

Spring 2023

SWCON253: Machine Learning

Lecture 03

Gradient Descent

Jinwoo Choi
Assistant Professor
CSE, Kyung Hee University



Contents

1. Vector Calculus
2. Iterative Optimization & Gradient Descent
3. Automatic Differentiation

References

- *Mathematics for Machine Learning* by Deisenroth, Faisal, and Ong (<https://mml-book.com>)
- *Intro to Deep Learning & Generative Models* by Sebastian Raschka (<http://pages.stat.wisc.edu/~sraschka/teaching/stat453-ss2020/>)
- 패턴 인식 by 오일석, 기계 학습 by 오일석



1. Vector Calculus

1. Derivative
2. Chain Rule
3. Gradient: *Collection of Partial Derivatives for a Scalar Function*
4. *Multivariate Chain Rule*
5. Jacobian: *Collection of Gradients for a Vector Function*
6. Hessian: *Collection of Second-order Partial Derivatives*
7. Gradients for a *Matrix Function*

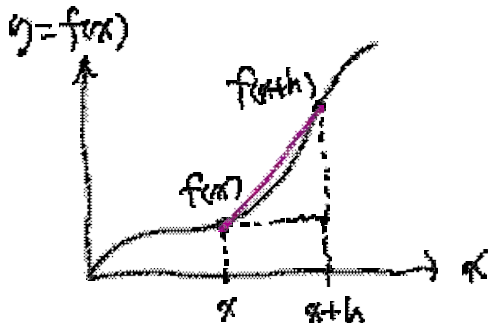


Derivative: Differentiation of univariate function

◆ Derivative of $f(x)$ 기움기

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

x a scalar



◆ Useful Formula

Product Rule: $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$ 곱미분

Quotient Rule: $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$

Sum Rule: $(f(x) + g(x))' = f'(x) + g'(x)$

★ Chain Rule: $(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x) = \frac{\partial g}{\partial y} \cdot \frac{\partial f}{\partial x}$
 함수의 함수 미분
 그냥 같은 말



or



Derivative (cont'd)

◆ Derivatives of Common Functions

	Function $f(x)$	Derivative with respect to x
1	a	0
2	x	1
3	ax	a
4	x^2	$2x$
5	x^a	ax^{a-1}
6	a^x	$\log(a)a^x$
7	$\log(x)$	$1/x$
8	$\log_a(x)$	$1/(x \log(a))$
9	$\sin(x)$	$\cos(x)$
10	$\cos(x)$	$-\sin(x)$
11	$\tan(x)$	$\sec^2(x)$



Chain Rule

◆ Composition of Functions

$$F(x) = f(g(x)).$$

$$F'(x) = f'(g(x))g'(x).$$

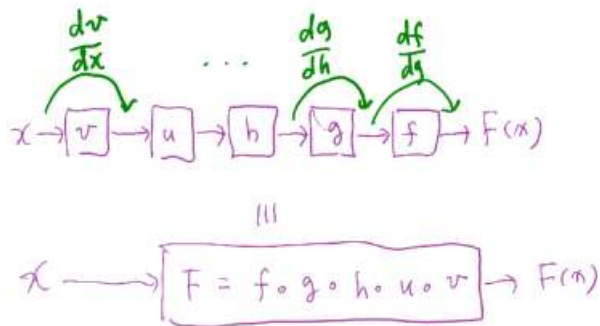
$$\frac{d}{dx}[f(g(x))] = \frac{df}{dg} \cdot \frac{dg}{dx}.$$



$$F(x) = f(g(h(u(v(x)))))$$

$$\frac{dF}{dx} = \frac{d}{dx}F(x) = \frac{d}{dx}f(g(h(u(v(x)))))$$

$$= \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{du} \cdot \frac{du}{dv} \cdot \frac{dv}{dx}$$



Example:

$$f(x) = \log(\sqrt{x}) \quad g(x) \equiv \sqrt{x}$$

$$\frac{df}{dx} = \frac{d}{dg} \log(g) \cdot \frac{d}{dx} \sqrt{x}$$

$$= \frac{1}{\sqrt{x}} \cdot \frac{1}{2\sqrt{x}} = \frac{1}{2x}$$

Gradient: Differentiation of *multivariate* function

◆ Partial Derivatives of $f(x)$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

다변수 함수

변수마다 미분한거
벡터로 씀

vector $\rightarrow x \in \mathbb{R}^n$

$$\frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h}$$

\vdots

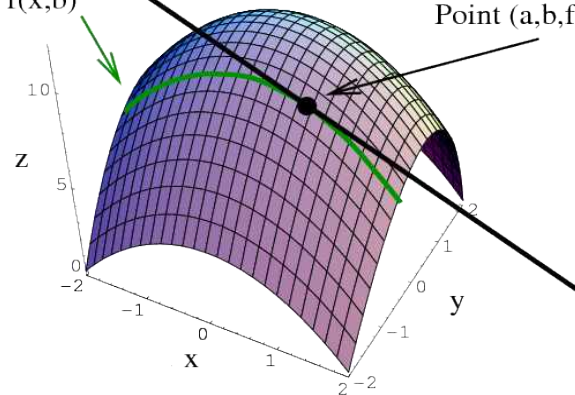
$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(x)}{h}$$

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Line has slope $\frac{\partial f}{\partial x}(a, b)$

Graph of $f(x, b)$

Point $(a, b, f(a, b))$



Gradient: (cont'd)

R

◆ **Gradient:** Collection of the Partial Derivatives

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} = \frac{d}{d\mathbf{x}} f(\mathbf{x})$$

Product Rule: $\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} g(\mathbf{x}) + f(\mathbf{x}) \frac{\partial g}{\partial \mathbf{x}}$

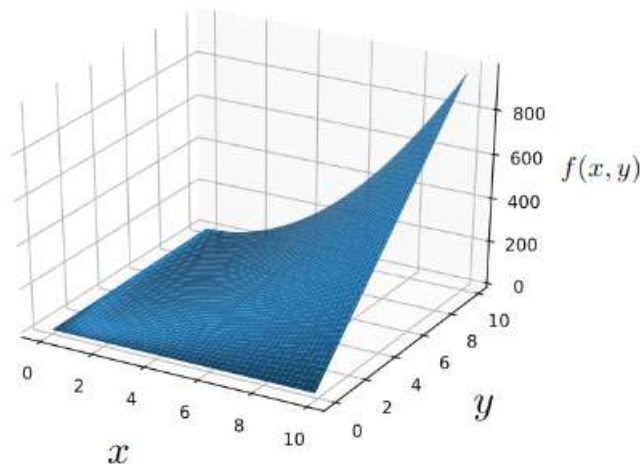
Sum Rule: $\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$

Chain Rule: $\frac{\partial}{\partial \mathbf{x}} (g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (g(f(\mathbf{x}))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}}$

Example:

$$f(x, y) = x^2 y + y$$

$$\nabla f(x, y) = \begin{bmatrix} \partial f / \partial x \\ \partial f / \partial y \end{bmatrix} = \begin{bmatrix} 2xy \\ x^2 + 1 \end{bmatrix}$$



Multivariate Chain Rule

h

◆ Two Variables Case

$$f(g(x), h(x))$$

$$\frac{d}{dx} [f(g(x), h(x))] = \frac{\partial f}{\partial g} \cdot \frac{dg}{dx} + \frac{\partial f}{\partial h} \cdot \frac{dh}{dx}$$

◆ $f(\mathbf{v}(x))$ For $\mathbf{v}(x) = \begin{bmatrix} g(x) \\ h(x) \end{bmatrix}$

$$\frac{df(\mathbf{v})}{dx} = \frac{\partial f(\mathbf{v})}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial x} = \nabla_{\mathbf{v}} f(\mathbf{v}) \frac{\partial \mathbf{v}}{\partial x}$$

$$= \begin{bmatrix} \partial f / \partial g \\ \partial f / \partial h \end{bmatrix} \begin{bmatrix} dg/dx \\ dh/dx \end{bmatrix} = \frac{\partial f}{\partial g} \cdot \frac{dg}{dx} + \frac{\partial f}{\partial h} \cdot \frac{dh}{dx}$$

Example:

$$f(g, h) = g^2 h + h$$

$$\text{where } g(x) = 3x, \text{ and } h(x) = x^2$$

$$\frac{\partial f}{\partial g} = 2gh$$

$$\frac{\partial f}{\partial h} = g^2 + 1$$

$$\frac{dg}{dx} = \frac{d}{dx} 3x = 3$$

$$\frac{dh}{dx} = \frac{d}{dx} x^2 = 2x$$

$$\frac{d}{dx} [f(g(x), h(x))] = [2gh \cdot 3] + [(g^2 + 1) \cdot 2x] = 2xg^2 + 6gh + 2x$$

$$\frac{d}{dx} [f(g(x), h(x))]$$



Jacobian: Gradients of **vector-valued** function

R

◆ Multivariate **Vector-Valued** Function

$$\underset{\text{vector}}{f}(x_1, x_2, \dots, x_n) = \begin{bmatrix} f_1(x_1, x_2, x_3, \dots, x_n) \\ f_2(x_1, x_2, x_3, \dots, x_n) \\ f_3(x_1, x_2, x_3, \dots, x_n) \\ \vdots \\ f_m(x_1, x_2, x_3, \dots, x_n) \end{bmatrix}$$

Input, output $\in \mathbb{R}^n, \mathbb{R}^m$
→ vector function

Example:

$$f: \mathbb{R}^2 \mapsto \mathbb{R}^3$$

$$f(x) = (2x_1 + x_2^2, -x_1^2 + 3x_2, 4x_1x_2)^T$$

$f_1(x)$
 $= f_1(x)$ $= f_2(x)$ $= f_3(x)$

$$J = \begin{pmatrix} 2 & 2x_2 \\ -2x_1 & 3 \\ 4x_2 & 4x_1 \end{pmatrix} \quad J|_{(2,1)}^T = \begin{pmatrix} 2 & 2 \\ -4 & 3 \\ 4 & 8 \end{pmatrix}$$

◆ **Jacobian: Gradients of Vector-Valued Function**

$J(f(x)) = \nabla_x f(x)$
Matrix of the partial derivatives

$$J(x_1, x_2, x_3, \dots, x_m) =$$

↓
m × n 행렬

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} & \dots & \frac{\partial f_2}{\partial x_n} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} & \dots & \frac{\partial f_3}{\partial x_n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \frac{\partial f_m}{\partial x_3} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} (\nabla f_1)^T$$

$$\begin{bmatrix} \nabla_x f_1(x)^T \\ \vdots \\ \nabla_x f_m(x)^T \end{bmatrix}$$



Hessian: 2nd-order differentiation of *multivariate* function

◆ Partial Derivatives of Gradient

2차 미분 ← 1차 미분 2차 미분

2차 미분

$$H(f(x)) = \frac{\partial}{\partial x} \nabla_x f(x) = J(\nabla_x f(x))$$

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

↪ 대칭행렬

Note that the Hessian is always symmetric, since

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}$$

Example:

$$f(x) = f(x_1, x_2)$$

$$= \left(4 - 2.1x_1^2 + \frac{x_1^4}{3} \right) x_1^2 + x_1 x_2 + (-4 + 4x_2^2) x_2^2$$

$$H = \begin{pmatrix} 10x_1^4 - 25.2x_1^2 + 8 & 1 \\ 1 & 48x_2^2 - 8 \end{pmatrix}$$

$$H|_{(0,1)^T} = \begin{pmatrix} 8 & 1 \\ 1 & 40 \end{pmatrix}$$

↪ x_1 1번, x_2 1번
↪ x_2 2번



Useful Formula



R

◆ For Linear Functions

$$\nabla_{\underline{x}} \underline{b}^T \underline{x} = \underline{b} \iff \frac{d}{dx} ax = a$$

$$\nabla_{\underline{x}}^2 \underline{b}^T \underline{x} = 0 \iff \frac{d^2}{dx^2} ax = 0$$

proof)

For $\underline{x} \in \mathbb{R}^n$, let $f(\underline{x}) = \underline{b}^T \underline{x}$ for some known vector $\underline{b} \in \mathbb{R}^n$. Then

$$f(\underline{x}) = \sum_{i=1}^n b_i x_i \quad \frac{\partial f(\underline{x})}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k.$$

대칭행렬
↑

$$\nabla_{\underline{x}} \underline{x}^T A \underline{x} = 2A \underline{x} \text{ (if } A \text{ symmetric)}$$

$$\nabla_{\underline{x}}^2 \underline{x}^T A \underline{x} = 2A \text{ (if } A \text{ symmetric)}$$

$$\begin{aligned} \iff \frac{d}{dx} ax^2 &= 2ax \\ \iff \frac{d^2}{dx^2} ax^2 &= 2a \end{aligned}$$

proof)

$f(\underline{x}) = \underline{x}^T A \underline{x}$ for $A \in \mathbb{S}^n$.

$$f(\underline{x}) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j.$$

$$= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

$$= \frac{\partial}{\partial x_k} \left[\sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right]$$

$$= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j = 2 \sum_{i=1}^n A_{ki} x_i$$

✗ If A is not symmetric,

$$\nabla_{\underline{x}} \underline{x}^T A \underline{x} = (A + A^T) \underline{x}$$



Useful Formula (cont'd)

R

◆ For Linear Functions

$$\nabla_{\underline{x}} \underline{b}^T \underline{x} = \underline{b}$$

$$\nabla_{\underline{x}}^2 \underline{b}^T \underline{x} = 0$$

$$\nabla_{\underline{x}} A \underline{x} = \nabla_{\underline{x}} \begin{bmatrix} a_1^T \underline{x} \\ \vdots \\ a_n^T \underline{x} \end{bmatrix} = \begin{bmatrix} -a_1^T \\ \vdots \\ -a_n^T \end{bmatrix} = A$$

$$A = \begin{bmatrix} -a_1^T \\ \vdots \\ -a_n^T \end{bmatrix}$$

$$\nabla_{\underline{x}} \begin{bmatrix} f_1(\underline{x}) \\ \vdots \\ f_n(\underline{x}) \end{bmatrix} \text{ (Jacobian)} = \begin{bmatrix} \nabla f_1(\underline{x})^T \\ \vdots \\ \nabla f_n(\underline{x})^T \end{bmatrix}$$

$$\nabla_{\underline{x}} \underline{x}^T A \underline{x} = 2A \underline{x} \text{ (if } A \text{ symmetric)}$$

$$\nabla_{\underline{x}}^2 \underline{x}^T A \underline{x} = 2A \text{ (if } A \text{ symmetric)}$$

$$\nabla_{\underline{x}} \|\underline{x}\|_2^2 = \nabla_{\underline{x}} \underline{x}^T \underline{x} = 2 \underline{x}$$

$$\|\underline{x}\|_2^2 = \underline{x}^T \underline{x}$$

$$(\underline{x}^T \underline{x}) = 2 \underline{x} = 2 \underline{x}$$

Exercises :

$$\nabla_{\underline{x}} (\underline{x}^T A^T A \underline{x}) = 2 A^T A \underline{x}$$

$$\nabla_{\underline{x}} (\underline{b}^T A \underline{x}) = \nabla_{\underline{x}} \underline{c}^T \underline{x} = \underline{c} = (\underline{b}^T A)^T = A^T \underline{b}$$

→ 행 벡터

* If A is not symmetric,

$$\nabla_{\underline{x}} \underline{x}^T A \underline{x} = (A + A^T) \underline{x}$$

$$\nabla_{\underline{x}}^2 \underline{x}^T A \underline{x} = (A + A^T)$$



Gradient of a Matrix variables

Suppose that $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a function that takes as input a matrix A of size $m \times n$ and returns a real value. Then the **gradient** of f (with respect to $A \in \mathbb{R}^{m \times n}$) is the matrix of partial derivatives, defined as:

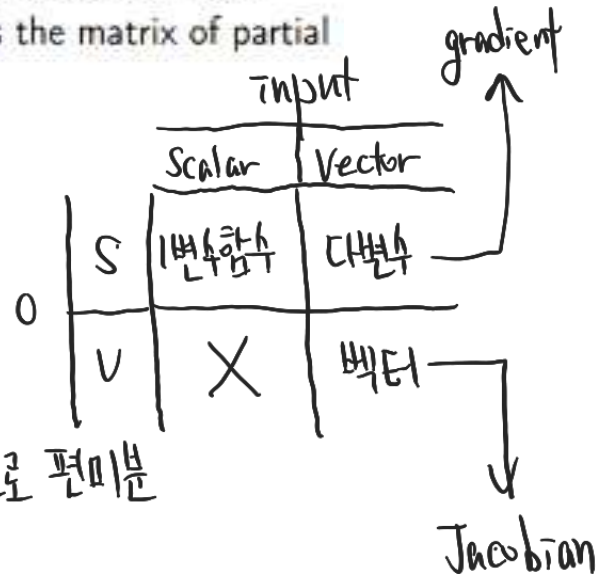
$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \dots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \dots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \dots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

i.e., an $m \times n$ matrix with

$$(\nabla_A f(A))_{ij} = \frac{\partial f(A)}{\partial A_{ij}}$$

j 번째 원소로 편미분




$f\left(\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}\right)$
 ex) 100x100 이미지를 10000 차원 벡터로 처리 = 벡터 함수
 → 그대로 처리 (벡터 함수와 원론적으로 같음)



2. Iterative Optimization & Gradient Descent

1. ML as an Optimization Problem

2. Iterative Optimization

- 
- 
- 

3. Gradient Descent (GD)

4. Stochastic Gradient Descent (SGD)

5. Minibatch Gradient Descent (Minibatch GD)



ML as an Optimization Problem

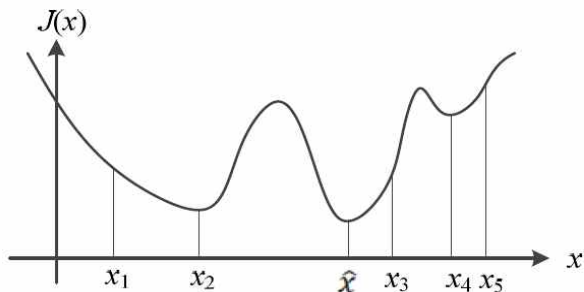
◆ 기계 학습이 해야 할 일을 식으로 정의하면,

주어진 Cost Function $J(\theta)$ 에 대해 θ : 파라미터 (벡터, 행렬)

$J(\theta)$ 를 최소로 하는 최적해 $\hat{\theta}$ 을 찾아라. 즉, $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$

$\min J(\theta): J(\theta)$ 의 최솟값

$\operatorname{argmin} J(\theta): J(\theta)$ 가 최솟일 때의 θ 값



◆ Global Optimum vs. Local Optima

- \hat{x} 은 전역 최적해
- x_2 와 x_4 는 지역 최적해 (근처에서만 최소)



Iterative Optimization – General Principles

◆ Training Dataset: D \rightarrow 이진 분류

● E.g., for binary classification case ($y \in \{0, 1\}$),

$$D = (\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, \langle \mathbf{x}^{[2]}, y^{[2]} \rangle, \dots, \langle \mathbf{x}^{[n]}, y^{[n]} \rangle) \in (\mathbb{R}^m \times \{0, 1\})^n$$

데이터 레이블 / 벡터 스칼라

◆ Model & Predicted Output: $\hat{y} = h_{\theta}(x)$

linear model
 $y = \theta_0 + \theta_1 x$

◆ Cost Function: $J(\theta)$

◆ General Principles

- 1) Initialize parameters (θ) $D \rightarrow$ 일단 랜덤하게 초기화
- 2) For every training epoch (i): \rightarrow 데이터를 한번씩 봄?

★ For every [**some set of training set**]:

- ① Predict output (\hat{y}) & Calculate the cost ($J(\theta)$)
- ② If the cost is satisfactory, then terminates.
- ③ Otherwise, **update parameters** (θ) and repeat

\rightarrow 어떻게 Update 할거냐?

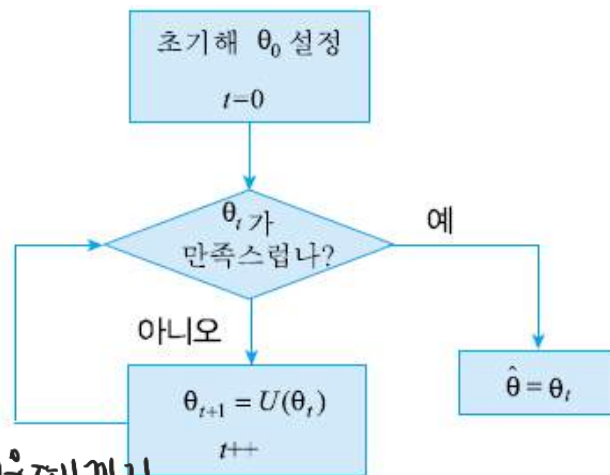


그림 11.6 최적해를 찾기 위한 반복 알고리즘

Iterative Optimization – *Parameter Update*

◆ Typical Form of Parameter Update

1) $J(\theta)$ 가 작아지도록 $\Delta\theta$ 를 구한다.

2) θ 를 업데이트 한다: $\theta = \theta + \Delta\theta \rightarrow$ Optimization 방법

● E.g., *Gradient Descent*: $\Delta\theta = -\alpha \nabla J(\theta)$

$\alpha \rightarrow$ Learning Rate



Iterative Optimization – *Learning Modes*

◆ Training Dataset

$$\mathcal{D} = (\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, \langle \mathbf{x}^{[2]}, y^{[2]} \rangle, \dots, \langle \mathbf{x}^{[n]}, y^{[n]} \rangle) \in (\mathbb{R}^m \times \{0, 1\})^n$$

◆ Learning Modes (Update Modes)

● **On-line** mode: Update parameters for every training **example** 하나 볼 때마다

● **Batch** mode: Update parameters for every training **epoch** 다 보고 평균

1) Initialize parameters (θ)

2) For every training **epoch** (\mathcal{D}):

★ For every [**some set of training set**]:

① Predict output (\hat{y}) & Calculate the cost ($J(\theta)$)

② If the cost is satisfactory, then terminates.

③ Otherwise, **update parameters** (θ) and repeat ★

f training set → Mini-Batch




Gradient Descent (GD)

◆ Parameter Update

- based on the **Gradient** of the Cost Function: $\theta = \theta - \alpha \nabla J(\theta)$

◆ Batch Learning Modes

-  Every training **epoch**
 - ★ Training set에 속한 모든 training example의 Gradient를 평균한 후 한꺼번에 갱신

알고리즘 2-4 배치 경사 하강 알고리즘(BGD)

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

출력: 최적해 $\hat{\theta}$


```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.
2  repeat
3       $\mathbb{X}$ 에 있는 샘플의 그레이디언트  $\nabla_1, \nabla_2, \dots, \nabla_n$ 을 계산한다.
4       $\nabla_{total} = \frac{1}{n} \sum_{i=1, n} \nabla_i$  // 그레이디언트 평균을 계산
5       $\theta = \theta - \rho \nabla_{total}$ 
6  until(멈춤 조건)
7   $\hat{\theta} = \theta$ 
```

업데이트 1번

Stochastic Gradient Descent (SGD)

◆ On-line Learning Modes

↳ Online or Mini-Batch

●  every training **example**

★ 각 training example의 Gradient를 계산한 후 즉시 갱신

★ 결과가 training example들의 순서에 의존하지 않도록 example들의 순서를 "**임의로(stochastic)**" 선택한다.

↳ shuffling (공통 x)
↳ sampling (선택 o)

알고리즘 2-5 스토케스틱 경사 하강 알고리즘(SGD)

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

출력: 최적해 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.  
2  repeat  
3   $\mathbb{X}$ 의 샘플의 순서를 섞는다. Shuffling 매 epoch마다  
4  for ( $i=1$  to  $n$ )  
5       $i$ 번째 샘플에 대한 그레이디언트  $\nabla_i$ 를 계산한다.  
6       $\theta = \theta - \rho \nabla_i$   
7  until(멈춤 조건)  
8   $\hat{\theta} = \theta$  업데이트 n번
```

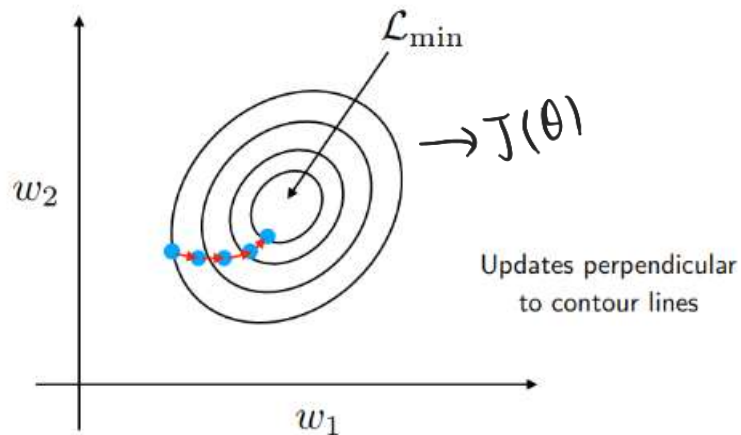
다른 방식의 구현 (독립 샘플링)

Sampling
Sampling

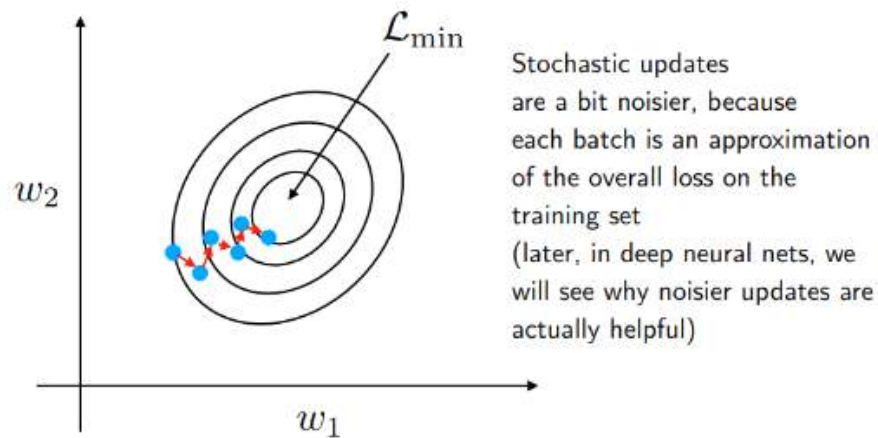
```
3   $\mathbb{X}$ 에서 임의로 샘플 하나를 뽑는다.  
4  뽑힌 샘플의 그레이디언트  $\nabla$ 를 계산한다.  
5   $\theta = \theta - \rho \nabla$ 
```

GD vs. SGD

◆ Batch GD



◆ SGD → Deep Learning에서는 더 도움 됨




3. Automatic Differentiation

실제 미분

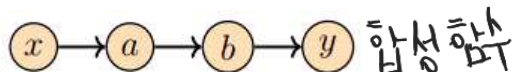
Automatic Differentiation

◆ Automatic Differentiation

-  numerically. derivative of gradient
수치적으로
- Note that it is **not a formula** but a procedure and it is **not symbolic** differentiation.
- Backpropagation is a special case of it.

Most automatic differentiation systems, including Autograd, construct the **computation graph**.

It has two modes: **Reverse mode** & **Forward mode**:



$$y = F(x) = b(a(x))$$

reverse mode $\frac{dy}{dx} = \left(\frac{dy}{db} \frac{db}{da} \right) \frac{da}{dx},$

forward mode, $\frac{dy}{dx} = \frac{dy}{db} \left(\frac{db}{da} \frac{da}{dx} \right).$

<https://arxiv.org/pdf/1502.05767.pdf>

https://www.cs.toronto.edu/~rgrosse/courses/csc321_2018/slides/lec10.pdf



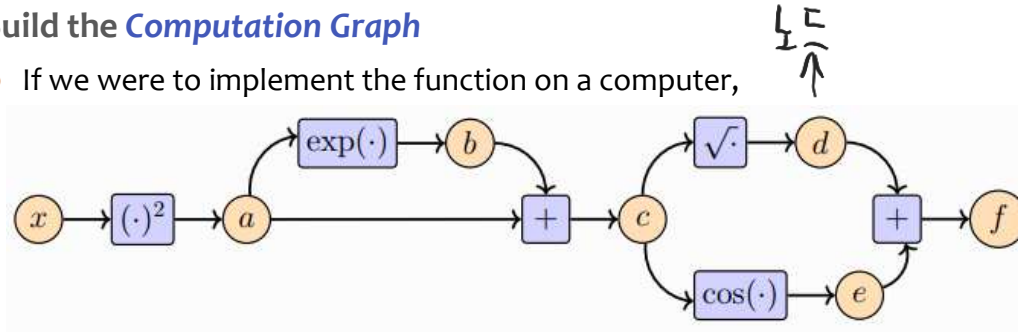
Automatic Differentiation – *Illustration*

Consider a function

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

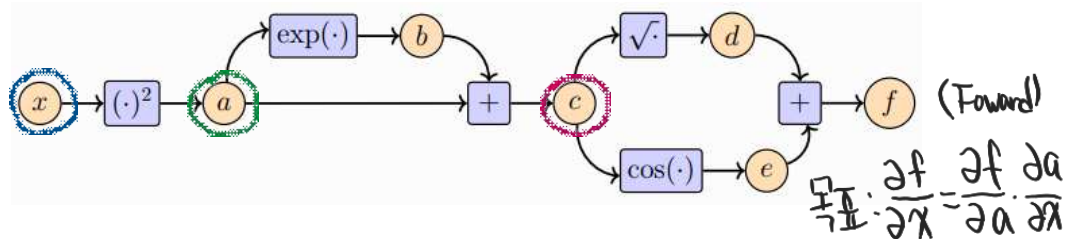
1) Build the *Computation Graph*

- If we were to implement the function on a computer,



Automatic Differentiation – Illustration (cont'd)

R



2) Calculate the values and the derivatives

$$a = x^2,$$

$$b = \exp(a),$$

$$c = a + b,$$

$$d = \sqrt{c},$$

$$e = \cos(c),$$

$$f = d + e.$$

$$\frac{\partial a}{\partial x} = 2x$$

$$\frac{\partial b}{\partial a} = \exp(a)$$

$$\frac{\partial d}{\partial c} = \frac{1}{2\sqrt{c}}$$

$$\frac{\partial e}{\partial c} = -\sin(c)$$

각 편미분을 변식으로 생각

3) Apply Chain Rule (in reverse mode)

$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial d} + \frac{\partial f}{\partial e}$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial c} \cdot \frac{\partial c}{\partial a} = \frac{\partial f}{\partial c} \cdot \left(\frac{\partial b}{\partial a} + 1 \right)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \cdot \frac{\partial a}{\partial x}$$

Complexity $\frac{1}{2}$ Forward
Reverse $\frac{1}{2}$ $\frac{1}{2}$

"Backpropagation"

- we observe that the computation required for calculating the derivatives is **of similar complexity** as the computation of the function itself.