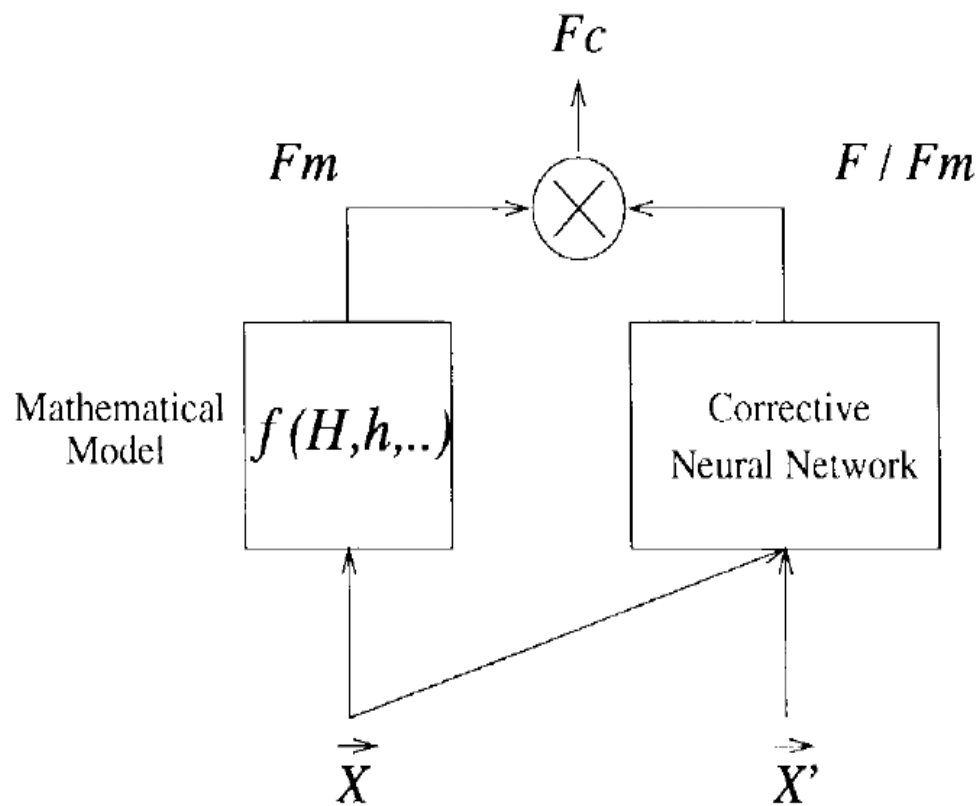
- 상온의 금속재료를, 회전하는 2개의 롤 사이로 통과시켜서 여러 가지 형태의 재료, 즉 판(板)·봉(棒)·관(管)·형재(形材) 등으로 가공
- 주어진 강판을 원하는 두께로 만들기 위해 적용해야 할 Roll의 최적 압력 계산 필요하나 현재 사용하는 물리적 모델은 최적 압력을 계산해내지 못함

# 냉간 압연

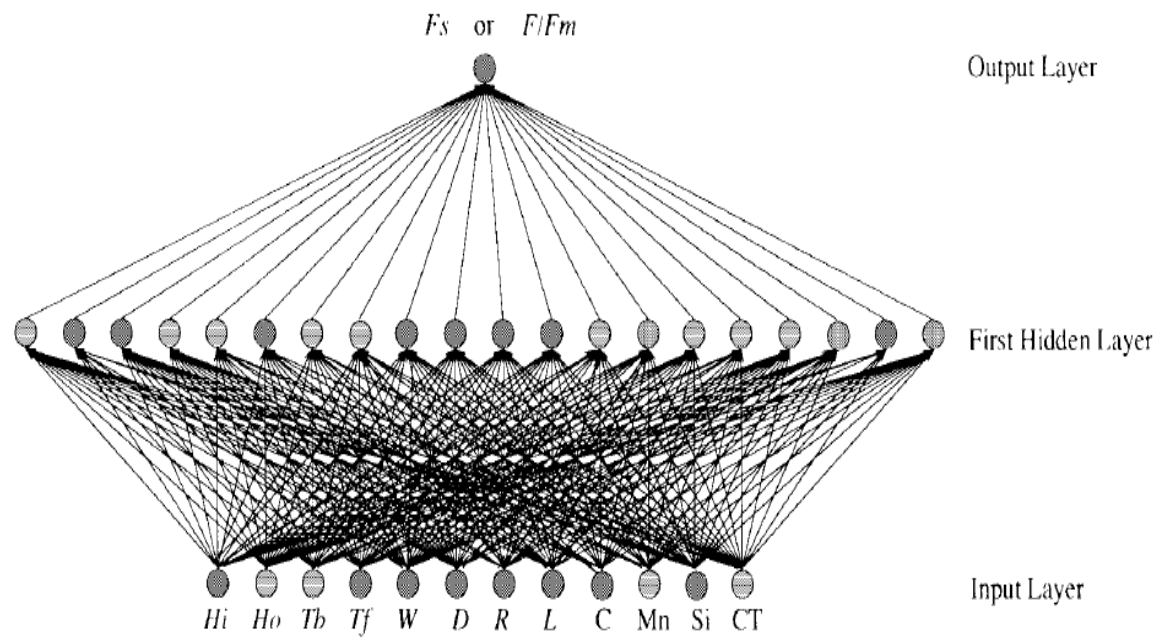


- 물리적 모델의 압력 보정 값을 타겟으로 하는 신경망 모델을 조업 데이터로부터 구축.
- 즉, 압력\_보정치 = NN (강판 두께, 전후 텐션, 넓이, 롤 지름, 성분-탄소, 망간, 실리콘)
- 이를 이용하여
  - 최적 Roll 압력치 (=수식기반 압력 + 압력\_보정치) 계산하여
  - 최적 조업 수행 가능

# 보정 모델



# 신경회로망



# 예측 결과

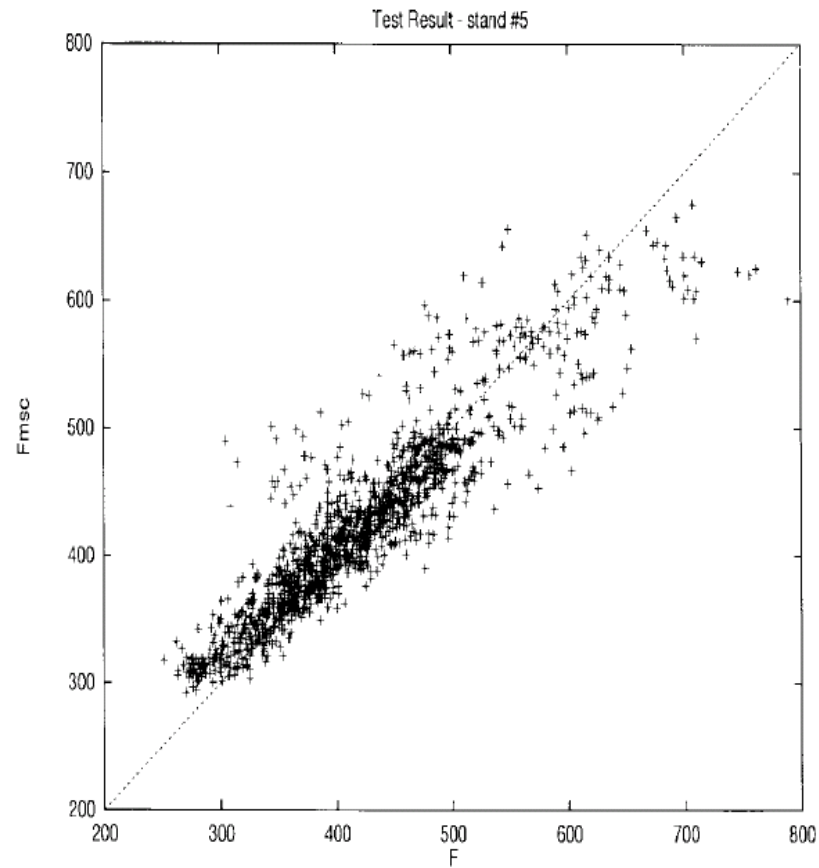
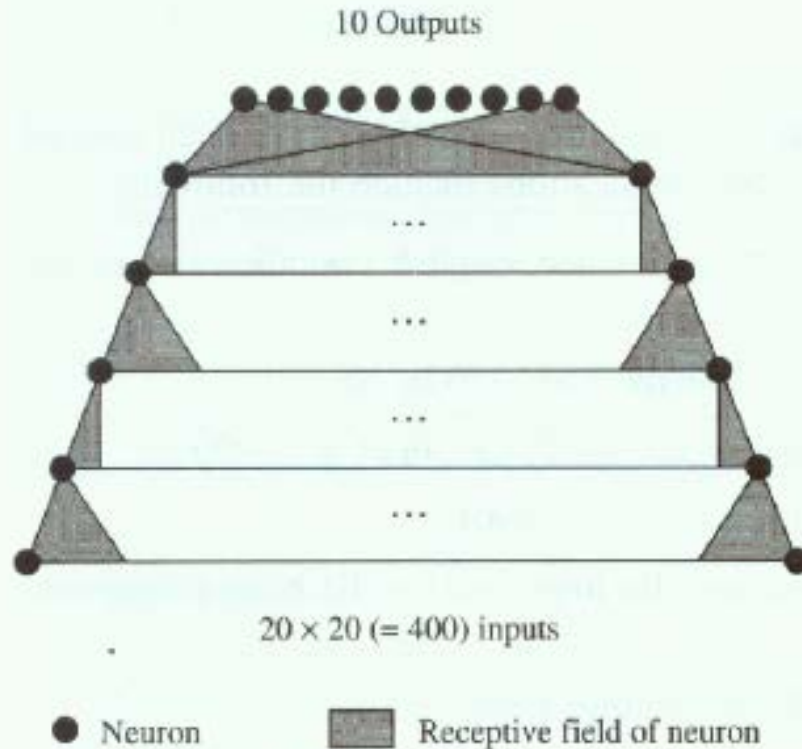


Fig. 11. Predictions of the M-S-C model on *FieldTestSet* at stand no. 5.

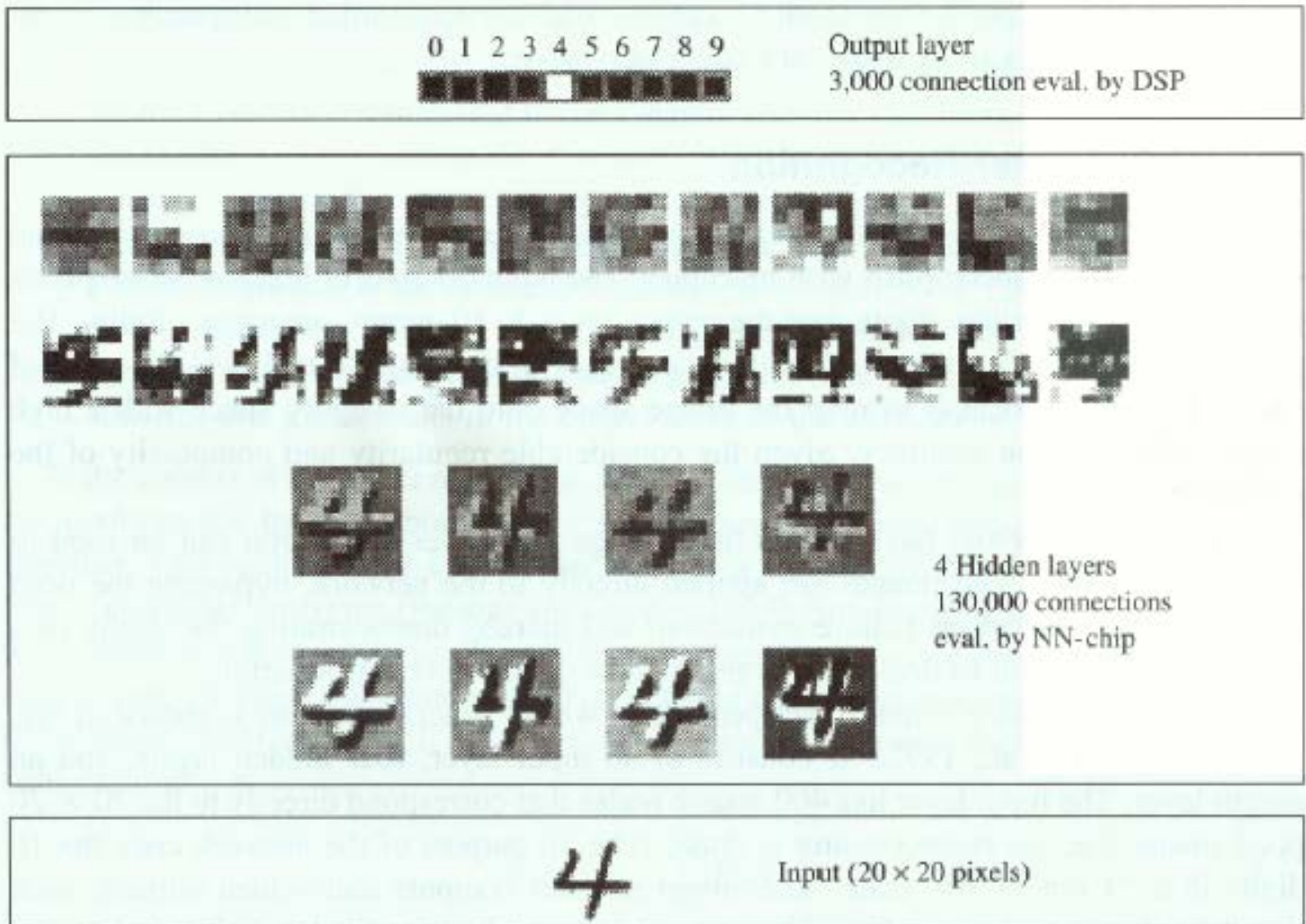
# 숫자 자동 인식



Layer	Neurons	Synapses
5	10	3,000
4	300	1,200
3	1,200	50,000
2	784	3,136
1	3,136	78,400

**FIGURE 6.30** General structure of the optical character recognition (OCR) network.  
(From E. Säckinger et al., 1992a, with permission of IEEE.)





**FIGURE 6.31** Example for the states of the OCR network with a number four as input. (From E. Säckinger et al., 1992a, with permission of IEEE.)

# 결론

# 장점

- 예측 능력이 좋음
- 복잡한 관계를 잡아낼 수 있음

# 단점

- “블랙박스” 예측변수와 결과 사이의 관계에 대한 직관을 제공하지 못한다.
- 변수선택 메커니즘이 없다.
- 변수가 많다면 계산량이 늘어 남 (특히 dummy var 은 시냅스의 수를 극적으로 증가시킨다).

# 요약

- 신경망은 분류와 예측에 사용될 수 있다.
- 타겟 변수와 일련의 예측변수들 사이의 매우 유연하고 복잡한 관계를 잡아낼 수 있다.
- 가장 강력한 예측 모델
- 주요 위험: 과적합
- 대용량 데이터가 필요하다.
- 예측성능이 좋지만, 본질상 “블랙박스”

# 요약

- 수 십 종류의 신경망 모델 존재함
- 이 가운데 가장 많이 사용되는 모델인 “다층 퍼셉트론”을 소개함.
- 다른 교사학습 신경망 모델들은
  - Single layer perceptron
  - Radial basis function networks
  - Probabilistic neural network
  - Recurrent network
- 비교사 학습 신경망 모델들은
  - Self Organizing Map network
  - Hopfield network , Boltzmann Machine