

기계학습 기초수학

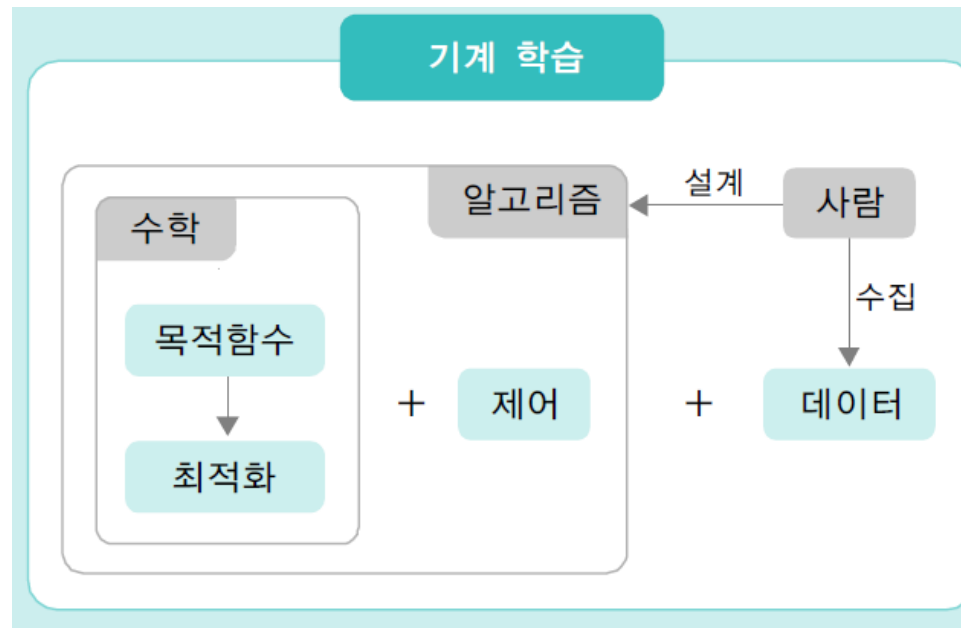
이현석 교수

2024-1

PREVIEW

■ 기계 학습에서 수학의 역할

- **수학**은 목적함수를 정의하고, 목적함수가 최저가 되는 점을 찾아주는 최적화 이론 제공
- 최적화 이론에 규제, 모멘텀, 학습률, 멈춤조건과 같은 제어를 추가하여 **알고리즘** 구축
- **사람**은 알고리즘을 설계하고 데이터를 수집함



PREVIEW

- 선형대수: 이 분야의 개념을 이용하면 학습 모델의 매개변수집합, 데이터, 선형연산의 결합 등을 행렬 또는 텐서로 간결하게 표현할 수 있다. 데이터를 분석하여 유용한 정보를 알아내거나 특징 공간을 변환하는 등의 과업을 수행하는 데 핵심 역할을 한다.
- 확률과 통계: 데이터에 포함된 불확실성을 표현하고 처리하는 데 활용한다. 베이즈 이론과 최대 우도 기법을 이용하여 확률 추론을 수행한다.
- 최적화: 목적함수를 최소화하는 최적해를 찾는 데 활용하며, 주로 미분을 활용한 방법을 사용한다. 수학자들이 개발한 최적화 방법을 기계 학습이라는 도메인에 어떻게 효율적으로 적용할지가 주요 관심사이다.

벡터와 행렬

■ 벡터

- 샘플을 특징 벡터로 feature vector 표현
- 예) Iris 데이터에서 꽃받침의 길이, 꽃받침의 너비, 꽃잎의 길이, 꽃잎의 너비라는 4개의 특징이 각각 5.1, 3.5, 1.4, 0.2인 샘플

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 5.1 \\ 3.5 \\ 1.4 \\ 0.2 \end{pmatrix}$$

- 여러 개의 특징 벡터를 첨자로 구분

$$\mathbf{x}_1 = \begin{pmatrix} 5.1 \\ 3.5 \\ 1.4 \\ 0.2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 4.9 \\ 3.0 \\ 1.4 \\ 0.2 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 4.7 \\ 3.2 \\ 1.3 \\ 0.2 \end{pmatrix}, \dots, \mathbf{x}_{150} = \begin{pmatrix} 5.9 \\ 3.0 \\ 5.1 \\ 1.8 \end{pmatrix}$$

벡터와 행렬

■ 행렬

- 여러 개의 벡터를 담음
- 훈련집합을 담은 행렬을 설계행렬이라 부름
- 예) Iris 데이터에 있는 150개의 샘플을 설계 행렬 \mathbf{X} 로 표현

$$\mathbf{X} = \begin{pmatrix} 5.1 & 3.5 & 1.4 & 0.2 \\ 4.9 & 3.0 & 1.4 & 0.2 \\ 4.7 & 3.2 & 1.3 & 0.2 \\ 4.6 & 3.1 & 1.5 & 0.2 \\ \vdots & \vdots & \vdots & \vdots \\ 6.2 & 3.4 & 5.4 & 2.3 \\ 5.9 & 3.0 & 5.1 & 1.8 \end{pmatrix} = \begin{pmatrix} x_{1,1} & x_{1,2} & x_{1,3} & x_{1,4} \\ x_{2,1} & x_{2,2} & x_{2,3} & x_{2,4} \\ x_{3,1} & x_{3,2} & x_{3,3} & x_{3,4} \\ x_{4,1} & x_{4,2} & x_{4,3} & x_{4,4} \\ \vdots & \vdots & \vdots & \vdots \\ x_{149,1} & x_{149,2} & x_{149,3} & x_{149,4} \\ x_{150,1} & x_{150,2} & x_{150,3} & x_{150,4} \end{pmatrix}$$

← 행row

↑
열column

벡터와 행렬

■ 행렬 A 의 전치행렬 A^T

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix}, \quad A^T = \begin{pmatrix} a_{11} & a_{21} & \cdots & a_{n1} \\ a_{12} & a_{22} & \cdots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1m} & a_{2m} & \cdots & a_{nm} \end{pmatrix}$$

예를 들어, $A = \begin{pmatrix} 3 & 4 & 1 \\ 0 & 5 & 2 \end{pmatrix}$ 라면 $A^T = \begin{pmatrix} 3 & 0 \\ 4 & 5 \\ 1 & 2 \end{pmatrix}$

- Iris의 설계 행렬을 전치행렬 표기에 따라 표현하면,

$$X = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_{150}^T \end{pmatrix}$$

벡터와 행렬

■ 행렬을 이용하면 수학을 간결하게 표현할 수 있음

- 예) 다항식의 행렬 표현

$$f(\mathbf{x}) = f(x_1, x_2, x_3)$$

$$= 2x_1x_1 - 4x_1x_2 + 3x_1x_3 + x_2x_1 + 2x_2x_2 + 6x_2x_3 - 2x_3x_1 + 3x_3x_2 + 2x_3x_3 + 2x_1 + 3x_2 - 4x_3 + 5$$

$$= (x_1 \quad x_2 \quad x_3) \begin{pmatrix} 2 & -4 & 3 \\ 1 & 2 & 6 \\ -2 & 3 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + (2 \quad 3 \quad -4) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + 5$$

$$= \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$$

■ 특수한 행렬들

$$\text{정사각행렬} \begin{pmatrix} 2 & 0 & 1 \\ 1 & 21 & 5 \\ 4 & 5 & 12 \end{pmatrix}, \quad \text{대각행렬} \begin{pmatrix} 50 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 8 \end{pmatrix},$$

$$\text{단위행렬} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \text{대칭행렬} \begin{pmatrix} 1 & 2 & 11 \\ 2 & 21 & 5 \\ 11 & 5 & 1 \end{pmatrix}$$

벡터와 행렬

■ 행렬 연산

■ 행렬 곱셈 $\mathbf{C} = \mathbf{AB}$, 이때 $c_{ij} = \sum_{k=1,s} a_{ik} b_{kj}$

2*3 행렬 $\mathbf{A} = \begin{pmatrix} 3 & 4 & 1 \\ 0 & 5 & 2 \end{pmatrix}$ 와 3*3 행렬 $\mathbf{B} = \begin{pmatrix} 2 & 0 & 1 \\ 1 & 0 & 5 \\ 4 & 5 & 1 \end{pmatrix}$ 을 곱하면 2*3 행렬 $\mathbf{C} = \mathbf{AB} = \begin{pmatrix} 14 & 5 & 24 \\ 13 & 10 & 27 \end{pmatrix}$

• 교환법칙 성립하지 않음: $\mathbf{AB} \neq \mathbf{BA}$

• 분배법칙과 결합법칙 성립: $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$ 이고 $\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}$

■ 벡터의 내적 $\mathbf{a} \cdot \mathbf{b} = \mathbf{a}^T \mathbf{b} = \sum_{k=1,d} a_k b_k$

$\mathbf{x}_1 = \begin{pmatrix} 5.1 \\ 3.5 \\ 1.4 \\ 0.2 \end{pmatrix}$ 와 $\mathbf{x}_2 = \begin{pmatrix} 4.9 \\ 3.0 \\ 1.4 \\ 0.2 \end{pmatrix}$ 의 내적 $\mathbf{x}_1 \cdot \mathbf{x}_2$ 는 37.49

벡터와 행렬

■ 텐서

- 3차원 이상의 구조를 가진 숫자 배열
- 예) 3차원 구조의 RGB 컬러 영상

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 0 & 3 & 2 & 2 \\ 2 & 0 & 2 & 2 & 3 & 1 \\ 3 & 0 & 1 & 2 & 6 & 7 \\ 3 & 1 & 2 & 3 & 5 & 6 \\ 1 & 2 & 2 & 2 & 2 & 3 \\ 3 & 0 & 0 & 1 & 1 & 0 \\ 5 & 4 & 1 & 3 & 3 & 3 \\ 2 & 2 & 1 & 2 & 2 & 1 \end{pmatrix} \begin{pmatrix} 6 \\ 3 \\ 0 \\ 3 \\ 1 \end{pmatrix}$$

놈과 유사도

■ 벡터와 행렬의 크기를 놈으로 측정

- 벡터의 p차 놈

$$p\text{차 놈: } \|\mathbf{x}\|_p = \left(\sum_{i=1,d} |x_i|^p \right)^{\frac{1}{p}}$$

$$\text{최대 놈: } \|\mathbf{x}\|_\infty = \max(|x_1|, |x_2|, \dots, |x_d|)$$

- 예) $\mathbf{x} = (3 \ -4 \ 1)$ 일 때, 2차 놈은 $\|\mathbf{x}\|_2 = (3^2 + (-4)^2 + 1^2)^{1/2} = 5.099$

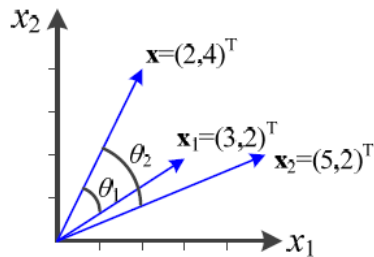
- 행렬의 프로베니우스 놈 $\|\mathbf{A}\|_F = \left(\sum_{i=1,n} \sum_{j=1,m} a_{ij}^2 \right)^{\frac{1}{2}}$

$$\text{예를 들어, } \left\| \begin{pmatrix} 2 & 1 \\ 6 & 4 \end{pmatrix} \right\|_F = \sqrt{2^2 + 1^2 + 6^2 + 4^2} = 7.550$$

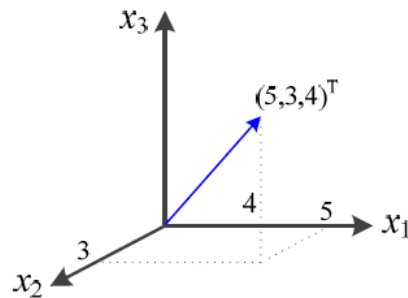
놈과 유사도

■ 유사도와 거리

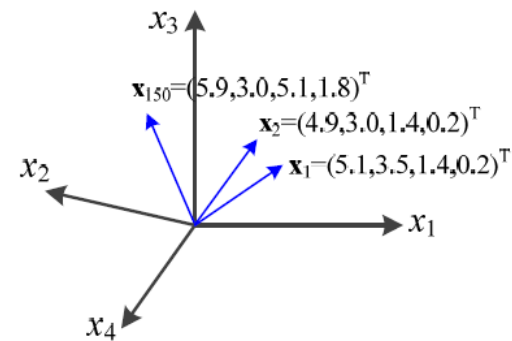
- 벡터를 기하학적으로 해석



(a) 2차원 벡터



(b) 3차원 벡터



(c) 4차원 벡터(Iris 데이터)

- 코사인 유사도

$$\text{cosine_similarity}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a}}{\|\mathbf{a}\|} \cdot \frac{\mathbf{b}}{\|\mathbf{b}\|} = \cos(\theta)$$

선형결합과 벡터공간

■ 벡터

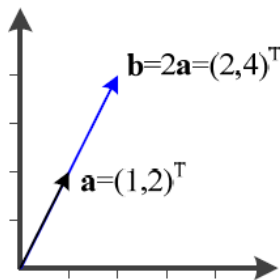
- 공간상의 한 점으로 화살표 끝이 벡터의 좌표에 해당

■ 선형결합이 만드는 벡터공간

- 기저벡터 \mathbf{a} 와 \mathbf{b} 의 선형결합

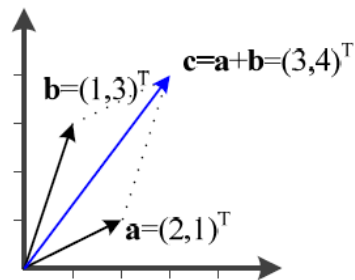
$$\mathbf{c} = \alpha_1 \mathbf{a} + \alpha_2 \mathbf{b}$$

- 선형결합으로 만들어지는 공간을 **벡터공간**이라 부름

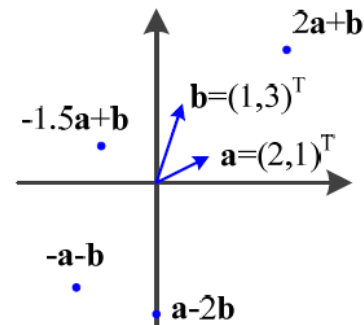


(a) 벡터에 스칼라 곱

그림 2-6 벡터의 연산

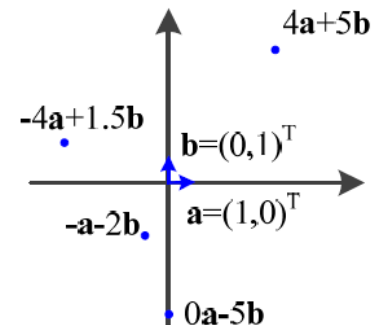


(b) 두 벡터의 덧셈



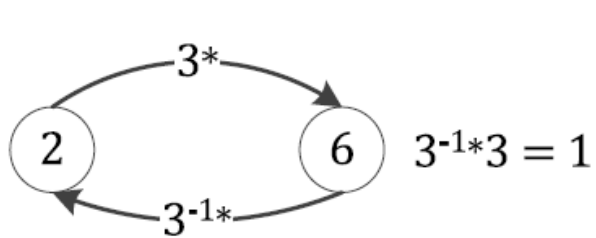
(a) 기저 벡터와 벡터공간

그림 2-7 벡터공간

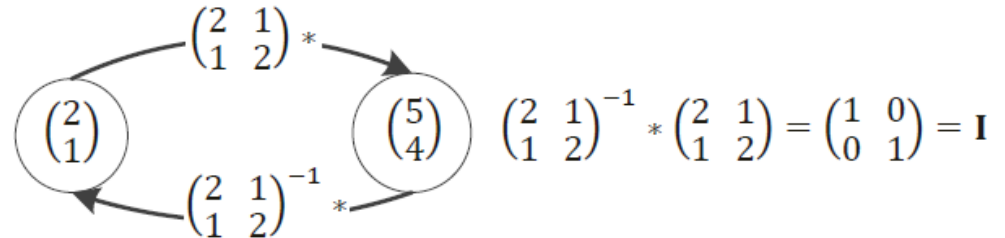


(b) 정규직교 기저 벡터

■ 역행렬의 원리



(a) 역수의 원리



(b) 역행렬의 원리

그림 2-9 역행렬

- 정사각행렬 \mathbf{A} 의 역행렬 \mathbf{A}^{-1}

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$$

- 예를 들어, $\begin{pmatrix} 2 & 1 \\ 6 & 4 \end{pmatrix}$ 의 역행렬은 $\begin{pmatrix} 2 & -0.5 \\ -3 & 1 \end{pmatrix}$

■ 정리

정리 2-1 다음 성질은 서로 필요충분조건이다.

- A 는 역행렬을 가진다. 즉, 특이행렬이 아니다.
- A 는 최대계수를 가진다.
- A 의 모든 행이 선형독립이다.
- A 의 모든 열이 선형독립이다.
- A 의 행렬식은 0이 아니다.
- $A^T A$ 는 양의 정부호 positive definite 대칭 행렬이다.
- A 의 고윳값은 모두 0이 아니다.

행렬 분해

■ 분해란?

- 정수 3717은 특성이 보이지 않지만, $3 \times 3 \times 7 \times 59$ 로 소인수 분해를 하면 특성이 보이듯이, 행렬도 분해하면 여러모로 유용함

■ 고윳값과 고유 벡터

- 고유 벡터 \mathbf{v} 와 고윳값 λ

$$A\mathbf{v} = \lambda\mathbf{v}$$

- 예를 들어, $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 3 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ 이고 $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 1 \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ 이므로, $\lambda_1 = 3, \lambda_2 = 1$ 이고 $\mathbf{v}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $\mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$

행렬 분해

■ 고윳값 분해 eigen value decomposition

$$A = Q\Lambda Q^{-1} \quad (2.21)$$

- Q 는 A 의 고유 벡터를 열에 배치한 행렬이고 Λ 는 고윳값을 대각선에 배치한 대각행렬
- 예를 들어, $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \end{pmatrix}$
- 고윳값 분해는 정사각행렬에만 적용 가능한데, 기계 학습에서는 정사각행렬이 아닌 경우의 분해도 필요하므로 고윳값 분해는 한계를 가짐

행렬 분해

- $n*m$ 행렬 \mathbf{A} 의 특잇값 분해 SVD(singular value decomposition)

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (2.22)$$

- 왼쪽 특이행렬 \mathbf{U} 는 $\mathbf{A}\mathbf{A}^T$ 의 고유 벡터를 열에 배치한 $n*n$ 행렬
- 오른쪽 특이행렬 \mathbf{V} 는 $\mathbf{A}^T\mathbf{A}$ 의 고유 벡터를 열에 배치한 $m*m$ 행렬
- $\mathbf{\Sigma}$ 는 $\mathbf{A}\mathbf{A}^T$ 의 고유값의 제곱근을 대각선에 배치한 $n*m$ 대각행렬

예를 들어, \mathbf{A} 를 $4*3$ 행렬이라고 했을 때 다음과 같이 특잇값 분해가 된다.

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 2 \\ 3 & 1 & 0 \\ 1 & 0 & 2 \end{pmatrix} = \begin{pmatrix} -0.1914 & -0.2412 & 0.1195 & -0.9439 \\ -0.5144 & 0.6990 & -0.4781 & -0.1348 \\ -0.6946 & -0.6226 & -0.2390 & 0.2697 \\ -0.4651 & 0.2560 & 0.8367 & 0.1348 \end{pmatrix}$$

$$\begin{pmatrix} 3.7837 & 0 & 0 \\ 0 & 2.7719 & 0 \\ 0 & 0 & 1.4142 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} -0.7242 & -0.4555 & -0.5177 \\ -0.6685 & 0.2797 & 0.6891 \\ 0.1690 & -0.8452 & 0.5071 \end{pmatrix}$$

확률과 통계

- 기계 학습이 처리할 데이터는 불확실한 세상에서 발생하므로, 불확실성을 다루는 확률과 통계를 잘 활용해야 함

■ 확률변수 random variable

■ 예) 윷



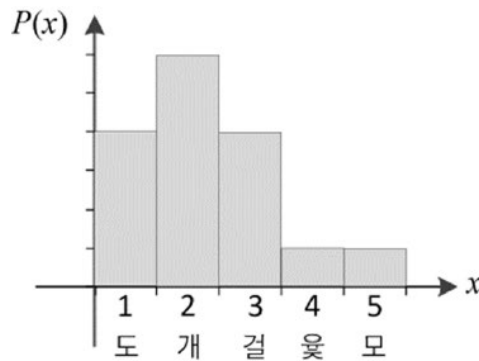
그림 2-13 윷을 던졌을 때 나올 수 있는 다섯 가지 경우(왼쪽부터 도, 개, 걸, 윷, 모)

- 다섯 가지 경우 중 한 값을 갖는 확률변수 x
- x 의 정의역은 {도, 개, 걸, 윷, 모}

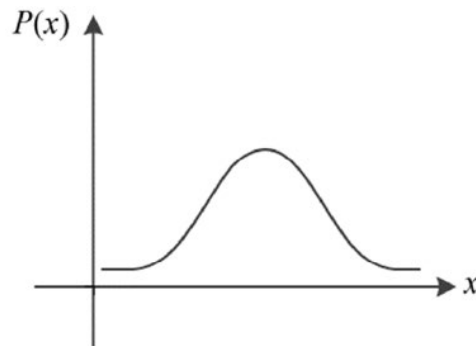
확률 기초

■ 확률분포

$$P(x = \text{도}) = \frac{4}{16}, P(x = \text{개}) = \frac{6}{16}, P(x = \text{걸}) = \frac{4}{16}, P(x = \text{웃}) = \frac{1}{16}, P(x = \text{모}) = \frac{1}{16}$$



(a) 이산인 경우의 확률질량함수



(b) 연속인 경우의 확률밀도함수

그림 2-14 확률분포

■ 확률벡터 random vector

- 예) Iris에서 확률벡터 \mathbf{x} 는 4차원 $\mathbf{x} = (x_1, x_2, x_3, x_4)^T = (\text{꽃받침 길이}, \text{꽃받침 너비}_1, \text{꽃잎 길이}, \text{꽃잎 너비})^T$

■ 간단한 확률실험 장치

- 주머니에서 번호를 뽑은 다음, 번호에 따라 해당 병에서 공을 뽑고 색을 관찰함
- 번호를 y , 공의 색을 x 라는 확률변수로 표현하면 정의역은 $y \in \{①, ②, ③\}$, $x \in \{\text{파랑, 하양}\}$

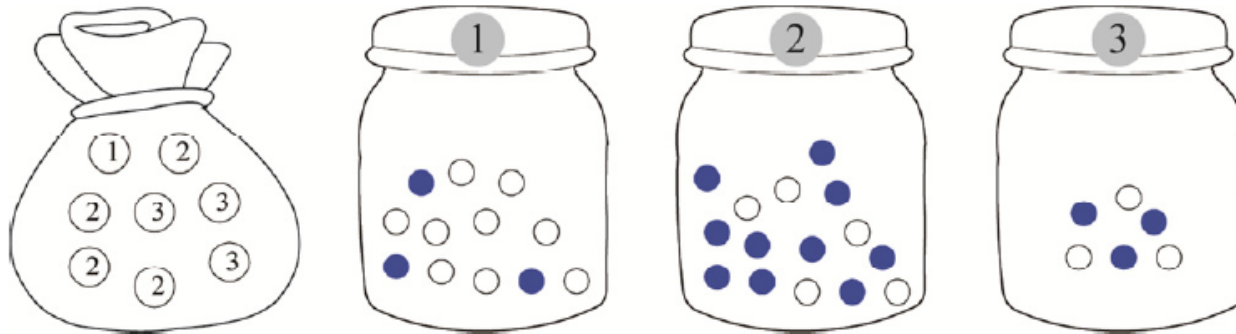


그림 2-15 확률 실험

확률 기초

■ 곱 규칙과 합 규칙

- ①번 카드를 뽑을 확률은 $P(y=①)=P(①)=1/8$
- 카드는 ①번, 공은 하양일 확률은 $P(y=①, x=하양)=P(①, 하양) \leftarrow$ 결합확률

$$P(y = ①, x = 하양) = P(x = 하양 | y = ①)P(y = ①) = \frac{9}{12} \frac{1}{8} = \frac{3}{32}$$

- 곱 규칙 $P(y, x) = P(x|y)P(y)$ (2.23)

- 하얀 공이 뽑힐 확률

$$\begin{aligned} P(\text{하양}) &= P(\text{하양}|①)P(①) + P(\text{하양}|②)P(②) + P(\text{하양}|③)P(③) \\ &= \frac{9}{12} \frac{1}{8} + \frac{5}{15} \frac{4}{8} + \frac{3}{6} \frac{3}{8} = \frac{43}{96} \end{aligned}$$

- 합 규칙 $P(x) = \sum_y P(y, x) = \sum_y P(x|y)P(y)$ (2.24)

베이즈 정리와 기계 학습

■ 베이즈 정리 (식 (2.26))

$$P(y, x) = P(x|y)P(y) = P(x, y) = P(y|x)P(x)$$

$$\longrightarrow P(y|x) = \frac{P(x|y)P(y)}{P(x)} \quad (2.26)$$

- 다음 질문을 식 (2.27)로 쓸 수 있음

“하얀 공이 나왔다는 사실만 알고 어느 병에서 나왔는지 모르는데, 어느 병인지 추정하라.”

$$\hat{y} = \operatorname{argmax}_y P(y|x) \quad (2.27)$$

$P([1;2;3] / \quad)$ 가 argmax

베이즈 정리와 기계 학습

■ 베이즈 정리 (식 (2.26))

- 베이즈 정리를 적용하면, $\hat{y} = \operatorname{argmax}_y P(y|x = \text{하양}) = \operatorname{argmax}_y \frac{P(x = \text{하양}|y)P(y)}{P(x = \text{하양})}$

- 세 가지 경우에 대해 확률을 계산하면,

$$P(\textcircled{1}|\text{하양}) = \frac{P(\text{하양}|\textcircled{1})P(\textcircled{1})}{P(\text{하양})} = \frac{\frac{9}{12} \frac{1}{8}}{\frac{43}{96}} = \frac{9}{43}$$

$$P(\textcircled{2}|\text{하양}) = \frac{P(\text{하양}|\textcircled{2})P(\textcircled{2})}{P(\text{하양})} = \frac{\frac{5}{15} \frac{4}{8}}{\frac{43}{96}} = \frac{16}{43} \longrightarrow \textcircled{3}\text{번 병일 확률이 가장 높음}$$

$$P(\textcircled{3}|\text{하양}) = \frac{P(\text{하양}|\textcircled{3})P(\textcircled{3})}{P(\text{하양})} = \frac{\frac{3}{6} \frac{3}{8}}{\frac{43}{96}} = \frac{18}{43}$$

■ 베이즈 정리의 해석

$$\overbrace{P(y|x)}^{\text{사후확률}} = \frac{\overbrace{P(x|y)}^{\text{우도}} \overbrace{P(y)}^{\text{사전확률}}}{P(x)}$$

베이즈 정리와 기계 학습

■ 기계 학습에 적용

- 예) Iris 데이터 분류 문제
 - 특징 벡터 \mathbf{x} , 부류 $y \in \{\text{setosa}, \text{versicolor}, \text{virginica}\}$
 - 분류 문제를 argmax 로 표현하면 식 (2.29)

$$\hat{y} = \underset{y}{\text{argmax}} P(y|\mathbf{x}) \quad (2.29)$$

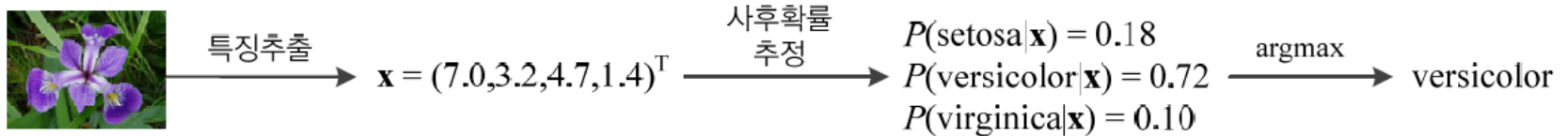


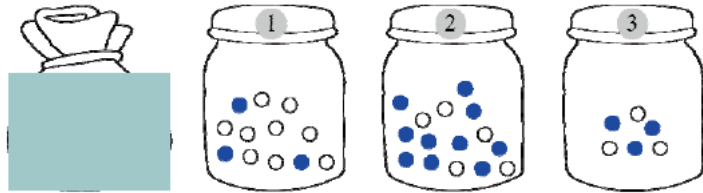
그림 2-16 붓꽃의 부류 예측 과정

- 사후확률 $P(y|\mathbf{x})$ 를 직접 추정하는 일은 아주 단순한 경우를 빼고 불가능
- 따라서 베이즈 정리를 이용하여 추정함

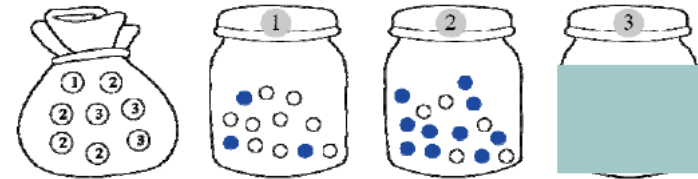
- 사전확률은 식 (2.30)으로 추정
- 우도는 밀도 추정 기법으로 추정

$$\text{사전확률: } P(y = c_i) = \frac{n_i}{n} \quad (2.30)$$

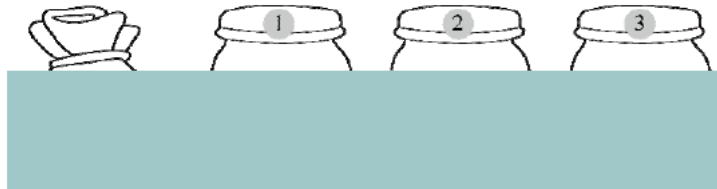
■ 매개변수 θ 를 모르는 상황에서 매개변수를 추정하는 문제



(a) $\theta = \{p_1, p_2\}$



(b) $\theta = \{q_3\}$



(c) $\theta = \{p_1, p_2, q_1, q_2, q_3\}$

그림 2-17 매개변수가 감추어진 여러 가지 상황

■ 예) [그림 2-17(b)] 상황

데이터집합 $\mathbb{X} = \{\bullet \circ \circ \bullet \circ \bullet \circ \circ \bullet \bullet \circ \circ\}$

“데이터 \mathbb{X} 가 주어졌을 때, \mathbb{X} 를 발생시켰을 가능성을 최대로 하는 매개변수 $\theta = \{q_3\}$ 의 값을 찾아라.”

최대 우도

■ 최대 우도법

- [그림 2-17(b)] 문제를 수식으로 쓰면,

$$\hat{q}_3 = \operatorname{argmax}_{q_3} P(\mathbb{X}|q_3) \quad (2.31)$$

- 일반화 하면,

$$\text{최대 우도 추정: } \hat{\Theta} = \operatorname{argmax}_{\Theta} P(\mathbb{X}|\Theta) \quad (2.32)$$

- 수치 문제를 피하기 위해 로그 표현으로 바꾸면,

$$\text{최대 로그우도 추정: } \hat{\Theta} = \operatorname{argmax}_{\Theta} \log P(\mathbb{X}|\Theta) = \operatorname{argmax}_{\Theta} \sum_{i=1}^n \log P(\mathbf{x}_i|\Theta) \quad (2.34)$$

평균과 분산

■ 데이터의 요약 정보로서 평균과 분산

$$\left. \begin{array}{l} \text{평균 } \mu = \frac{1}{n} \sum_{i=1}^n x_i \\ \text{분산 } \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \end{array} \right\} \quad (2.36)$$

■ 평균 벡터와 공분산 행렬

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (2.37)$$

$$\boldsymbol{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \quad (2.39)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{pmatrix}$$

■ 평균 벡터와 공분산 행렬 예제

예제 2-7

Iris 데이터베이스의 샘플 중 8개만 가지고 공분산 행렬을 계산하자.

$$\mathbb{X} = \{\mathbf{x}_1 = \begin{pmatrix} 5.1 \\ 3.5 \\ 1.4 \\ 0.2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 4.9 \\ 3.0 \\ 1.4 \\ 0.2 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 4.7 \\ 3.2 \\ 1.3 \\ 0.2 \end{pmatrix}, \mathbf{x}_4 = \begin{pmatrix} 4.6 \\ 3.1 \\ 1.5 \\ 0.2 \end{pmatrix}, \mathbf{x}_5 = \begin{pmatrix} 5.0 \\ 3.6 \\ 1.4 \\ 0.2 \end{pmatrix}, \mathbf{x}_6 = \begin{pmatrix} 5.4 \\ 3.9 \\ 1.7 \\ 0.4 \end{pmatrix}, \mathbf{x}_7 = \begin{pmatrix} 4.6 \\ 3.4 \\ 1.4 \\ 0.3 \end{pmatrix}, \mathbf{x}_8 = \begin{pmatrix} 5.0 \\ 3.4 \\ 1.5 \\ 0.2 \end{pmatrix}\}$$

먼저 평균벡터를 구하면 $\boldsymbol{\mu} = (4.9125, 3.3875, 1.45, 0.2375)^T$ 이다. 첫 번째 샘플 \mathbf{x}_1 을 식 (2.39)에 적용하면 다음과 같다.

$$\begin{aligned} (\mathbf{x}_1 - \boldsymbol{\mu})(\mathbf{x}_1 - \boldsymbol{\mu})^T &= \begin{pmatrix} 0.1875 \\ 0.1125 \\ -0.05 \\ -0.0375 \end{pmatrix} \begin{pmatrix} 0.1875 & 0.1125 & -0.05 & -0.0375 \end{pmatrix} \\ &= \begin{pmatrix} 0.0325 & 0.0211 & -0.0094 & -0.0070 \\ 0.0211 & 0.0127 & -0.0056 & -0.0042 \\ -0.0094 & -0.0056 & 0.0025 & 0.0019 \\ -0.0070 & -0.0042 & 0.0019 & 0.0014 \end{pmatrix} \end{aligned}$$

나머지 7개 샘플도 같은 계산을 한 다음, 결과를 모두 더하고 8로 나누면 다음과 같은 공분산 행렬을 얻는다.

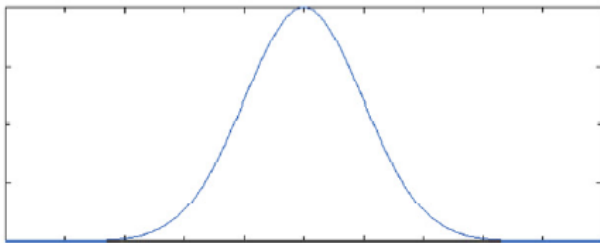
$$\boldsymbol{\Sigma} = \begin{pmatrix} 0.0661 & 0.0527 & 0.0181 & 0.0083 \\ 0.0527 & 0.0736 & 0.0181 & 0.0130 \\ 0.0181 & 0.0181 & 0.0125 & 0.0056 \\ 0.0083 & 0.0130 & 0.0056 & 0.0048 \end{pmatrix}$$

유용한 확률분포

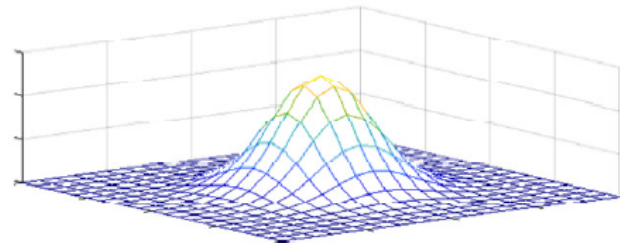
■ 가우시안 분포

- 평균 μ 와 분산 σ^2 으로 정의

$$N(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right)$$



(a) 1차원



(b) 2차원

그림 2-19 가우시안 분포

- 다차원 가우시안 분포: 평균벡터 $\boldsymbol{\mu}$ 와 공분산행렬 $\boldsymbol{\Sigma}$ 로 정의

$$N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|\boldsymbol{\Sigma}|}\sqrt{(2\pi)^d}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

유용한 확률분포

■ 베르누이 분포

- 성공($x=1$) 확률 p 이고 실패($x=0$) 확률이 $1-p$ 인 분포

$$Ber(x; p) = p^x(1-p)^{1-x} = \begin{cases} p, & x = 1 \text{ 일 때} \\ 1-p, & x = 0 \text{ 일 때} \end{cases}$$

■ 이항 분포

- 성공 확률이 p 인 베르누이 실험을 m 번 수행할 때 성공할 횟수의 확률분포

$$B(x; m, p) = C_m^x p^x (1-p)^{m-x} = \frac{m!}{x! (m-x)!} p^x (1-p)^{m-x}$$

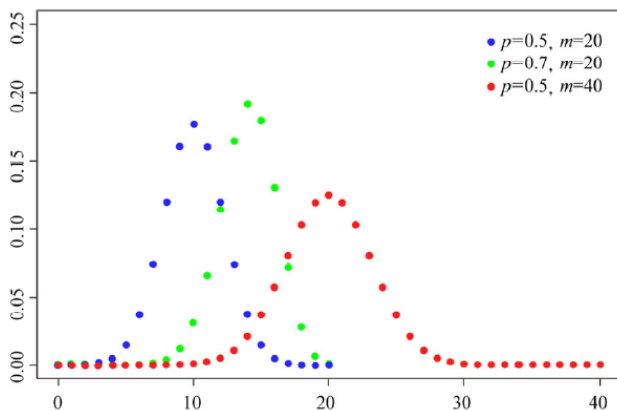


그림 2-20 이항 분포

최적화

■ 순수 수학 최적화와 기계 학습 최적화의 차이

- 순수 수학의 최적화 예) $f(x_1, x_2) = -(\cos(x_1^2) + \sin(x_2^2))^2$ 의 최저점을 찾아라.
- 기계 학습의 최적화는 단지 **훈련집합**이 주어지고, 훈련집합에 따라 정해지는 목적함수의 최저점을 찾아야 함
 - 데이터로 미분하는 과정 필요 → 오류 역전파 알고리즘 (딥러닝)
 - 주로 SGD(스토캐스틱 경사 하강법) 사용

최적화 이론

■ 최적의 선택을 찾는 것

- 무엇이 “최적”인가?
- 모든 가능한 $\mathbf{x} \in \Omega$ 중 함수 $f(\mathbf{x})$ 를 최소화 (혹은 최대화)하는 \mathbf{x} 를 찾는

■ 최적화 문제의 수학적 표현

$$\begin{array}{ll} \text{minimize or maximize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \Omega \end{array}$$

- $f(\mathbf{x})$: 목적함수 (objective function)
- $\mathbf{x} = [x_1, \dots, x_n]^T$: 결정변수벡터 (vector of decision variables)
- $\Omega \subset \mathbb{R}^n$: 가능해집합 (feasible set)
- 목적함수 f 를 최소화(최대화)하는 \mathbf{x} 를 최적해 (optimal solution) 라 함

■ 예1)

$$\begin{array}{ll}\text{minimize} & x^2 - 2x + 1 \\ \text{subject to} & x \in [-2, 2]\end{array}$$

■ 예2)

$$\begin{array}{ll}\text{minimize} & x^2 - 2x + 1 \\ \text{subject to} & x \in [-6, -2]\end{array}$$

최적화 이론

■ 예3) 선형회귀

- n 개의 점 $(x_1, y_1), \dots, (x_n, y_n)$
- 위의 점을 가장 적은 오차로 표현하는 직선 $y = ax + b$ 를 찾고 싶을 경우
- 목적 \rightarrow 평균제곱오차를 최소화하는 직선의 매개변수 a, b 를 찾는 최적화 문제

$$\underset{a,b}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n (ax_i + b - y_i)^2$$

- n, x_i, y_i : 상수
- a, b : 결정변수

매개변수 공간의 탐색

■ 학습 모델의 매개변수 공간

- 특징 공간보다 수 배~수만 배 넓음 (딥러닝의 경우)
 - 선형회귀에서는 특징 공간은 1차원, 매개변수 공간은 2차원
 - MNIST 인식하는 딥러닝 모델은 784차원 특징 공간, 수십만~수백만 차원의 매개변수 공간
- [그림 2-23] 개념도의 매개변수 공간: \hat{x} 은 전역 최적해, x_2 와 x_4 는 지역 최적해
- x_2 와 같이 전역 최적해에 가까운 지역 최적해를 찾고 만족하는 경우 많음

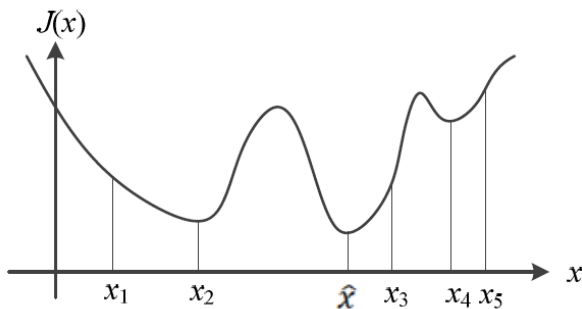
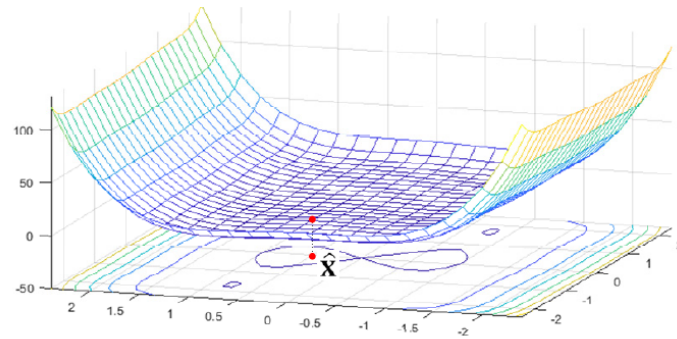


그림 2-23 최적해 탐색



■ 기계 학습이 해야 할 일을 식으로 정의하면,

$$J(\theta) \text{를 최소로 하는 최적해 } \hat{\theta} \text{을 찾아라. 즉, } \hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) \quad (2.50)$$

- θ 는 매개변수, $J(\theta)$ 는 목적함수

매개변수 공간의 탐색

■ 최적화 문제 해결

■ **날날탐색** exhaustive search 알고리즘

- 차원이 조금만 높아져도 적용 불가능
- 예) 4차원 Iris에서 각 차원을 1000구간으로 나눈다면 총 1000^4 개의 점을 평가해야 함

■ **무작위 탐색** 알고리즘

- 아무 전략이 없는 순진한 알고리즘

알고리즘 2-1 날날탐색 알고리즘

입력: 훈련집합 \mathbb{X} 와 \mathbb{Y}

출력: 최적해 $\hat{\theta}$

```

1 가능한 해를 모두 생성하여 집합  $S$ 에 저장한다.
2  $min$ 을 충분히 큰 값으로 초기화한다.
3 for ( $S$ 에 속하는 각 점  $\theta_{current}$ 에 대해)
4     if( $J(\theta_{current}) < min$ )  $min = J(\theta_{current})$ ,  $\theta_{best} = \theta_{current}$ 
5  $\hat{\theta} = \theta_{best}$ 
```

알고리즘 2-2 무작위 탐색 알고리즘

입력: 훈련집합 \mathbb{X} 와 \mathbb{Y}

출력: 최적해 $\hat{\theta}$

```

1  $min$ 을 충분히 큰 값으로 초기화한다.
2 repeat
3     무작위로 해를 하나 생성하고  $\theta_{current}$ 라 한다.
4     if( $J(\theta_{current}) < min$ )  $min = J(\theta_{current})$ ,  $\theta_{best} = \theta_{current}$ 
5 until(멈춤 조건)
6  $\hat{\theta} = \theta_{best}$ 
```

매개변수 공간의 탐색

- [알고리즘 2-3]은 기계 학습이 사용하는 전형적인 알고리즘
 - 라인 3에서는 목적함수가 작아지는 방향을 주로 미분으로 찾아냄

알고리즘 2-3 기계 학습이 사용하는 전형적인 탐색 알고리즘(1장의 [알고리즘 1-1]과 같음)

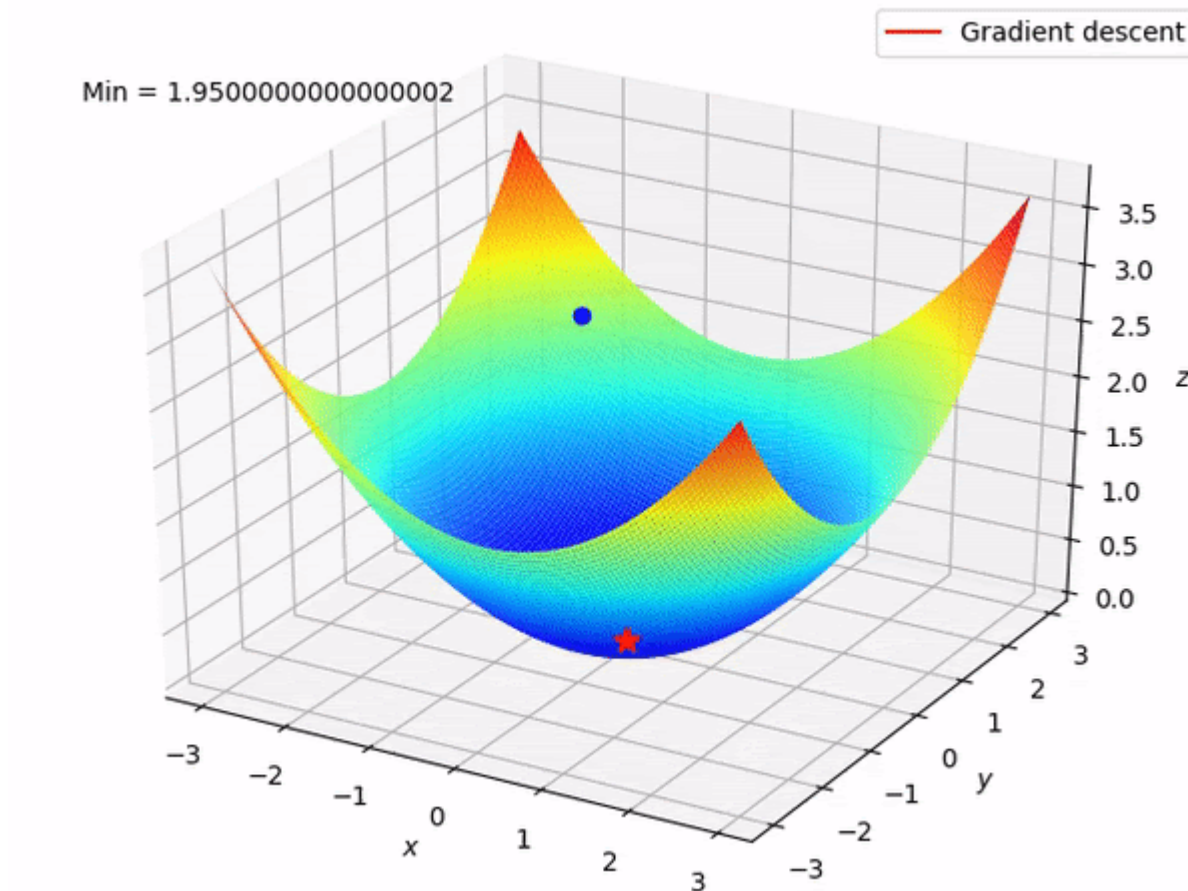
입력: 훈련집합 \mathbb{X} 와 \mathbb{Y}

출력: 최적해 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 을 설정한다.  
2  repeat  
3       $J(\theta)$ 가 작아지는 방향  $d\theta$ 를 구한다.  
4       $\theta = \theta + d\theta$   
5  until(멈춤 조건)  
6   $\hat{\theta} = \theta$ 
```

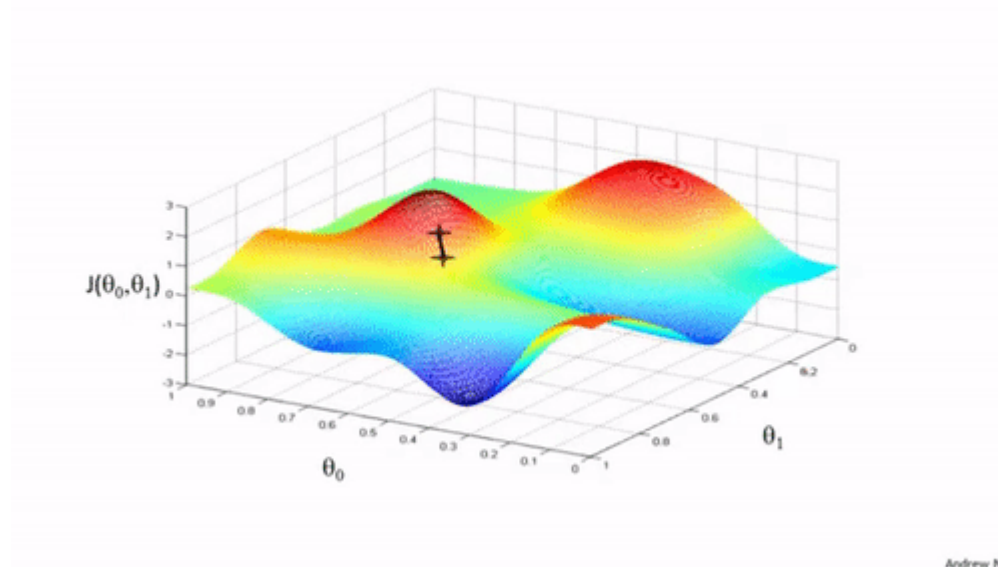
매개변수 공간의 탐색

■ 미분을 이용한 경사하강법



매개변수 공간의 탐색

■ 미분을 이용한 경사하강법



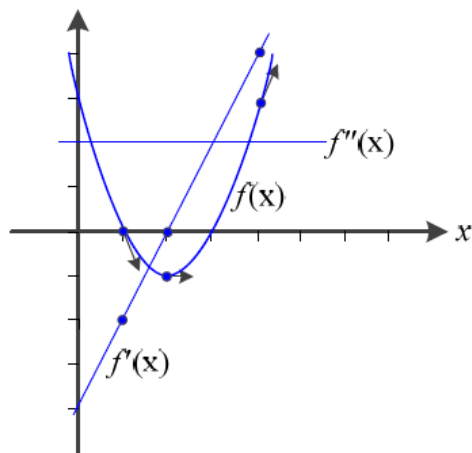
Andrew Ng

■ 미분에 의한 최적화

■ 미분의 정의

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}, \quad f''(x) = \lim_{\Delta x \rightarrow 0} \frac{f'(x + \Delta x) - f'(x)}{\Delta x} \quad (2.51)$$

- 1차 도함수 $f'(x)$ 는 함수의 기울기, 즉 값이 커지는 방향을 지시함
- 따라서 $-f'(x)$ 방향에 목적함수의 최저점이 존재
- [알고리즘 2-3]에서 $d\Theta$ 로 $-f'(x)$ 를 사용함 ← 경사 하강 알고리즘의 핵심 원리



$$y = f(x) = x^2 - 4x + 3$$

$$y' = f'(x) = 2x - 4$$

그림 2-24 간단한 미분 예제

미분

■ 편미분

- 변수가 여러 개인 함수의 미분
- 미분값이 이루는 벡터를 [그래디언트](#)라 부름

- 여러 가지 표기: $\nabla f, \frac{\partial f}{\partial \mathbf{x}}, \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)^T$

- 예)

$$\left. \begin{aligned} f(\mathbf{x}) &= f(x_1, x_2) = \left(4 - 2.1x_1^2 + \frac{x_1^4}{3} \right) x_1^2 + x_1 x_2 + (-4 + 4x_2^2) x_2^2 \\ \nabla f = f'(\mathbf{x}) &= \frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)^T = (2x_1^5 - 8.4x_1^3 + 8x_1 + x_2, 16x_2^3 - 8x_2 + x_1)^T \end{aligned} \right\} \quad (2.52)$$

■ 기계 학습에서 편미분

- 매개변수 집합 Θ 에 많은 변수가 있으므로 편미분을 많이 사용