

# 다층 퍼셉트론

- 퍼셉트론은 선형 분류기라는 한계
- 선형 분리 불가능한 상황에서는 일정한 양의 오류

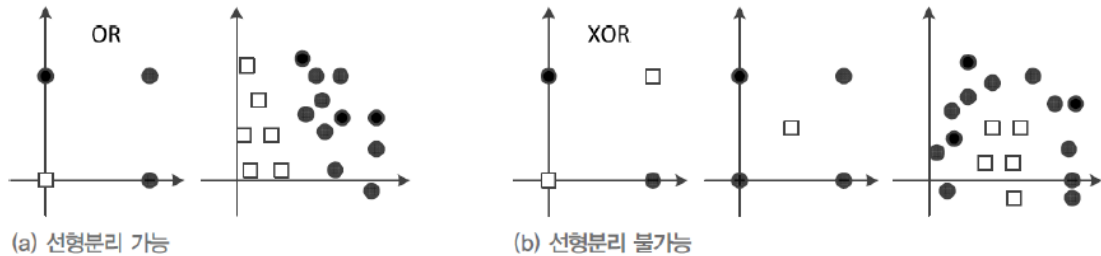


그림 3-7 선형분리가 가능한 상황과 불가능한 상황

- XOR 문제에서는 75%가 정확률 한계
- 1969년 민스키의 Perceptrons : 극복 방안 제시, 당시 기술로 실현 불가능
- 1974년 웨어보스의 박사 논문에서 오류 역전파 알고리즘 제안
- 1986년 루멜하트의 저서 Parallel Distributed Processing 다층 퍼셉트론 이론 정립, 신경망 부활
- **핵심 아이디어**
  - 은닉층을 둔다. 은닉층은 원래 특징 공간을 분류하는 데 훨씬 유리한 새로운 특징 공간으로 변환
  - 시그모이드 활성화함수를 도입한다
    - 계단함수를 활성화함수로 사용하면 경성 의사결정에 해당
    - 다층 퍼셉트론은 연성 의사결정이 가능한 **시그모이드 활성화함수** 사용
      - **출력이 연속값**이고, 출력을 신뢰도로 간주함으로써 더 융통성 있게 의사결정
  - 오류 역전파 알고리즘 사용
    - 역방향으로 진행하면서 한 번에 한 층 씩 **그레디언트를 계산하고 가중치를 갱신**

## 특징 공간 변환

- 퍼셉트론 2개를 사용한 XOR 문제의 해결
  - 퍼셉트론 1과 퍼셉트론 2가 모두 +1이면 ○ 부류, 아니면 □ 부류

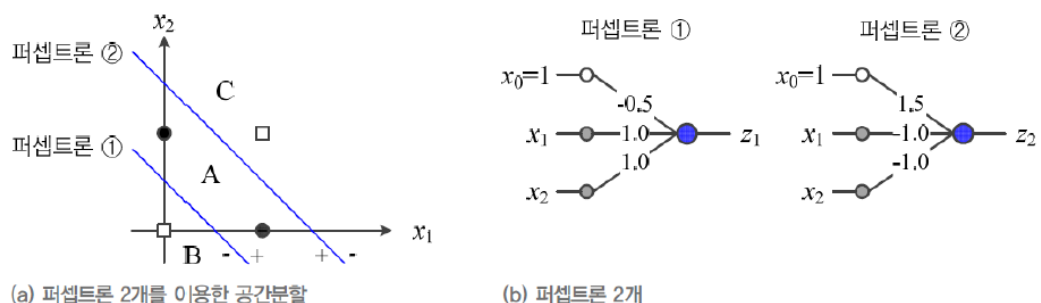
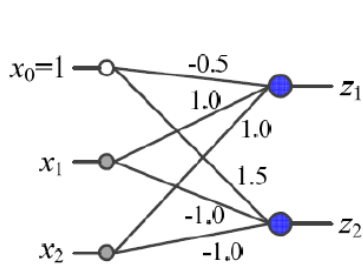
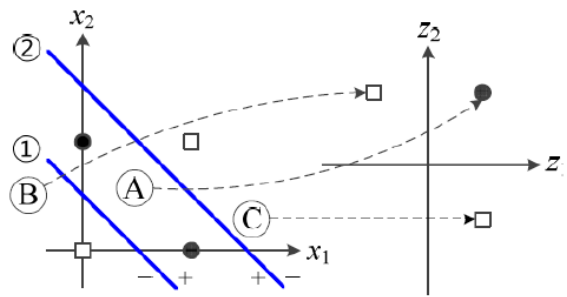


그림 3-8 XOR 문제의 해결

- 퍼셉트론 2개를 병렬 결합하면
  - 원래 공간  $\mathbf{x} = (x_1, x_2)^T$ 를 새로운 특징 공간  $\mathbf{z} = (z_1, z_2)^T$ 로 변환
  - 새로운 특징 공간  $\mathbf{z}$ 에서는 선형 분리 가능함



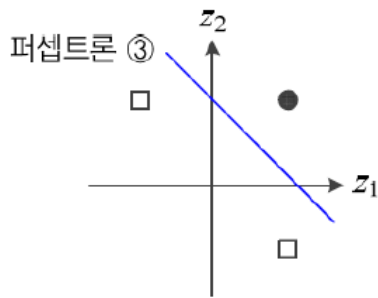
(a) 두 퍼셉트론을 병렬로 결합



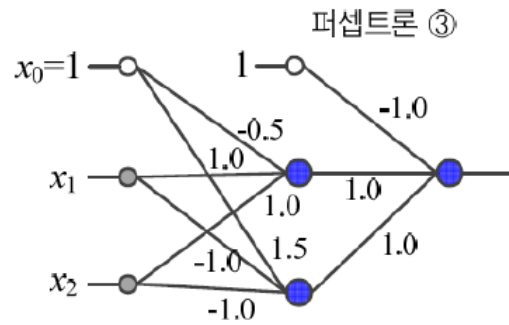
(b) 원래 특징 공간  $x$ 를 새로운 특징 공간  $z$ 로 변환

그림 3-9 특징 공간의 변환

- 사람이 수작업으로 특징 학습을 수행한 것과 유사함
- 이후, 퍼셉트론 1개를 순차 결합하면
  - 새로운 특징 공간  $z$ 에서 선형 분리를 수행하는 퍼셉트론 3을 순차 결합하면 다층 퍼셉트론



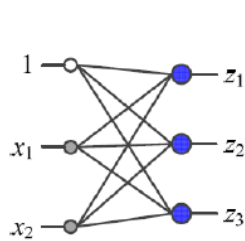
(a) 새로운 특징 공간에서 분할



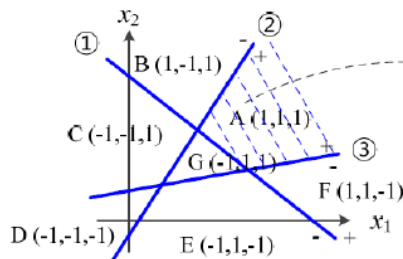
(b) 퍼셉트론 3개를 결합한 다층 퍼셉트론

그림 3-10 다층 퍼셉트론

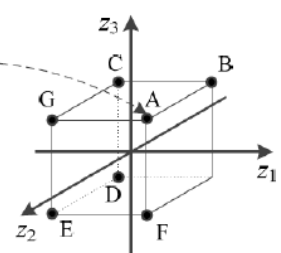
- 이 다층 퍼셉트론은 훈련집합에 있는 4개 샘플 (0,0) (0,1) (1,0) (1,1)을 제대로 분류하는가?
- 다층 퍼셉트론의 용량
  - 3개 퍼셉트론을 결합하면, 2차원 공간을 7개 영역으로 나누고 각 영역을 3차원 점으로 변환
  - 활성화함수  $\tau$ 로 계단함수를 사용하므로 영역을 점으로 변환



(a) 퍼셉트론 3개를 결합



(b) 7개 부분공간으로 나눔



(c) 3차원 공간의 점으로 매핑

그림 3-11 퍼셉트론을 3개 결합했을 때 공간 변환

- 일반화하여,  $p$ 개 퍼셉트론을 결합하면  $p$ 차원 공간으로 변환
  - $1 + \sum_{i=1}^p i$  개의 영역으로 분할

구 조	결정 구역	Exclusive-or	Classes with Meshed Regions	Most General Region Shapes
Single-layer j i	Half Plane Bounded by Hyperplane			
Two-layer k j i	Convex Open or Closed Regions			
Three-layer l k j i	Arbitrary (Complexity limited by Number of Units)			

## 활성 함수

- 딱딱한 공간 분할과 부드러운 공간 분할
  - 계단함수는 딱딱한 의사결정(영역을 점으로 변환)
  - 나머지 활성함수는 부드러운 의사결정(영역을 영역으로 변환)

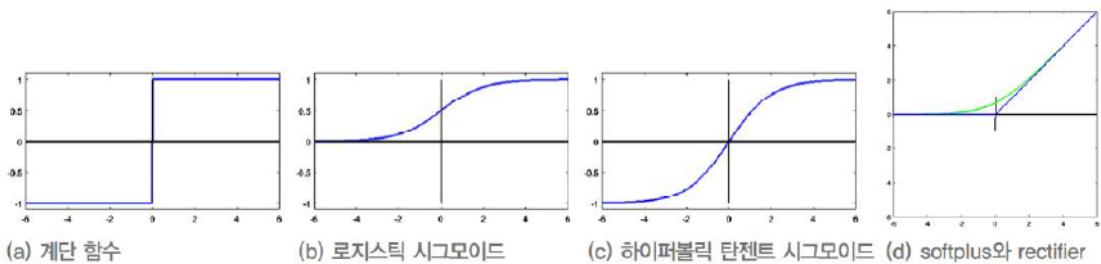


그림 3-12 신경망이 사용하는 활성함수

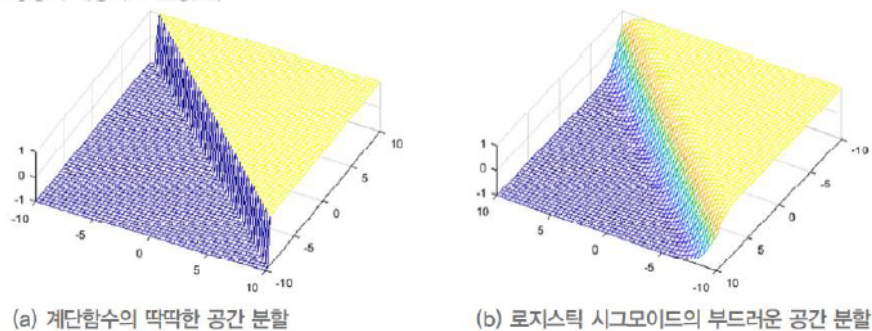


그림 3-13 퍼셉트론의 공간 분할 유형

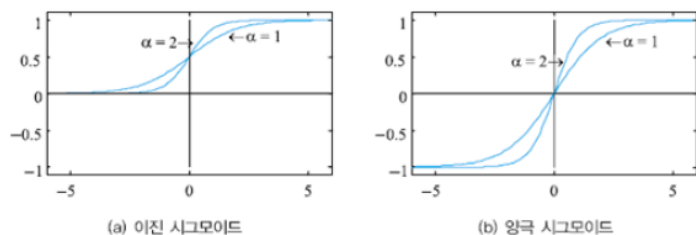
- 대표적인 비선형 함수인 시그모이드를 활성함수로 사용

이진 시그모이드 함수:

$$\left. \begin{aligned} \tau_1(x) &= \frac{1}{1 + e^{-\alpha x}} \\ \tau_1'(x) &= \alpha \tau_1(x)(1 - \tau_1(x)) \end{aligned} \right\}$$

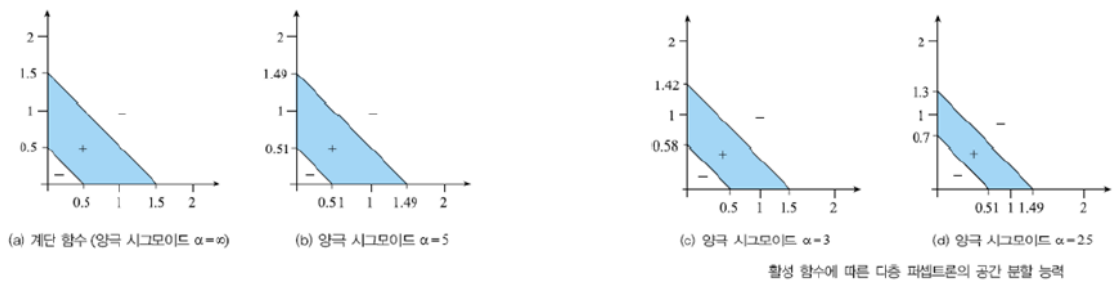
양극 시그모이드 함수:

$$\left. \begin{aligned} \tau_2(x) &= \frac{2}{1 + e^{-\alpha x}} - 1 \\ \tau_2'(x) &= \frac{\alpha}{2} (1 + \tau_2(x))(1 - \tau_2(x)) \end{aligned} \right\}$$



활성 함수로 널리 사용되는 두 가지 시그모이드 함수

- 활성 함수에 따른 다층 퍼셉트론의 공간 분할 능력 변화



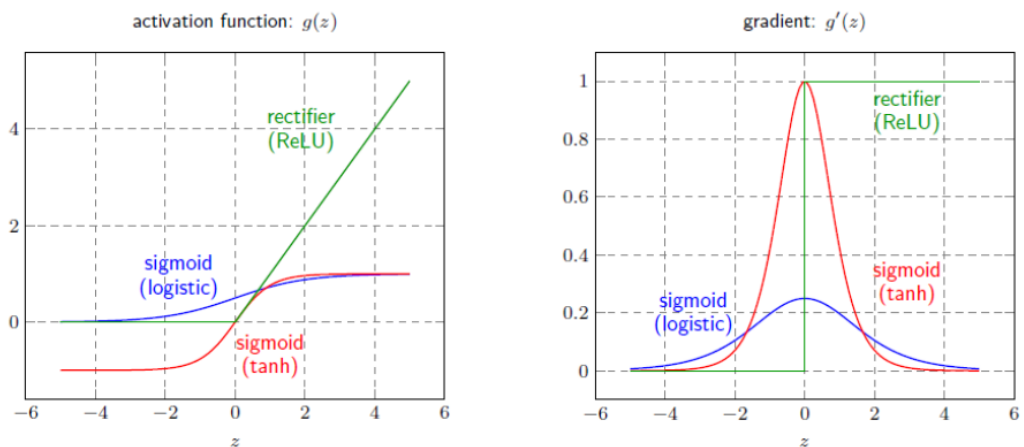
## ● 신경망이 사용하는 다양한 활성화 함수

- 로지스틱 시그모이드와 하이퍼볼릭 탄젠트는  $a$ 가 커질수록 계단함수에 가까워짐
- 모두 1차 도함수 계산이 빠름 (특히 ReLU는 비교 연산 한 번)

표 3-1 활성화 함수로 사용되는 여러 함수

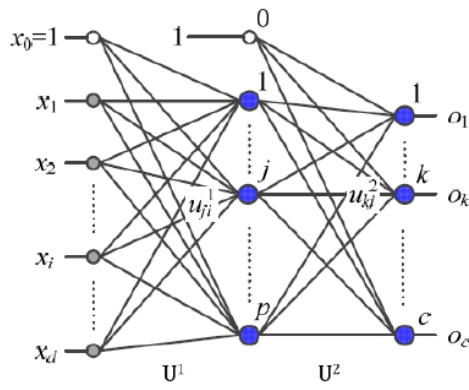
함수 이름	함수	1차 도함수	범위
계단	$\tau(s) = \begin{cases} 1 & s \geq 0 \\ -1 & s < 0 \end{cases}$	$\tau'(s) = \begin{cases} 0 & s \neq 0 \\ \text{불가} & s = 0 \end{cases}$	-1과 1
로지스틱 시그모이드	$\tau(s) = \frac{1}{1 + e^{-as}}$	$\tau'(s) = a\tau(s)(1 - \tau(s))$	(0,1)
하이퍼볼릭 탄젠트	$\tau(s) = \frac{2}{1 + e^{-as}} - 1$	$\tau'(s) = \frac{a}{2}(1 - \tau(s)^2)$	(-1,1)
소프트플러스	$\tau(s) = \log_e(1 + e^s)$	$\tau'(s) = \frac{1}{1 + e^{-s}}$	$[0, \infty)$
렉티파이어(ReLU)	$\tau(s) = \max(0, s)$	$\tau'(s) = \begin{cases} 0 & s < 0 \\ 1 & s > 0 \\ \text{불가} & s = 0 \end{cases}$	$[0, \infty)$

- 퍼셉트론은 계단함수, 다층 퍼셉트론은 로지스틱 시그모이드와 하이퍼볼릭 탄젠트, 딥러닝은 **ReLU**(Rectified Linear Activation)를 주로 사용
- 일반적으로 은닉층에서 **로지스틱 시그모이드**를 활성화 함수로 많이 사용
  - 시그모이드의 넓은 포화곡선은 그래디언트 기반한 학습을 어렵게 함
  - ReLU 대두됨

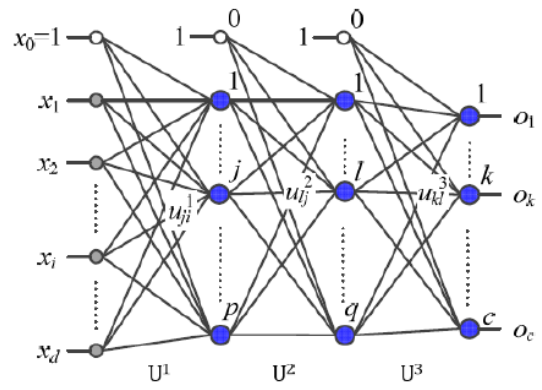


## 구조

- 입력층-은닉층-출력층의 2층 구조
  - $d+1$ 개의 입력 노드( $d$ 는 특징의 개수),  $c$ 개의 출력 노드( $c$ 는 부류 개수)
  - $p$ 개의 은닉 노드 :  $p$ 는 하이퍼 매개변수(사용자가 정해주는 매개변수)
    - $p$ 가 너무 크면 과잉적합, 너무 작으면 과소적합 ▶ 하이퍼 매개변수 최적화



(a) 2층 퍼셉트론



(b) 3층 퍼셉트론

- b는 3층구조
- 다층 퍼셉트론의 매개변수 (가중치)

- 입력층 - 은닉층을 연결하는

$U^1$  ( $u^1_{ji}$ 은 입력층의  $i$ 번째 노드를 은닉층의  $j$ 번째 노드와 연결)

- 은닉층 - 출력층을 연결하는

$U^2$  ( $u^2_{kj}$ 은 은닉층의  $j$ 번째 노드를 출력층의  $k$ 번째 노드와 연결)

2층 퍼셉트론의 가중치 행렬:

$$U^1 = \begin{pmatrix} u^1_{10} & u^1_{11} & \cdots & u^1_{1d} \\ u^1_{20} & u^1_{21} & \cdots & u^1_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ u^1_{p0} & u^1_{p1} & \cdots & u^1_{pd} \end{pmatrix}, \quad U^2 = \begin{pmatrix} u^2_{10} & u^2_{11} & \cdots & u^2_{1p} \\ u^2_{20} & u^2_{21} & \cdots & u^2_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ u^2_{c0} & u^2_{c1} & \cdots & u^2_{cp} \end{pmatrix} \quad (3.11)$$

- 일반화하면  $u^l_{ji}$ 은  $l-1$ 번째 은닉층의  $i$ 번째 노드를  $l$ 번째 은닉층의  $j$ 번째 노드와 연결하는 가중치
  - 입력층을 0번째 은닉층, 출력층을 마지막 은닉층으로 간주

## 동작

- 특징 벡터  $x$ 를 출력 벡터  $o$ 로 매핑하는 함수로 간주할 수 있음

$$\left. \begin{aligned} \text{2층 퍼셉트론: } \mathbf{o} &= \mathbf{f}(\mathbf{x}) = \mathbf{f}_2(\mathbf{f}_1(\mathbf{x})) \\ \text{3층 퍼셉트론: } \mathbf{o} &= \mathbf{f}(\mathbf{x}) = \mathbf{f}_3(\mathbf{f}_2(\mathbf{f}_1(\mathbf{x}))) \end{aligned} \right\} \quad (3.12)$$

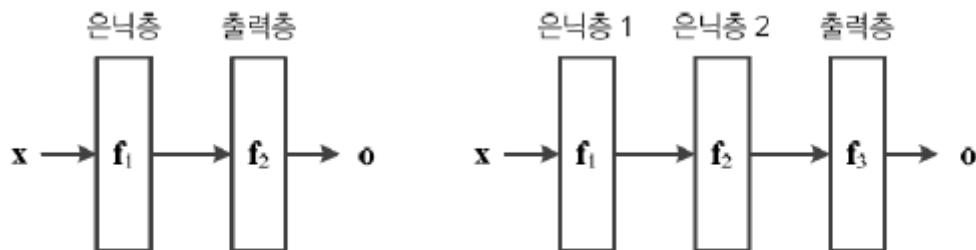


그림 3-15 다층 퍼셉트론을 간략화한 구조

- 깊은 신경망은  $\mathbf{o} = \mathbf{f}_L(\cdots \mathbf{f}_2(\mathbf{f}_1(\mathbf{x})))$ ,  $L \geq 4$  (딥러닝)

- 노드가 수행하는 연산을 구체적으로 쓰면

$j$ 번째 은닉 노드의 연산:

$$z_j = \tau(zsum_j), j = 1, 2, \dots, p \quad (3.13)$$

이때  $zsum_j = \mathbf{u}_j^1 \mathbf{x}$ 이고  $\mathbf{u}_j^1 = (u_{j0}^1, u_{j1}^1, \dots, u_{jd}^1)$ ,  $\mathbf{x} = (1, x_1, x_2, \dots, x_d)^T$

$k$ 번째 출력 노드의 연산:

$$o_k = \tau(osum_k), k = 1, 2, \dots, c \quad (3.14)$$

이때  $osum_k = \mathbf{u}_k^2 \mathbf{z}$ 이고  $\mathbf{u}_k^2 = (u_{k0}^2, u_{k1}^2, \dots, u_{kp}^2)$ ,  $\mathbf{z} = (1, z_1, z_2, \dots, z_p)^T$

- $\mathbf{u}_j^1$ 은  $j$ 번째 은닉 노드에 연결된 가중치 벡터 (식 (3.11)의  $\mathbf{U}^1$ 의  $j$ 번째 행)
- $\mathbf{u}_k^2$ 은  $k$ 번째 출력 노드에 연결된 가중치 벡터 (식 (3.11)의  $\mathbf{U}^2$ 의  $k$ 번째 행)

- 다층 퍼셉트론의 동작을 행렬로 표기하면

$$\mathbf{o} = \tau(\mathbf{U}^2 \tau_h(\mathbf{U}^1 \mathbf{x})) \quad (3.15)$$

- 은닉층은 특징 추출기

- 은닉층은 특징 벡터를 분류에 더 유리한 새로운 특징 공간으로 변환
- 현대 기계 학습에서는 특징 학습이라 (feature learning 혹은 data-driven features) 부름
  - 딥러닝은 더 많은 단계를 거쳐 특징학습을 함

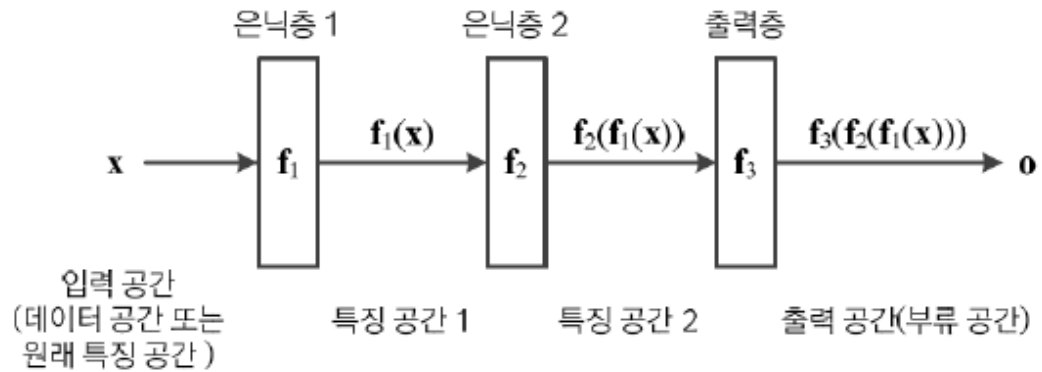
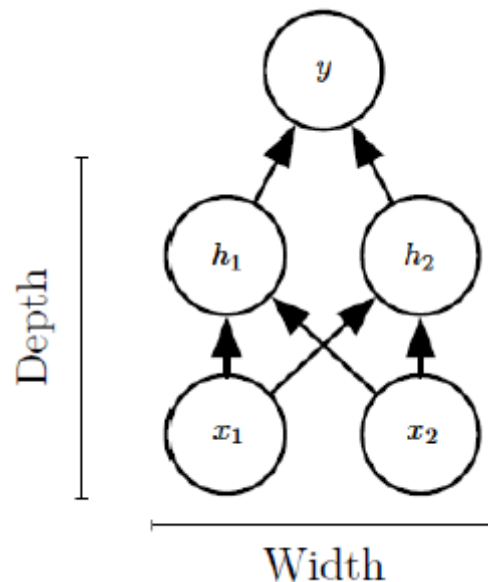


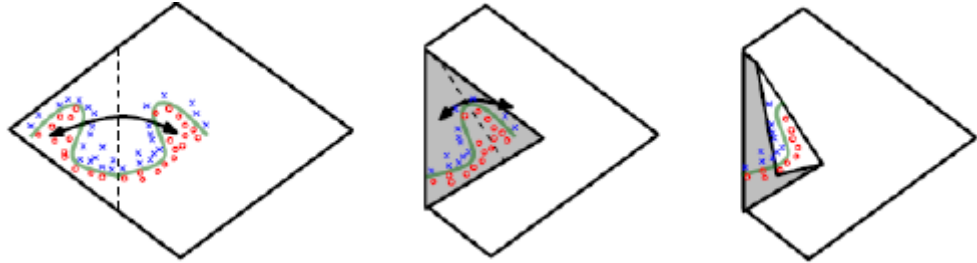
그림 3-16 특징 추출기로서의 은닉층

- 기본 구조

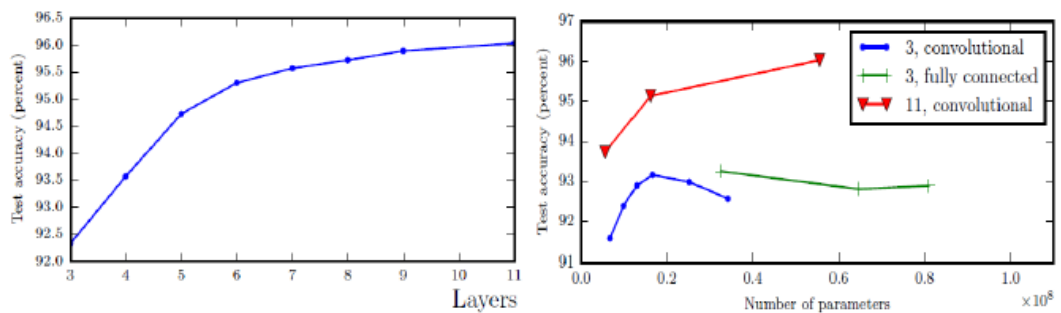


- 범용적 근사 이론

- 하나의 은닉층은 함수의 근사를 표현
- 다층 퍼셉트론도 공간을 변환하는 함수를 근사함
- 얇은 은닉층의 구조보다
  - 지수적으로 더 넓은 폭이 필요할 수 있음
  - 더 과잉적합 되기 쉬움
  - 일반적으로 깊은 은닉층의 구조가 좋은 성능을 가짐
- 은닉층의 깊이에 따른 이점
  - 지수의 표현



- 각 은닉층은 입력 공간을 어디서 접을지 지정 ▶ 지수적으로 많은 선형적인 영역 조각들
- 일반화 성능 향상과 과잉적합 해소



## 오류 역전파 알고리즘

### 목적함수의 정의

- 훈련집합
  - 특징 벡터 집합  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ 과 부류 벡터 집합  $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$
  - 부류 벡터는 원핫 코드로 표현됨, 즉  $\mathbf{y}_i = (0, 0, \dots, 1, \dots, 0)^T$
  - 설계 행렬로 쓰면

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} \mathbf{y}_1^T \\ \mathbf{y}_2^T \\ \vdots \\ \mathbf{y}_n^T \end{pmatrix} \quad (3.16)$$

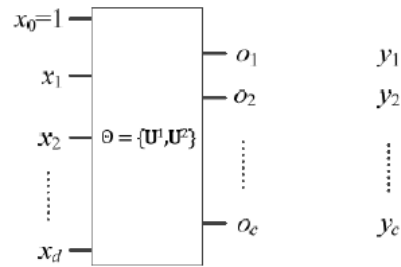
- 기계 학습의 목표
  - 모든 샘플을 옳게 분류하는 함수  $f$ 를 찾는 일

$$\left. \begin{array}{l} \mathbf{Y} = \mathbf{f}(\mathbf{X}) \\ \text{풀어 쓰면 } \mathbf{y}_i = \mathbf{f}(\mathbf{x}_i), i = 1, 2, \dots, n \end{array} \right\} \quad (3.17)$$

- 목적함수
  - 평균 제곱 오차로 정의

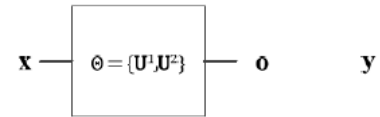
$$\left. \begin{array}{l} \text{온라인 모드: } e = \frac{1}{2} \|\mathbf{y} - \mathbf{o}\|_2^2 \\ \text{배치 모드: } e = \frac{1}{2n} \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{o}_i\|_2^2 \end{array} \right\} \quad (3.19)$$

입력 (특징 벡터  $\mathbf{x}$ )      출력 (출력 벡터  $\mathbf{o}$ )      기댓값 (부류 벡터  $\mathbf{y}$ )



(a) 블록 다이어그램

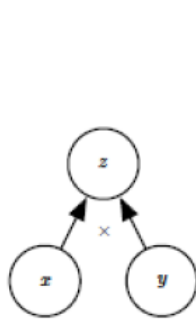
특징 벡터      출력 벡터      부류 벡터



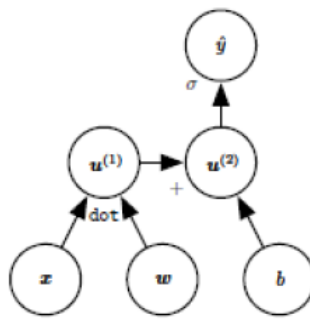
(b) 축약형 블록 다이어그램

그림 3-17 목적함수 정의에 사용하는 입력과 출력, 기댓값

- 연산 그래프의 예 : 연산을 그래프로 표현



(a)



(b)

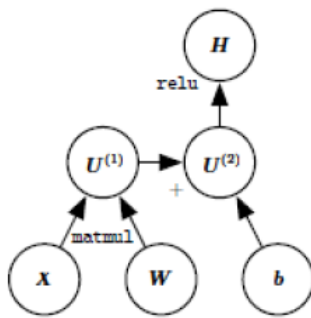
(a)  $z = xy$

(b) logistic regression

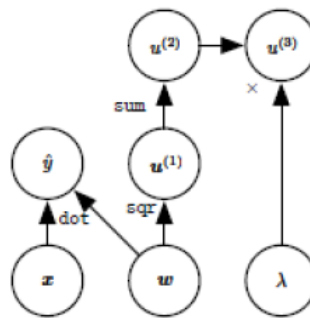
$$\hat{y} = \sigma(x^T w + b)$$

(c) ReLU activation

$$H = \max\{0, XW + b\}$$



(c)



(d)

(d) linear regression

▶  $w$  : weights

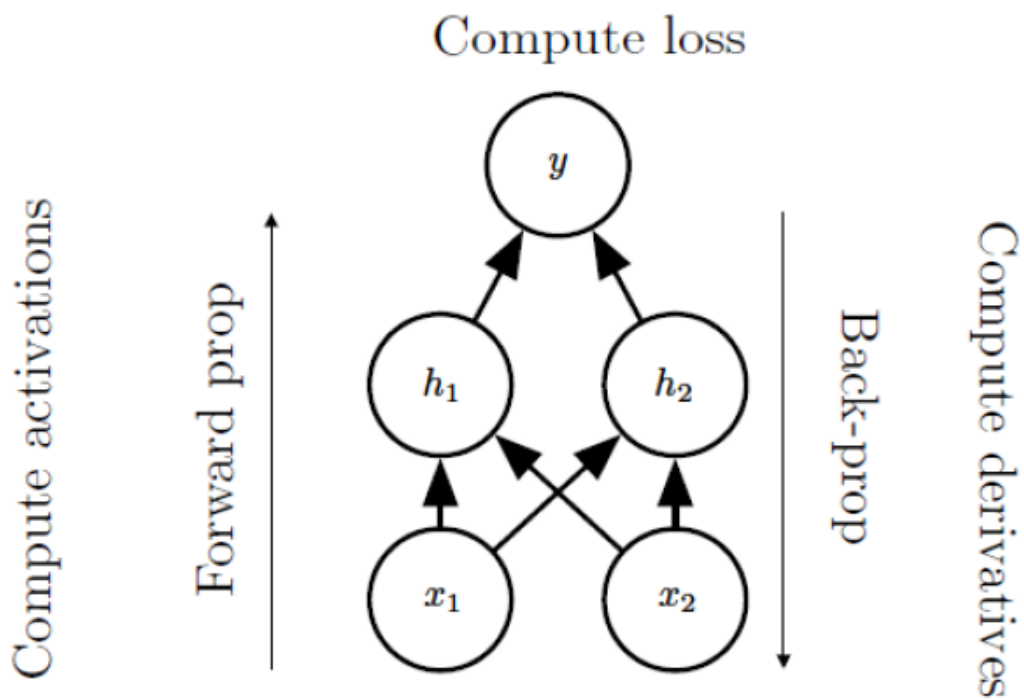
▶  $\hat{y}$  : prediction

▶  $\lambda \sum_i w_i^2$  : weight decay penalty

## 오류 역전파 알고리즘 설계

- 간단한 오류 역전파의 연산 그래프 예





- 오류 역전파 미분의 연쇄 법칙을 이용
  - 연쇄 법칙

- 수인 경우 :  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ .

- 벡터인 경우 :  $\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^{\top} \nabla_{\mathbf{y}} z$ .

$\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n$  라고 가정하면

$$g : \mathbb{R}^m \mapsto \mathbb{R}^n$$

$$f : \mathbb{R}^n \mapsto \mathbb{R}$$

$$\rightarrow \mathbf{y} = g(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^m \mapsto \mathbb{R}^n \ni \mathbf{y}$$

$$\begin{aligned} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \in \mathbb{R}^{n \times m} \end{aligned}$$

$$\begin{aligned}\nabla_y z &= \begin{bmatrix} \frac{\partial z}{\partial y_1} & \cdots & \frac{\partial z}{\partial y_n} \end{bmatrix}^T \in \mathbb{R}^{n \times 1} \\ \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z &= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \begin{bmatrix} \frac{\partial z}{\partial y_1} \\ \vdots \\ \frac{\partial z}{\partial y_n} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x_1} + \cdots + \frac{\partial z}{\partial y_n} \frac{\partial y_n}{\partial x_1} \\ \vdots \\ \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x_m} + \cdots + \frac{\partial z}{\partial y_n} \frac{\partial y_n}{\partial x_m} \end{bmatrix} = \begin{bmatrix} \frac{\partial z}{\partial x_1} \\ \vdots \\ \frac{\partial z}{\partial x_m} \end{bmatrix} = \nabla_x z\end{aligned}$$

- 야코비안 행렬과 그레디언트를 곱한 연쇄 법칙을 얻어서 구해짐

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z \quad \underbrace{x \in \mathbb{R}^m \xrightarrow{g} y \in \mathbb{R}^n \xrightarrow{f} \mathbb{R} \ni z}_{\text{Jacobian: } \frac{\partial y}{\partial x} \quad \text{gradient: } \nabla_y z}$$

- 연쇄법칙의 구현

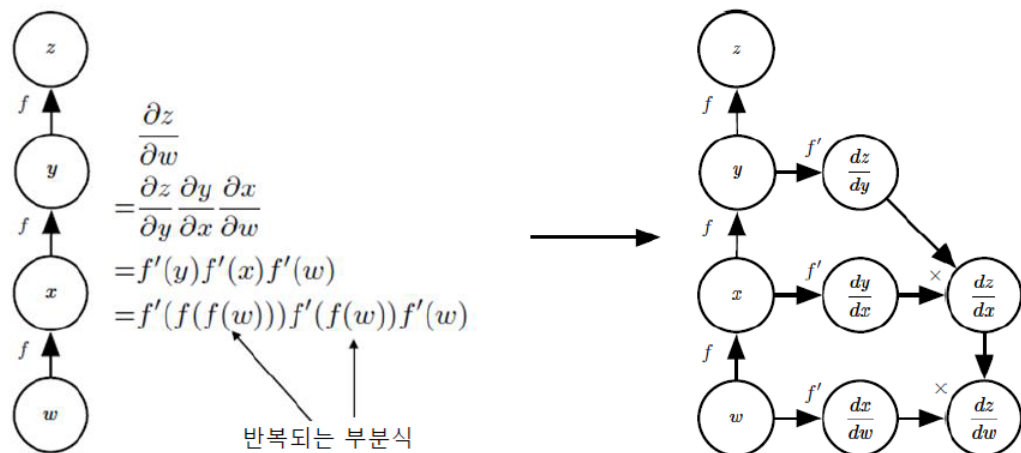
- 반복되는 부분식들을 저장하거나 재연산을 최소화

- ex) 동적 프로그래밍
- 연산속도와 저장 공간의 Trade-off

$$(f \circ g \circ h \circ i)'(x) = (f' \circ g \circ h \circ i)(x) \times (g' \circ h \circ i)(x) \times (h' \circ i)(x) \times i'(x)$$



- 그레디언트 계산을 위한 연산 그래프 예



- 식 (3.19)의 목적함수를 다시 쓰면

- 2층 퍼셉트론의 경우  $\Theta = \{U^1, U^2\}$

$$J(\Theta) = \frac{1}{2} \|\mathbf{y} - \mathbf{o}(\Theta)\|_2^2 \quad (3.20)$$

- 경사하강법

$J(\Theta) = J(\{\mathbf{U}^1, \mathbf{U}^2\})$ 의 최저점을 찾아주는 경사 하강법

$$\left. \begin{aligned} \mathbf{U}^1 &= \mathbf{U}^1 - \rho \frac{\partial J}{\partial \mathbf{U}^1} \\ \mathbf{U}^2 &= \mathbf{U}^2 - \rho \frac{\partial J}{\partial \mathbf{U}^2} \end{aligned} \right\} \quad (3.21)$$

- 식 (3.21)을 알고리즘 형태로 쓰면

### 알고리즘 3-3 다층 퍼셉트론을 위한 스토케스틱 경사 하강법

입력: 훈련집합  $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ , 학습률  $\rho$

출력: 가중치 행렬  $\mathbf{U}^1$ 과  $\mathbf{U}^2$

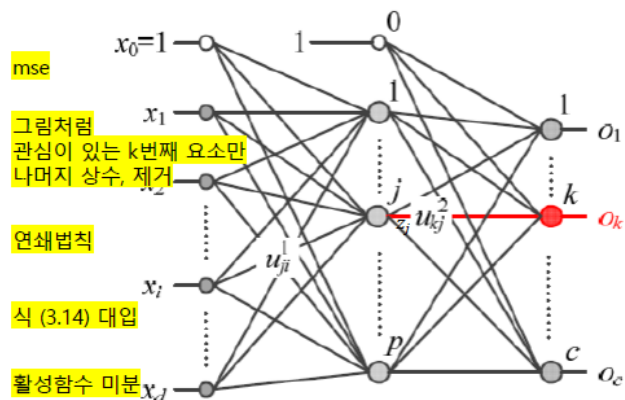
- 1  $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 를 초기화한다.
- 2 repeat
- 3    $\mathbb{X}$ 의 순서를 섞는다.
- 4   for ( $\mathbb{X}$ 의 샘플 각각에 대해)
- 5       식 (3.15)로 전방 계산을 하여  $\mathbf{o}$ 를 구한다.
- 6        $\frac{\partial J}{\partial \mathbf{U}^1}$ 와  $\frac{\partial J}{\partial \mathbf{U}^2}$ 를 계산한다.
- 7       식 (3.21)로  $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 를 갱신한다.
- 8 until (멈춤 조건)

- 오류 역전파의 유도

- 알고리즘 3-3의 라인 6을 위한 도함수 값  $\frac{\partial J}{\partial \mathbf{U}^1}$ 와  $\frac{\partial J}{\partial \mathbf{U}^2}$ 의 계산 과정

- 먼저  $\mathbf{U}^2$ 를 구성하는  $u_{kj}^2$ 로 미분하면,

$$\begin{aligned} \frac{\partial J}{\partial u_{kj}^2} &= \frac{\partial (0.5 \|\mathbf{y} - \mathbf{o}(\mathbf{U}^1, \mathbf{U}^2)\|_2^2)}{\partial u_{kj}^2} \\ &= \frac{\partial (0.5 \sum_{q=1}^c (y_q - o_q)^2)}{\partial u_{kj}^2} \\ &= \frac{\partial (0.5 (y_k - o_k)^2)}{\partial u_{kj}^2} \\ &= -(y_k - o_k) \frac{\partial o_k}{\partial u_{kj}^2} \\ &= -(y_k - o_k) \frac{\partial \tau(\text{osum}_k)}{\partial u_{kj}^2} \\ &= -(y_k - o_k) \tau'(\text{osum}_k) \frac{\partial \text{osum}_k}{\partial u_{kj}^2} \\ &= -(y_k - o_k) \tau'(\text{osum}_k) z_j \end{aligned}$$



(a)  $u_{kj}^2$ 가 미치는 영향

osum 미분 그림 3-18 매개변수가 미치는 영향

$$o_k = \tau(\text{osum}_k), k = 1, 2, \dots, c$$

$$\text{이때 } \text{osum}_k = \mathbf{u}_k^2 \mathbf{z} \text{이고 } \mathbf{u}_k^2 = (u_{k0}^2, u_{k1}^2, \dots, u_{kp}^2), \mathbf{z} = (1, z_1, z_2, \dots, z_p)^T$$

(3.14)

- $\mathbf{U}^1$ 을 구성하는  $u_{ji}^1$ 로 미분하면

$$o_k = \tau(osum_k), k = 1, 2, \dots, c \quad (3.14)$$

$$\text{이때 } osum_k = \mathbf{u}_k^2 \odot \mathbf{1} \text{고 } \mathbf{u}_k^2 = (u_{k0}^2, u_{k1}^2, \dots, u_{kp}^2), \mathbf{z} = (1, z_1, z_2, \dots, z_p)^T$$

$$z_j = \tau(zsum_j), j = 1, 2, \dots, p \quad (3.13)$$

$$\text{이때 } zsum_j = \mathbf{u}_j^1 \mathbf{x} \text{고 } \mathbf{u}_j^1 = (u_{j0}^1, u_{j1}^1, \dots, u_{jd}^1), \mathbf{x} = (1, x_1, x_2, \dots, x_d)^T$$

$$\frac{\partial J}{\partial u_{ji}^1} = \frac{\partial (0.5 \|\mathbf{y} - \mathbf{o}(\mathbf{U}^1, \mathbf{U}^2)\|_2^2)}{\partial u_{ji}^1}$$

$$= \frac{\partial \left( 0.5 \sum_{q=1}^c (y_q - o_q)^2 \right)}{\partial u_{ji}^1}$$

$$= - \sum_{q=1}^c (y_q - o_q) \frac{\partial o_q}{\partial u_{ji}^1}$$

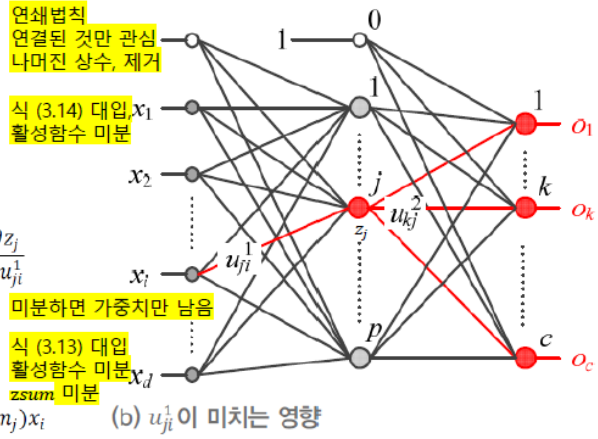
$$= - \sum_{q=1}^c (y_q - o_q) \tau'(osum_q) \frac{\partial osum_q}{\partial u_{ji}^1}$$

$$= - \sum_{q=1}^c (y_q - o_q) \tau'(osum_q) \frac{\partial osum_q}{\partial z_j} \frac{\partial z_j}{\partial u_{ji}^1}$$

$$= - \sum_{q=1}^c (y_q - o_q) \tau'(osum_q) u_{qj}^2 \frac{\partial z_j}{\partial u_{ji}^1}$$

$$= - \sum_{q=1}^c (y_q - o_q) \tau'(osum_q) u_{qj}^2 \tau'(zsum_j) x_i$$

$$= - \tau'(zsum_j) x_i \sum_{q=1}^c (y_q - o_q) \tau'(osum_q) u_{qj}^2$$



○ 지금까지 유도한 식을 정리하면

$$\delta_k = (y_k - o_k) \tau'(osum_k), \quad 1 \leq k \leq c \quad (3.22)$$

반복되는 부분식

$$\frac{\partial J}{\partial u_{kj}^2} = \Delta u_{kj}^2 = -\delta_k z_j, \quad 0 \leq j \leq p, 1 \leq k \leq c \quad (3.23)$$

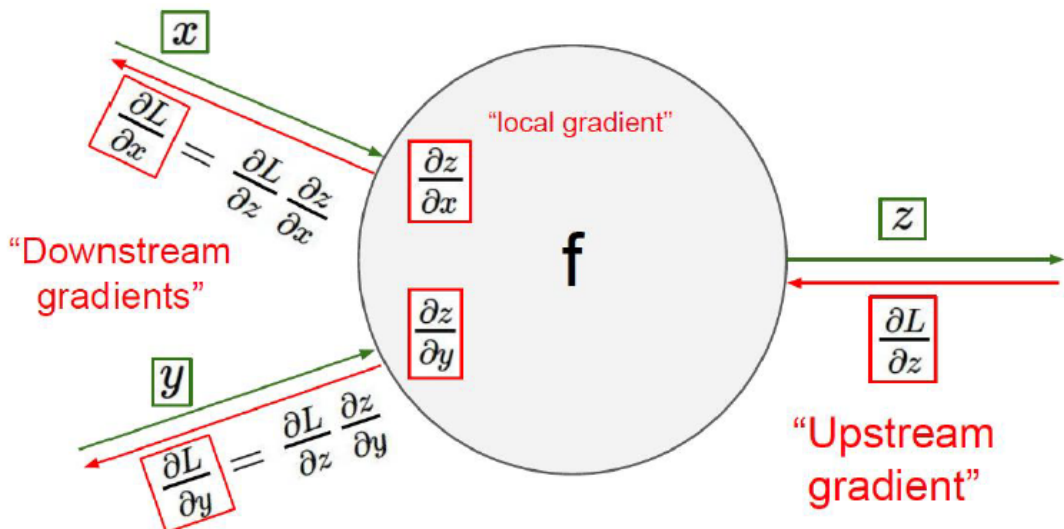
$$\eta_j = \tau'(zsum_j) \sum_{q=1}^c \delta_q u_{qj}^2, \quad 1 \leq j \leq p \quad (3.24)$$

$$\frac{\partial J}{\partial u_{ji}^1} = \Delta u_{ji}^1 = -\eta_j x_i, \quad 0 \leq i \leq d, 1 \leq j \leq p \quad (3.25)$$

○ 오류 역전파 알고리즘

- 식 (3.22) ~ (3.25)를 이용하여 출력층의 오류를 역방향(왼쪽)으로 전파하며 그레이디언트를 계산하는 알고리즘

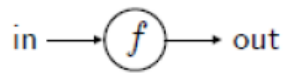
● 역전파 분해



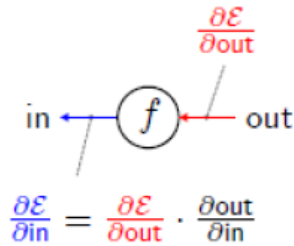
- 단일 노드의 역전파 예

$$\text{out} = f(\text{in})$$

forward



backprop

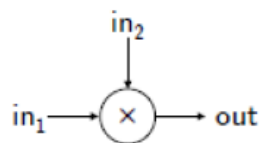


$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}} &= \underbrace{\frac{\partial \mathcal{E}}{\partial \text{out}}}_{\text{output gradient}} \cdot \underbrace{\frac{\partial \text{out}}{\partial \text{in}}}_{\text{local gradient}} \\ &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot f'(\text{in}) \end{aligned}$$

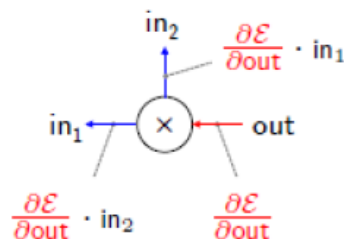
- 곱셈의 역전파 예

$$\text{out} = \text{in}_1 \cdot \text{in}_2$$

forward



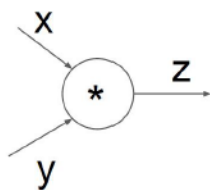
backprop



$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}_1} &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial \text{in}_1} \\ &= \underbrace{\frac{\partial \mathcal{E}}{\partial \text{out}}}_{\text{output gradient}} \cdot \underbrace{\text{in}_2}_{\text{local gradient}} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}_2} &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial \text{in}_2} \\ &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \text{in}_1 \end{aligned}$$

- 곱셈의 역전파 PyTorch 구현 예



(x,y,z are scalars)

```
class Multiply(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x, y):
        ctx.save_for_backward(x, y)
        z = x * y
        return z
    @staticmethod
    def backward(ctx, grad_z):
        x, y = ctx.saved_tensors
        grad_x = y * grad_z
        grad_y = x * grad_z
        return grad_x, grad_y
```

Need to stash some values for use in backward

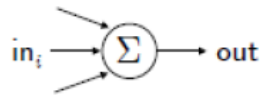
Upstream gradient

Multiply upstream and local gradients

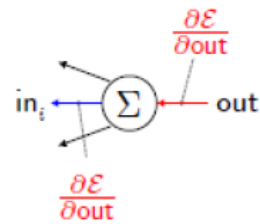
- 덧셈의 역전파 예

$$\text{out} = \sum_i \text{in}_i$$

forward



backprop



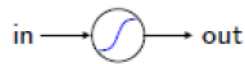
$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}_i} &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial \text{in}_i} \\ &= \underbrace{\frac{\partial \mathcal{E}}{\partial \text{out}}}_{\text{output gradient}} \cdot \underbrace{1}_{\text{local gradient}} \\ &= \frac{\partial \mathcal{E}}{\partial \text{out}} \\ &\triangleq \delta \end{aligned}$$

sum (forward)  $\Leftrightarrow$  fanout (backprop)

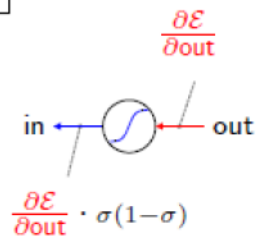
- 시그모이드의 역전파 예

$$\text{out} = \sigma(\text{in})$$

forward



backprop



$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}} &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial \text{in}} \\ &= \underbrace{\frac{\partial \mathcal{E}}{\partial \text{out}}}_{\text{output gradient}} \cdot \underbrace{\sigma'(\text{in})}_{\text{local gradient}} \\ &= \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot [\sigma(\text{in}) (1 - \sigma(\text{in}))] \end{aligned}$$

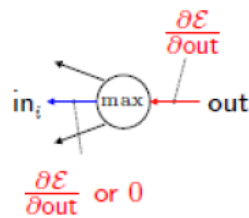
- 최대화의 역전파 예

$$\text{out} = \max_i \{ \text{in}_i \}$$

forward



backprop



$$\frac{\partial \mathcal{E}}{\partial \text{in}_i} = \frac{\partial \mathcal{E}}{\partial \text{out}} \cdot \underbrace{\frac{\partial \text{out}}{\partial \text{in}}}_{1 \text{ or } 0}$$

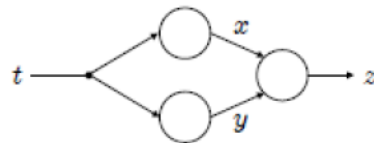
$$= \begin{cases} \frac{\partial \mathcal{E}}{\partial \text{out}} & \text{if } \text{in}_i \text{ is max} \\ 0 & \text{otherwise} \end{cases}$$

max (forward)  $\Leftrightarrow$  mux (backprop)

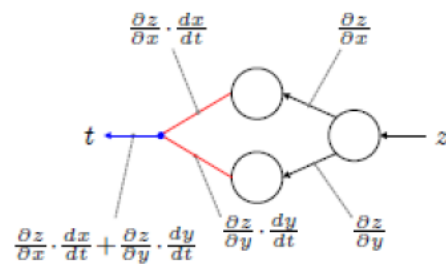
- 전개의 역전파 예

multivariable chain rule

forward



backprop



let

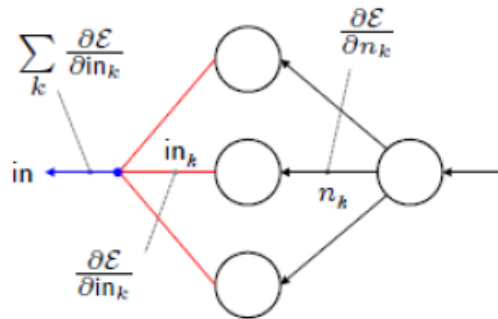
$$x = x(t), \quad y = y(t)$$

$$z = f(x, y)$$

then

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial t} + \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial t}$$

fanout



assuming

$$\mathcal{E} = f(n_1, \dots, n_k, \dots)$$

and

$$n_k = n_k(\text{in})$$

gives

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \text{in}} &= \sum_k \frac{\partial \mathcal{E}}{\partial n_k} \cdot \frac{\partial n_k}{\partial \text{in}} \\ &= \sum_k \frac{\partial \mathcal{E}}{\partial \text{in}_k} \end{aligned}$$

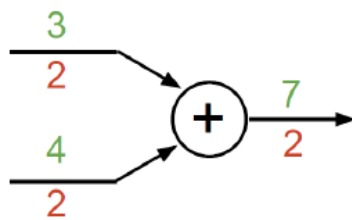
where

$$\text{in}_k \triangleq \text{input to } n_k$$

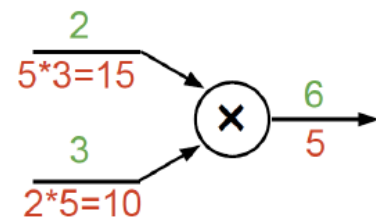
fanout (forward)  $\Leftrightarrow$  sum (backprop)

- 역전파 주요 예

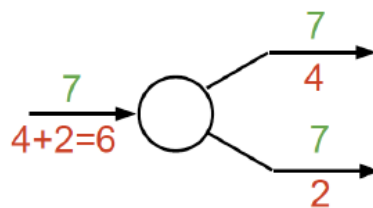
**add** gate: gradient distributor



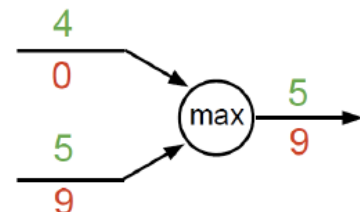
**mul** gate: “swap multiplier”



**copy** gate: gradient adder



**max** gate: gradient router

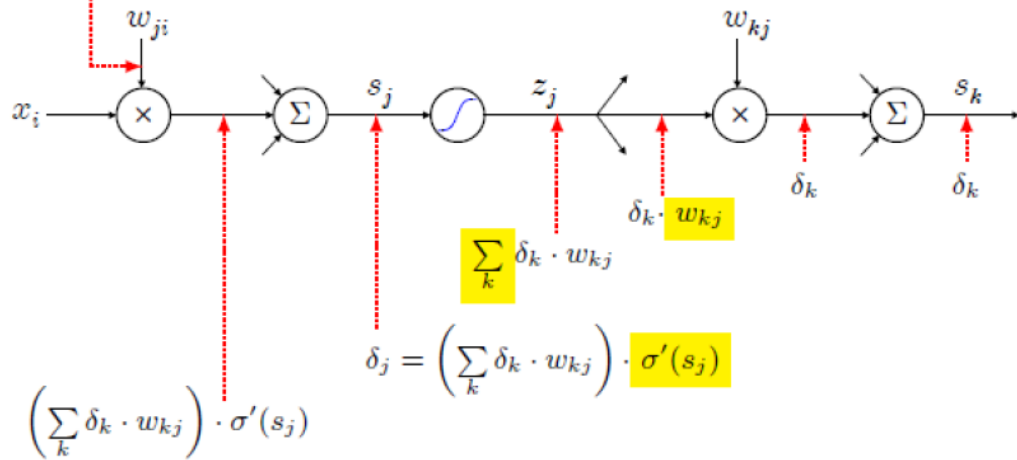


- 실제 역전파 예

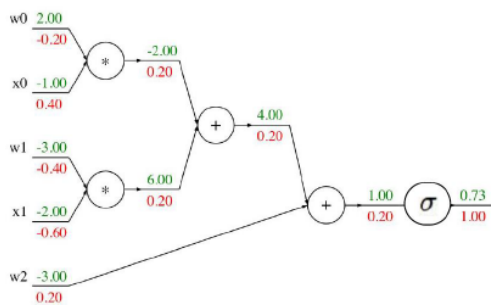


computing  $\frac{\partial \mathcal{E}}{\partial w_{ji}}$

$$\left( \sum_k \delta_k \cdot w_{kj} \right) \cdot \sigma'(s_j) \cdot x_i \Rightarrow \frac{\partial \mathcal{E}}{\partial w_{ji}} = \delta_j \cdot x_i = \left( \sum_k \delta_k \cdot w_{kj} \right) \cdot \sigma'(s_j) \cdot x_i$$



- 역전파의 간단한 구현



Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

Backward pass:  
Compute grads

```
    grad_L = 1.0
    grad_s3 = grad_L * (1 - L) * L
    grad_w2 = grad_s3
    grad_s2 = grad_s3
    grad_s0 = grad_s2
    grad_s1 = grad_s2
    grad_w1 = grad_s1 * x1
    grad_x1 = grad_s1 * w1
    grad_w0 = grad_s0 * x0
    grad_x0 = grad_s0 * w0
```

## 오류 역전파를 이용한 학습 알고리즘

- 식 (3.22) ~ (3.25)를 이용한 스토캐스틱 경사 하강법

### 알고리즘 3-4 다층 퍼셉트론 학습을 위한 스토캐스틱 경사 하강법

입력: 훈련집합  $\mathbb{X}$ 와  $\mathbb{Y}$ , 학습률  $\rho$

출력: 가중치 행렬  $\mathbf{U}^1$ 과  $\mathbf{U}^2$

```

1   $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 를 초기화한다.
2  repeat
3     $\mathbb{X}$ 의 순서를 섞는다.
4    for ( $\mathbb{X}$ 의 샘플 각각에 대해)
5      현재 처리하는 샘플을  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_c)^T$ 라 표기한다.
6       $x_0$ 과  $z_0$ 을 1로 설정한다. // 바이어스
          // 전방 계산
7      for ( $j=1$  to  $\rho$ )  $zsum_j = \mathbf{u}_j^1 \mathbf{x}$ ,  $z_j = \tau(zsum_j)$  // 식 (3.13)
8      for ( $k=1$  to  $c$ )  $osum_k = \mathbf{u}_k^2 \mathbf{z}$ ,  $o_k = \tau(osum_k)$  // 식 (3.14)
          // 오류 역전파
9      for ( $k=1$  to  $c$ )  $\delta_k = (y_k - o_k)\tau'(osum_k)$  // 식 (3.22)
10     for ( $k=1$  to  $c$ ) for ( $j=0$  to  $\rho$ )  $\Delta u_{kj}^2 = -\delta_k z_j$  // 식 (3.23)
11     for ( $j=1$  to  $\rho$ )  $\eta_j = \tau'(zsum_j) \sum_{q=1}^c \delta_q u_{qj}^2$  // 식 (3.24)
12     for ( $j=1$  to  $\rho$ ) for ( $i=0$  to  $d$ )  $\Delta u_{ji}^1 = -\eta_j x_i$  // 식 (3.25)
          // 가중치 갱신
13     for ( $k=1$  to  $c$ ) for ( $j=0$  to  $\rho$ )  $u_{kj}^2 = u_{kj}^2 - \rho \Delta u_{kj}^2$  // 식 (3.21)
14     for ( $j=1$  to  $\rho$ ) for ( $i=0$  to  $d$ )  $u_{ji}^1 = u_{ji}^1 - \rho \Delta u_{ji}^1$  // 식 (3.21)
15 until (멈춤 조건)

```

- 임의 샘플링 방식으로 바꾸려면

$$\left\{ \begin{array}{l} 3. \mathbb{X} \text{의 순서를 섞는다.} \\ 4. \text{for } (\mathbb{X} \text{의 샘플 각각에 대해}) \\ 5. \quad \text{현재 처리하는 샘플을 } \mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T, \mathbf{y} = (y_1, y_2, \dots, y_c)^T \text{라 표기한다.} \\ 6. \quad x_0 \text{과 } z_0 \text{을 1로 설정한다.} \end{array} \right.$$

$$\rightarrow \left\{ \begin{array}{l} 3. \mathbb{X} \text{에서 임의로 샘플 하나를 뽑는다.} \\ 4. \quad \text{뽑힌 샘플을 } \mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T, \mathbf{y} = (y_1, y_2, \dots, y_c)^T \text{라 표기한다.} \\ 5. \quad x_0 \text{과 } z_0 \text{을 1로 설정한다.} \end{array} \right.$$

- 도함수의 종류

#### Scalar to Scalar

$$x \in \mathbb{R}, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If  $x$  changes by a small amount, how much will  $y$  change?

#### Vector to Scalar

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left( \frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

For each element of  $x$ , if it changes by a small amount then how much will  $y$  change?

#### Vector to Vector

$$x \in \mathbb{R}^N, y \in \mathbb{R}^M$$

Derivative is **Jacobian**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left( \frac{\partial y}{\partial x} \right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

For each element of  $x$ , if it changes by a small amount then how much will each element of  $y$  change?

- 행렬 표기 : GPU를 사용한 고속 행렬 연산에 적합

#### 알고리즘 3-5 다층 퍼셉트론 학습을 위한 스토캐스틱 경사 하강법(행렬 표기)

입력: 훈련집합  $\mathbb{X}$ 와  $\mathbb{Y}$ , 학습률  $\rho$

출력: 가중치 행렬  $\mathbf{U}^1$ 과  $\mathbf{U}^2$

```

1   $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 를 초기화한다.
2  repeat
3     $\mathbb{X}$ 의 순서를 섞는다.
4    for ( $\mathbb{X}$ 의 샘플 각각에 대해)
5      현재 처리하는 샘플을  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_c)^T$ 라 표기한다.
6       $x_0$ 과  $z_0$ 을 1로 설정한다. // 바이어스
          // 전방 계산
7       $\mathbf{zsum} = \mathbf{U}^1 \mathbf{x}$ ,  $\tilde{\mathbf{z}} = \boldsymbol{\tau}(\mathbf{zsum})$  // 식 (3.13),  $\mathbf{zsum}_{p \times 1}$ ,  $\mathbf{U}^1_{p \times (d+1)}$ ,  $\mathbf{x}_{(d+1) \times 1}$ ,  $\tilde{\mathbf{z}}_{p \times 1}$ 
8       $\mathbf{osum} = \mathbf{U}^2 \tilde{\mathbf{z}}$ ,  $\mathbf{o} = \boldsymbol{\tau}(\mathbf{osum})$  // 식 (3.14),  $\mathbf{osum}_{c \times 1}$ ,  $\mathbf{U}^2_{c \times (p+1)}$ ,  $\tilde{\mathbf{z}}_{(p+1) \times 1}$ ,  $\mathbf{o}_{c \times 1}$ 
          // 오류 역전파
9       $\boldsymbol{\delta} = (\mathbf{y} - \mathbf{o}) \odot \boldsymbol{\tau}'(\mathbf{osum})$  // 식 (3.22),  $\boldsymbol{\delta}_{c \times 1}$ 
10      $\Delta \mathbf{U}^2 = -\boldsymbol{\delta} \mathbf{z}^T$  // 식 (3.23),  $\Delta \mathbf{U}^2_{c \times (p+1)}$ 
11      $\boldsymbol{\eta} = (\boldsymbol{\delta}^T \tilde{\mathbf{U}}^2)^T \odot \boldsymbol{\tau}'(\mathbf{zsum})$  // 식 (3.24),  $\tilde{\mathbf{U}}^2_{c \times p}$ ,  $\boldsymbol{\eta}_{p \times 1}$ 
12      $\Delta \mathbf{U}^1 = -\boldsymbol{\eta} \mathbf{x}^T$  // 식 (3.25),  $\Delta \mathbf{U}^1_{p \times (d+1)}$ 
          // 가중치 갱신
13      $\mathbf{U}^2 = \mathbf{U}^2 - \rho \Delta \mathbf{U}^2$  // 식 (3.21)
14      $\mathbf{U}^1 = \mathbf{U}^1 - \rho \Delta \mathbf{U}^1$  // 식 (3.21)
15 until (멈춤 조건)
```

## 미니배치 스토캐스틱 경사 하강법

- 미니배치 방식
  - 한번에  $t$ 개의 샘플을 처리함 ( $t$ 는 미니배치 크기)
    - $t = 1$ 이면 스토캐스틱 경사 하강법
    - $t = n$ 이면 배치 경사 하강법
  - 미니배치 방식은 보통  $t =$  수십 ~ 수백
    - 그레이디언트의 잡음을 줄여주는 효과 때문에 수렴이 빨라짐
    - GPU를 사용한 병렬처리에 유리함
  - 현대 기계 학습은 미니배치를 표준처럼 여겨 널리 사용함

### 알고리즘 3-6 다층 퍼셉트론 학습을 위한 '미니배치' 스토캐스틱 경사 하강법

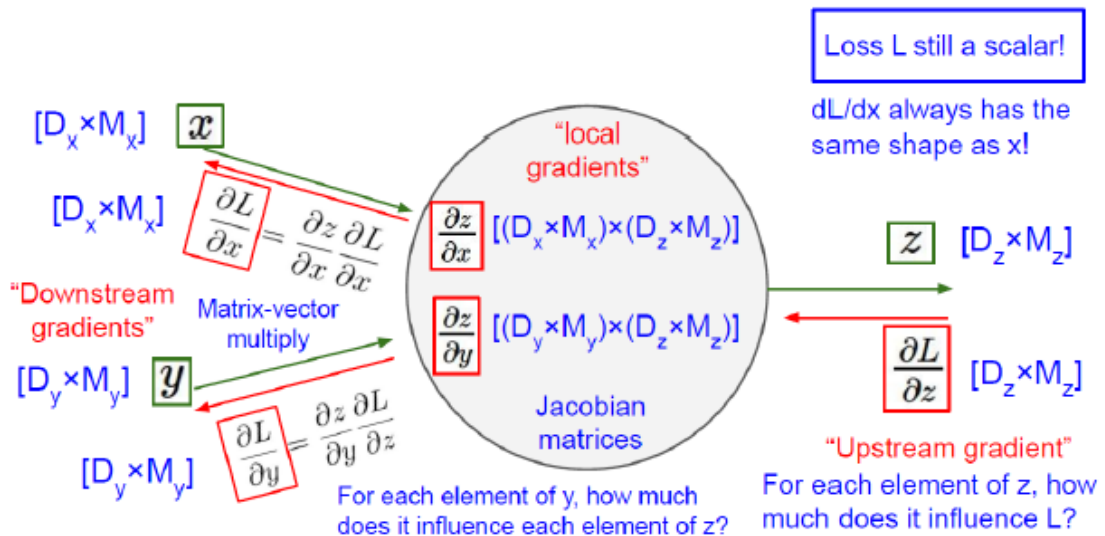
입력: 훈련집합  $\mathbb{X}$ 와  $\mathbb{Y}$ , 학습률  $\rho$ , 미니배치 크기  $t$

출력: 가중치 행렬  $\mathbf{U}^1$ 과  $\mathbf{U}^2$

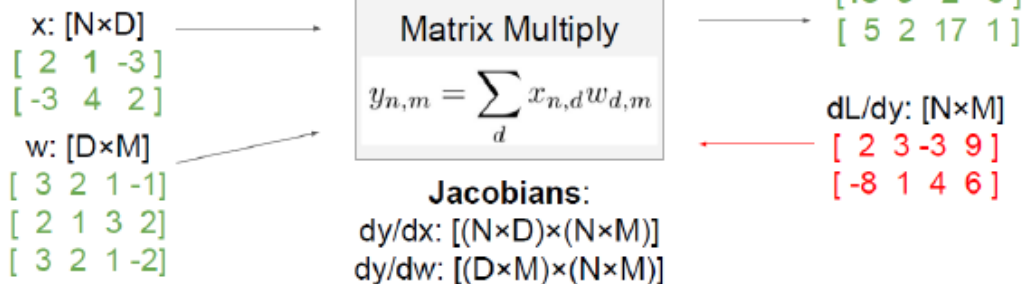
```

1   $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 를 초기화한다.
2  repeat
3       $\mathbb{X}$ 와  $\mathbb{Y}$ 에서  $t$ 개의 샘플을 무작위로 뽑아 미니배치  $\mathbb{X}'$ 와  $\mathbb{Y}'$ 를 만든다.
4       $\Delta \mathbf{U}^2 = \mathbf{0}$ ,  $\Delta \mathbf{U}^1 = \mathbf{0}$ 
5      for ( $\mathbb{X}'$ 의 샘플 각각에 대해)
6          현재 처리하는 샘플을  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_c)^T$ 라 표기한다.
7           $x_0$ 와  $z_0$ 를 1로 설정한다. // 바이어스
              // 전방 계산
8           $\mathbf{zsum} = \mathbf{U}^1 \mathbf{x}$ ,  $\tilde{\mathbf{z}} = \tau(\mathbf{zsum})$  // 식 (3.13),  $\mathbf{zsum}_{p+1}$ ,  $\mathbf{U}^1_{p \times (d+1)}$ ,  $\mathbf{x}_{(d+1) \times 1}$ ,  $\tilde{\mathbf{z}}_{p \times 1}$ 
9           $\mathbf{osum} = \mathbf{U}^2 \tilde{\mathbf{z}}$ ,  $\mathbf{o} = \tau(\mathbf{osum})$  // 식 (3.14),  $\mathbf{osum}_{c+1}$ ,  $\mathbf{U}^2_{c \times (p+1)}$ ,  $\tilde{\mathbf{z}}_{(p+1) \times 1}$ ,  $\mathbf{o}_{c \times 1}$ 
              // 오류 역전파
10          $\delta = (\mathbf{y} - \mathbf{o}) \odot \tau'(\mathbf{osum})$  // 식 (3.22),  $\delta_{c+1}$ 
11          $\Delta \mathbf{U}^2 = \Delta \mathbf{U}^2 + (-\delta \tilde{\mathbf{z}}^T)$  // 식 (3.23)을 누적,  $\Delta \mathbf{U}^2_{c \times (p+1)}$ 
12          $\eta = (\delta^T \tilde{\mathbf{U}}^2)^T \odot \tau'(\mathbf{zsum})$  // 식 (3.24),  $\tilde{\mathbf{U}}^2_{c \times p}$ ,  $\eta_{p \times 1}$ 
13          $\Delta \mathbf{U}^1 = \Delta \mathbf{U}^1 + (-\eta \mathbf{x}^T)$  // 식 (3.25)를 누적,  $\Delta \mathbf{U}^1_{p \times (d+1)}$ 
              // 가중치 갱신
14          $\mathbf{U}^2 = \mathbf{U}^2 - \rho \left( \frac{1}{t} \right) \Delta \mathbf{U}^2$  // 식 (3.21) - 평균 그래디언트로 갱신
15          $\mathbf{U}^1 = \mathbf{U}^1 - \rho \left( \frac{1}{t} \right) \Delta \mathbf{U}^1$  // 식 (3.21) - 평균 그래디언트로 갱신
16 until (멈춤 조건)
    
```

- 행렬의 역전파를 위한 도함수



#### Backprop with Matrices



For a neural net we may have

$N=64$ ,  $D=M=4096$

Each Jacobian takes 256 GB of memory!

## 다층 퍼셉트론에 의한 인식

- 예측 단계 (또는 테스트 단계)
  - 학습을 마친 후 현장 설치하여 사용(또는 테스트 집합으로 성능 테스트)

### 알고리즘 3-7 다층 퍼셉트론을 이용한 인식

입력: 테스트 샘플  $\mathbf{x}$  // 신경망의 가중치  $\mathbf{U}^1$ 과  $\mathbf{U}^2$ 는 이미 설정되었다고 가정함.

출력: 부류  $y$

- $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T$ 로 확장하고,  $x_0$ 과  $z_0$ 을 1로 설정한다.
- $\mathbf{zsum} = \mathbf{U}^1 \mathbf{x}$
- $\hat{\mathbf{z}} = \tau(\mathbf{zsum})$  // 식 (3.13)
- $\mathbf{osum} = \mathbf{U}^2 \hat{\mathbf{z}}$
- $\mathbf{o} = \tau(\mathbf{osum})$  // 식 (3.14)
- $\mathbf{o}$ 에서 가장 큰 값을 가지는 노드에 해당하는 부류 번호를  $y$ 에 대입한다.

- 라인 6을 수식으로 표현하면 
$$y = \underset{k}{\operatorname{argmax}} o_k$$
- 전방 계산 한번만 사용하므로 빠름

## 다층 퍼셉트론의 특성

### 오류 역전파 알고리즘의 빠른 속도

- 연산 횟수 비교

표 3-2 전방 계산과 오류 역전파 과정이 사용하는 연산 횟수

과정	수식	덧셈	곱셈
전방 계산	식 (3.13)	$dp$	$dp$
	식 (3.14)	$pc$	$pc$
오류 역전파	식 (3.22)	$c$	$c$
	식 (3.23)		$cp$
	식 (3.24)	$cp$	$cp$
	식 (3.25)		$dp$
	식 (3.21)	$cp+dp$	$cp+dp$

- 오류 역전파는 전방 계산보다 약 1.5~2배의 시간 소요 ▶ 빠른 계산 가능
- 하지만 학습 알고리즘은 수렴할 때 까지 오류 역전파를 반복해야 하므로 점근적 시간복잡도는

$$\theta((dp + pc)nq) \quad (n \text{은 훈련집합의 크기, } q \text{는 에포크 epoch 수})$$

- 에포크는 전체 학습 집합을 수행한 단위

### 모든 함수를 정확하게 근사할 수 있는 능력

- 호닉의 주장[Hornik1989]
  - 은닉층을 하나만 가진 다층 퍼셉트론은 범용근사자

- 은닉 노드가 충분히 많다면 활성화함수로 무엇을 사용하든 표준 다층 퍼셉트론은 어떤 함수라도 원하는 정확도만큼 근사화할 수 있다
- 은닉 노드를 무수히 많게 할 수 없으므로, 실질적으로는 복잡한 구조의 데이터에서는 성능 한계

## 성능 향상을 위한 휴리스틱의 중요성

- 순수한 최적화 알고리즘으로 높은 성능 불가능
  - 데이터 희소성, 잡음, 미숙한 신경망 구조 등의 이유
  - 성능 향상을 위한 갖가지 휴리스틱을 개발하고 공유함
- 휴리스틱 개발에서 중요 쟁점
  - 아키텍처
    - 은닉층과 은닉노드의 개수를 정해야 한다. 은닉층과 은닉 노드를 늘리면 신경망의 용량은 커지는 대신, 추정할 매개변수가 많아지고 학습 과정에서 과잉적합할 가능성이 커진다.
    - 현대 기계학습은 복잡한 모델을 사용하되, 적절한 규제기법을 적용하는 경향이 있음
  - 초깃값
    - 가중치를 초기화할때 보통 난수를 생성하여 설정하는데, 값의 범위와 분포가 중요
  - 학습률
    - 처음부터 끝까지 같은 학습률을 사용하는 방식과 처음에는 큰 값으로 시작하고 점점 줄이는 적응적 방식이 있음
  - 활성화함수
    - 초창기 다층 퍼셉트론은 주로 로지스틱 시그모이드나 tanh 함수를 사용했는데, 은닉층의 개수를 늘림에 따라 그레이디언트 소멸과 같은 몇가지 문제가 발생한다
    - 깊은 신경망은 주로 ReLU 함수를 사용
- 실용적인 성능
  - 1980 ~ 1990년대에 다층 퍼셉트론은 실용 시스템 제작에 크게 기여
    - 인쇄/필기 문자 인식으로 우편물 자동 분류기, 전표 인식기, 자동차 번호판 인식기 등
    - 음성 인식, 게임, 주가 예측, 정보 검색, 의료 진단, 유전자 검색, 반도체 결함 검사 등
- 하지만 한계 노출
  - 잡음이 섞인 상황에서 음성인식 성능 저하
  - 필기 주소 인식 능력 저하
  - 바둑 등의 복잡한 문제에서의 한계
- 이러한 한계를 딥러닝은 극복함