

# Ch. 4

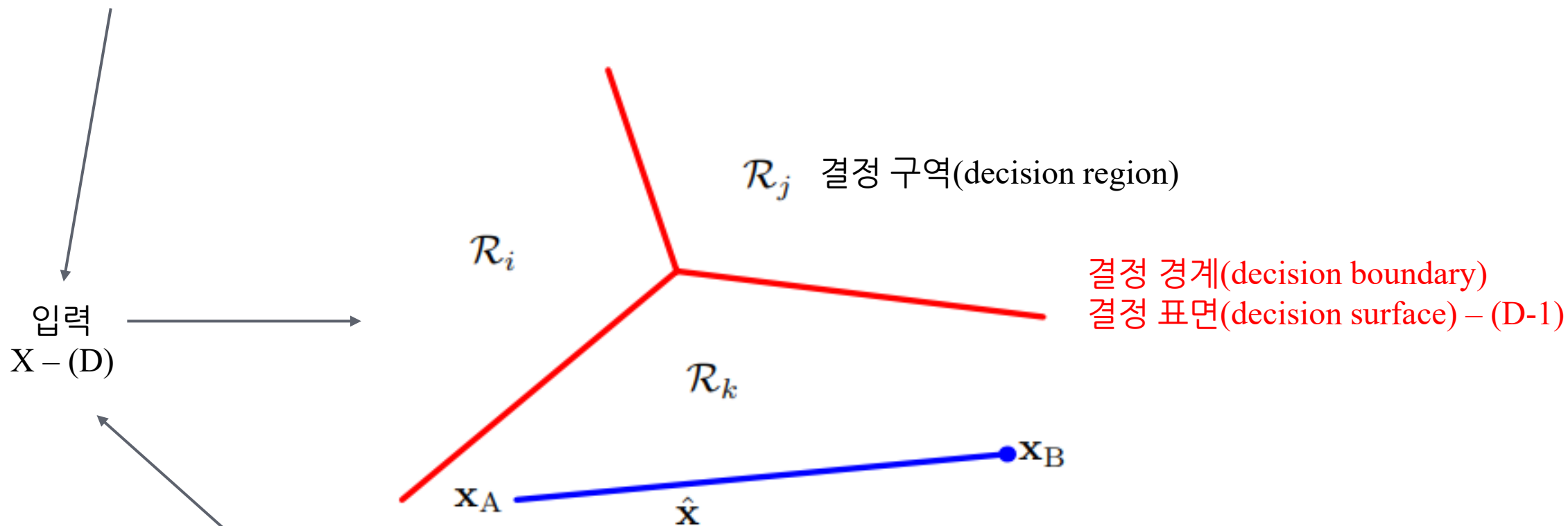
## Linear Models for Classification

조윤정

# Contents

<b>4</b>	<b>Linear Models for Classification</b>	<b>179</b>
4.1	Discriminant Functions . . . . .	181
4.1.1	Two classes . . . . .	181
4.1.2	Multiple classes . . . . .	182
4.1.3	Least squares for classification . . . . .	184
4.1.4	Fisher's linear discriminant . . . . .	186
4.1.5	Relation to least squares . . . . .	189
4.1.6	Fisher's discriminant for multiple classes . . . . .	191
4.1.7	The perceptron algorithm . . . . .	192
4.2	Probabilistic Generative Models . . . . .	196
4.2.1	Continuous inputs . . . . .	198
4.2.2	Maximum likelihood solution . . . . .	200
4.2.3	Discrete features . . . . .	202
4.2.4	Exponential family . . . . .	202
4.3	Probabilistic Discriminative Models . . . . .	203
4.3.1	Fixed basis functions . . . . .	204
4.3.2	Logistic regression . . . . .	205
4.3.3	Iterative reweighted least squares . . . . .	207
4.3.4	Multiclass logistic regression . . . . .	209
4.3.5	Probit regression . . . . .	210
4.3.6	Canonical link functions . . . . .	212
4.4	The Laplace Approximation . . . . .	213
4.4.1	Model comparison and BIC . . . . .	216
4.5	Bayesian Logistic Regression . . . . .	217
4.5.1	Laplace approximation . . . . .	217
4.5.2	Predictive distribution . . . . .	218

## 4 선형 분류 모델



데이터 집합  $X$  : 선형 분리 가능(linearly separable)한 집합

# 4 선형 분류 모델

One-hot-encoding

$$\mathbf{t} = (0, 1, 0, 0, 0)^T$$

$t_k$  : 클래스  $C_k$  일 확률

타깃변수  $t$   
해석

확률적 모델

- 1)  $t \in \{0, 1\}$   
 $t=1 \ C_1$   
 $t=0 \ C_2$
- 2)  $t : C_1$ 에 속할 확률

# 4 선형 분류 모델

1장 - 분류 문제 푸는 3가지 방법

- 1) 판별 함수 (discriminant function) 만들어 활용
- 2) 추론 단계에 조건부 확률 분포 모델 만들고 이를 활용(추론, 결정 분리)
  - 1) 직접 모델 - 조건부 확률을 매개변수 모델로 표현
  - 2) 생성적인 방법 - 베이저안 정리 이용

일반화된 선형 모델(generalized linear model)

$f$ 는 비선형이더라도 결정 경계면은  $x$ 에 대한 선형 함수  
→ 매개변수에 대해 선형적  $X$

$$y(\mathbf{x}) = \underline{f}(\mathbf{w}^T \mathbf{x} + w_0).$$

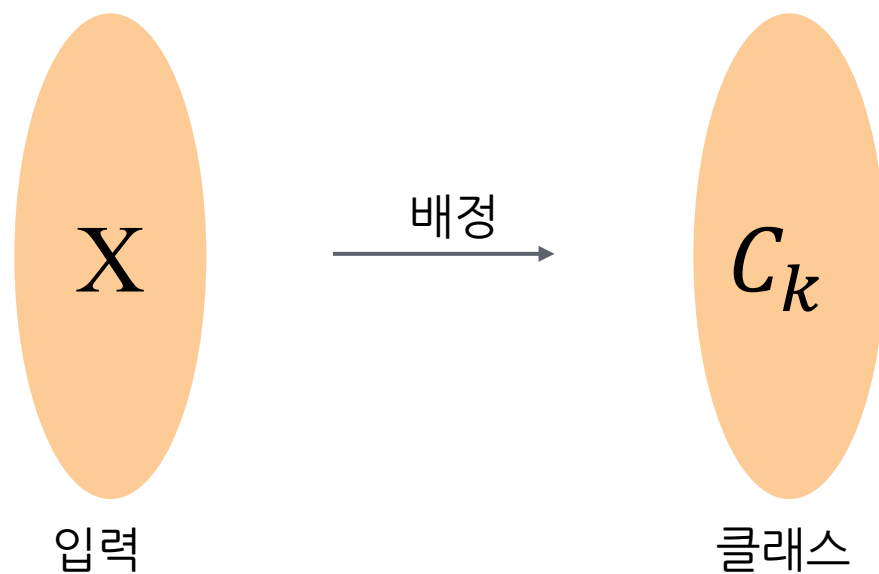
활성화 함수(activation function)

역함수

연결 함수(link function)

결정 경계면 :  $y(\mathbf{x}) = \text{constant} \longrightarrow \mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$

## 4.1 판별함수



선형 판별(linear discriminant) : 결정표면이 초평면

## 4.1.1 두 개의 클래스

선형 판별 함수 가장 단순하게 표현

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$\mathbf{w}$  : 가중 벡터

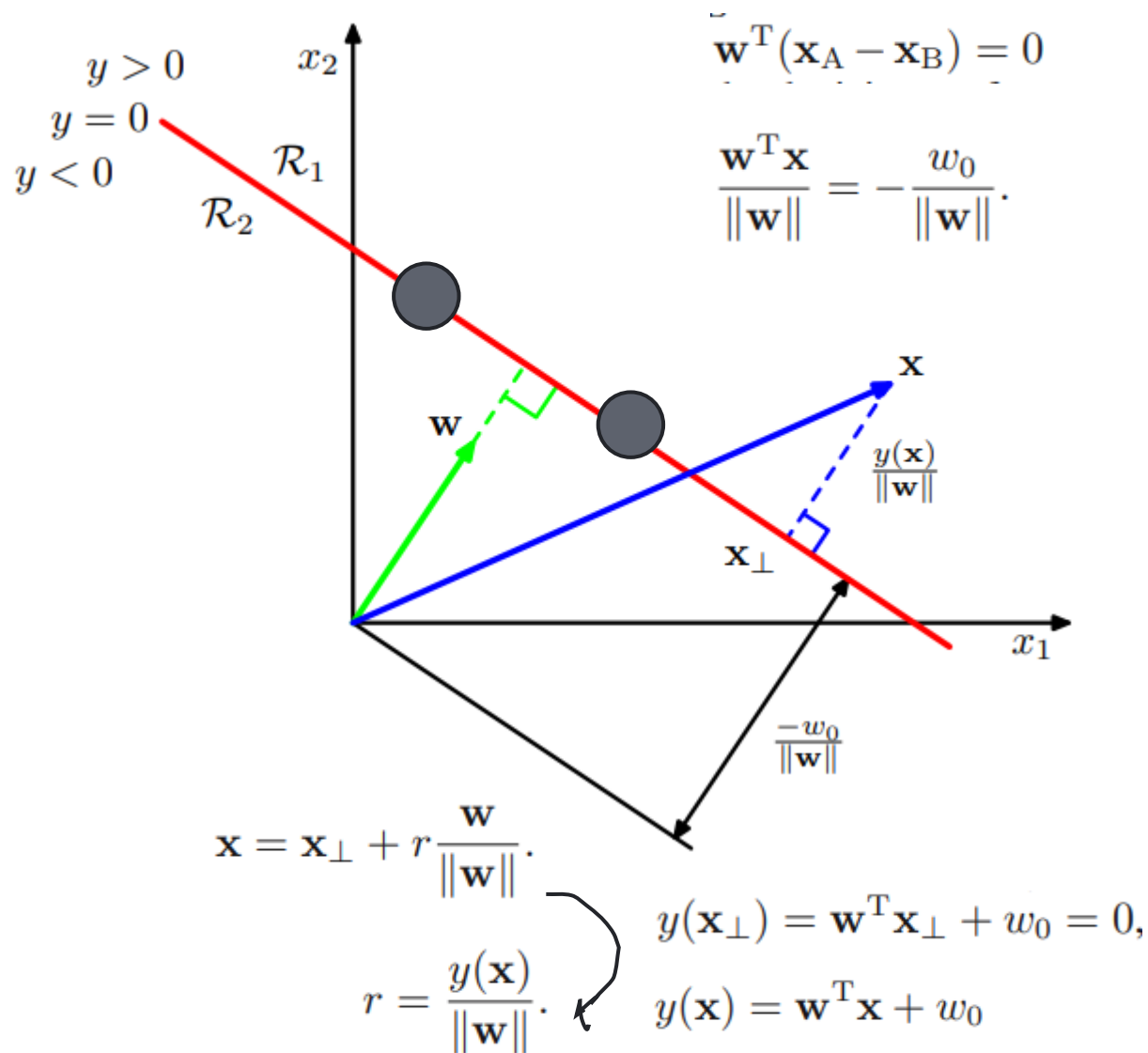
$w_0$  : 편향

threshold(임계값) : 편향의 음의 값

$y(\mathbf{x}) = 0$ , 결정경계

입력 차원 :  $D$

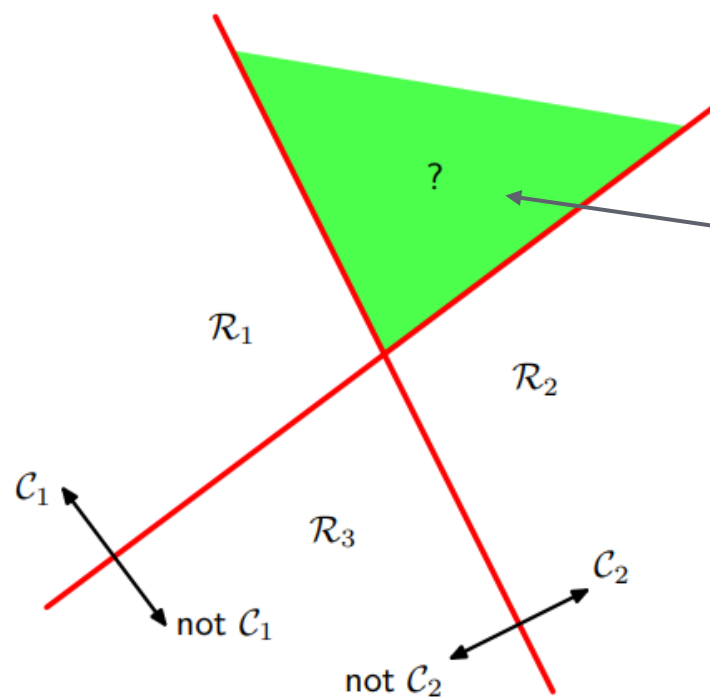
초평면 :  $D$ 차원 입력 공간 상의  $(D-1)$ 차원



## 4.1.2 다중 클래스

일대다 분류기

2클래스 판별 함수들  $\xrightarrow{(K-1)\text{개}}$   $K > 2$  K클래스 판별 함수



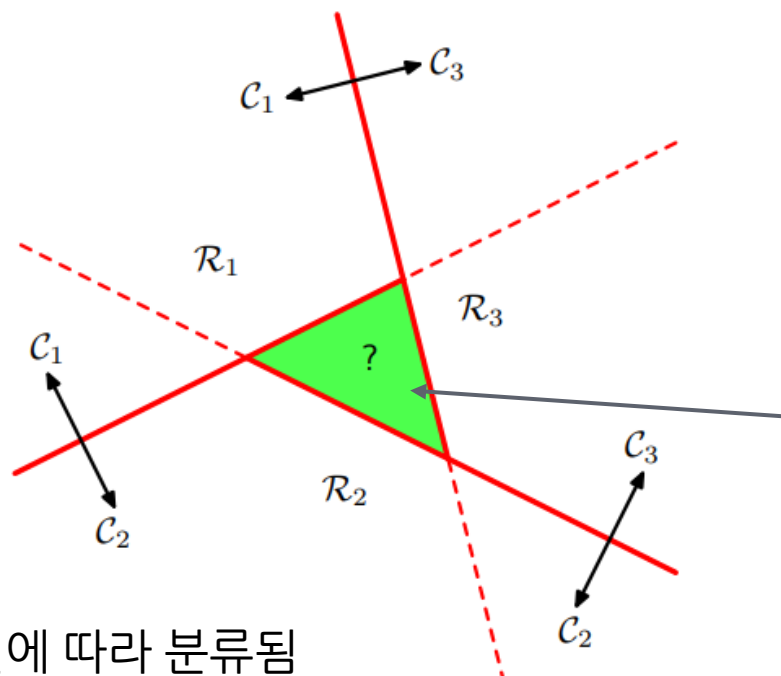
$\mathcal{C}_1$ 이자  $\mathcal{C}_2$  인 영역이 생김,  
분류해야하는데 중복,  
즉 불확실한 영역이 발생



## 4.1.2 다중 클래스

일대일 분류기

2클래스 판별 함수들  $\xrightarrow{K(K-1)/2 \text{ 개}}$   $K > 2$  K클래스 판별 함수



$c_1$ 이자  $c_2$  이자  $c_3$  인 영역이 생김,  
분류해야하는데 중복,  
즉 불확실한 영역이 발생

각 점은 판별 함수들의 다수결에 따라 분류됨

## 4.1.2 다중 클래스

문제 해결

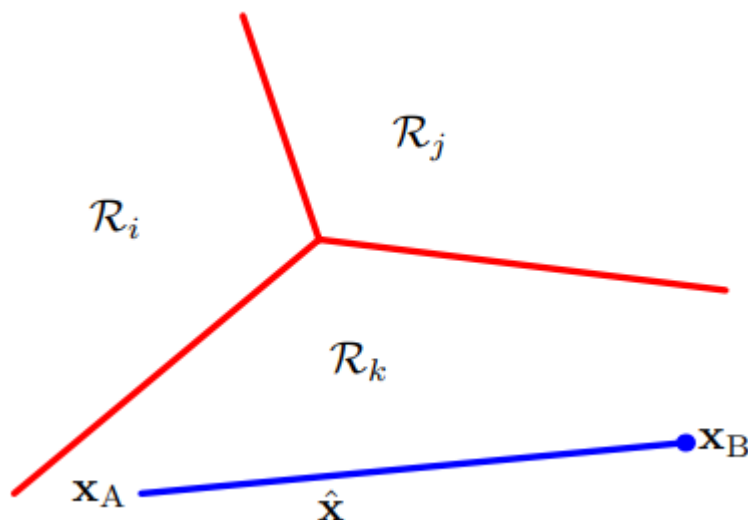
선형 함수들  $\xrightarrow{K\text{개}}$

$K > 2$   
하나의 K클래스 판별 함수

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ .

포인트  $\mathbf{x}$ 를  $C_k$ 에 배정



클래스  $C_k$ 와 클래스  $C_j$ 의 결정경계:  $y_k(\mathbf{x}) = y_j(\mathbf{x}) \longrightarrow$  초평면:  $(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$ .

## 4.1.2 다중 클래스

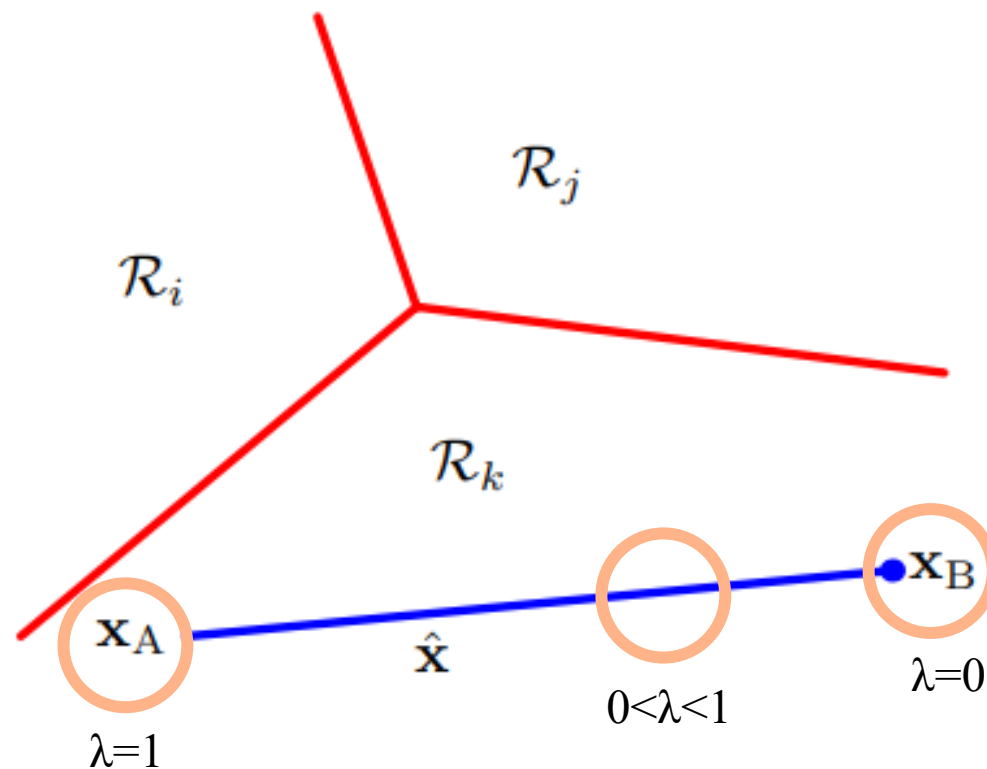
$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda)y_k(\mathbf{x}_B).$$

**볼록 집합 (Convex Set):**

정의: 집합이 볼록하다는 것은,  
그 집합 내의 임의의 두 점을 선택했을 때,  
그 두 점을 잇는 선분이 항상 그 집합 내에 존재하는 것을 의미함

$$\begin{array}{ccc} y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A) & \longrightarrow & y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}}) \\ y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B) & & \end{array}$$

$R_k$ 는 Convex



$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$$

**선형 보간(linear interpolation) :**  
두 점 사이 선분 위의 점을 찾는 법

## 4.1.3 ~ 4.1.7

선형 판별 함수의 ‘매개변수’ 학습

4.1.3 분류를 위한 **최소 제곱법**

4.1.4 **피셔의 선형 판별**

(4.1.5 두 클래스 문제에서 피셔는 최소 제곱법의 특별 케이스)

4.1.6 다중 클래스에 대한 피셔 판별식

4.1.7 **퍼셉트론 알고리즘**

## 4.1.3 분류를 위한 최소 제곱법

사용하는 이유 : 입력 벡터가 주어졌을 때 표적 벡터의 조건부 기댓값( $E[t|x]$ )의 근삿값을 구하는 방법이기

각각의 클래스  $C_k$  들을 선형 모델로 표현

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

↓ 벡터 표기 이용

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}}$$

제곱합의 오류 함수

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \right\}.$$

미분

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T}$$

이진 부호화 :  $E[t|x] \rightarrow$  사후 클래스 확률의 벡터

성능 ↓ (선형 모델의 제한적인 유연성  
← (0,1) 범위 밖 가질 수 있음)

최종 판별 함수

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} = \mathbf{T}^T \left( \widetilde{\mathbf{X}}^\dagger \right)^T \widetilde{\mathbf{x}}.$$

## 4.1.3 분류를 위한 최소 제곱법

$$\mathbf{a}^T \mathbf{t}_n + b = 0$$

어떤 훈련 집합의 표적 벡터들이 전부 아래 식을 만족한다면



$$\mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0.$$

어떤  $\mathbf{x}$ 값에 대한 모델 예측값 등 같은 제약 조건 (왼쪽 식)을 만족

최종 판별 함수

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} = \mathbf{T}^T (\widetilde{\mathbf{X}}^\dagger)^T \widetilde{\mathbf{x}}.$$

원 핫 인코딩을 한 K개의 클래스  $\mathbf{y}(\mathbf{x})$ 의 원소들을 전부 합하면 1

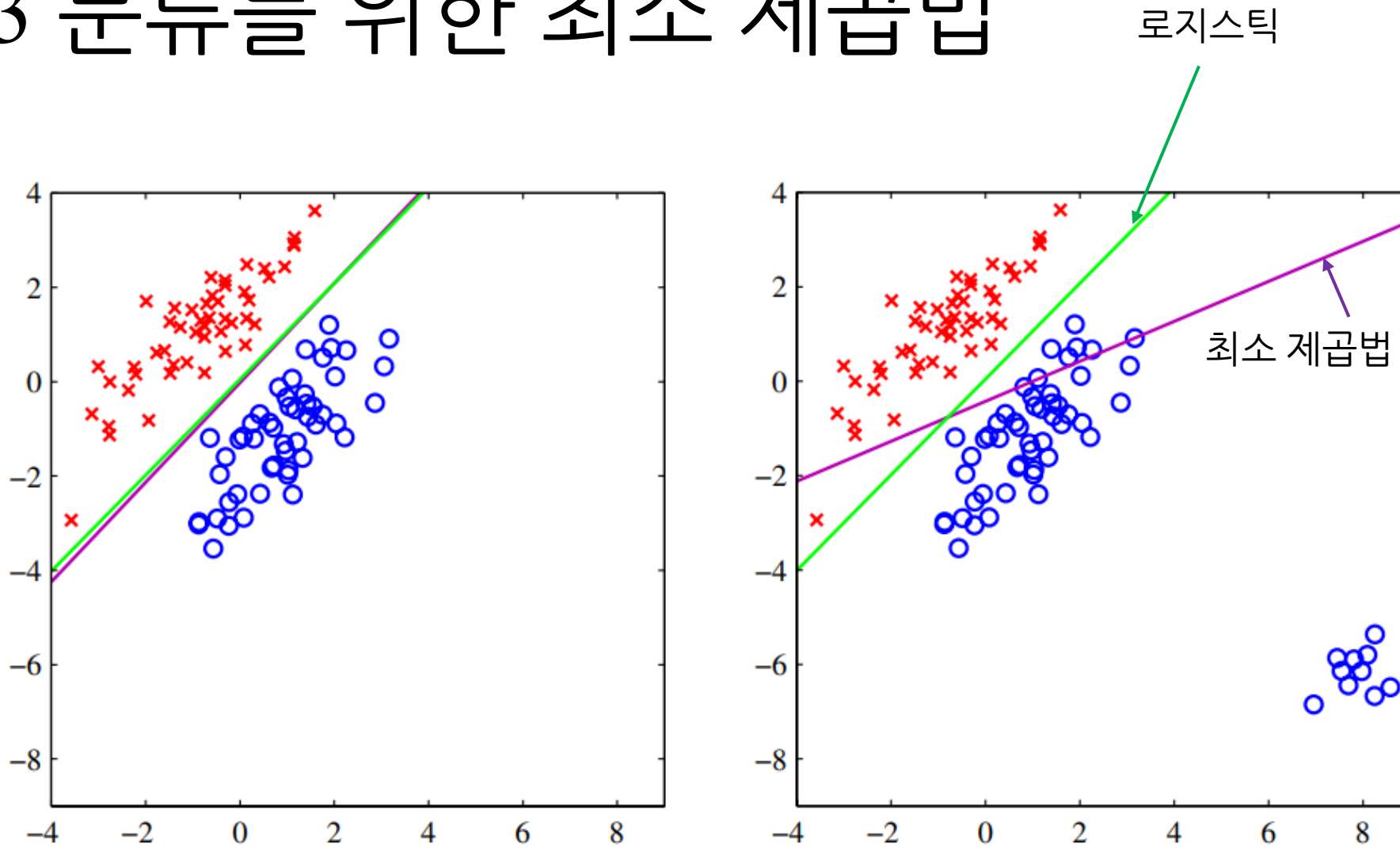
Ex) (0,1,0,0,0), K=5



Why? (0,1) 구간 사이에 값이 존재해야 한다는 '제약 조건' 이 없음

모델의 출력값을 확률로 해석

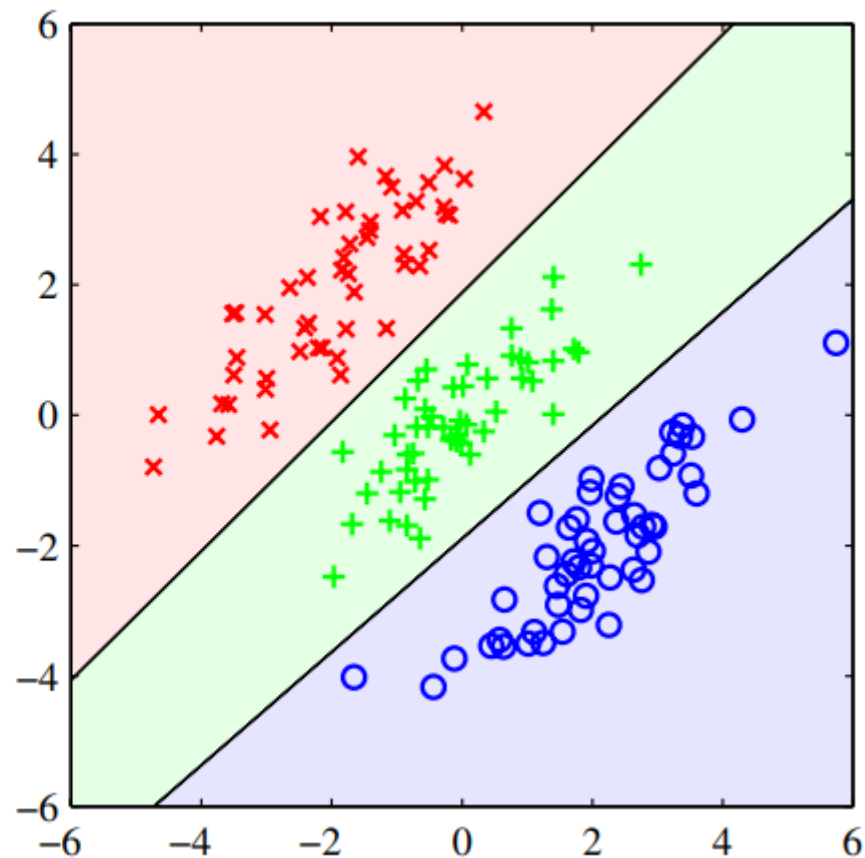
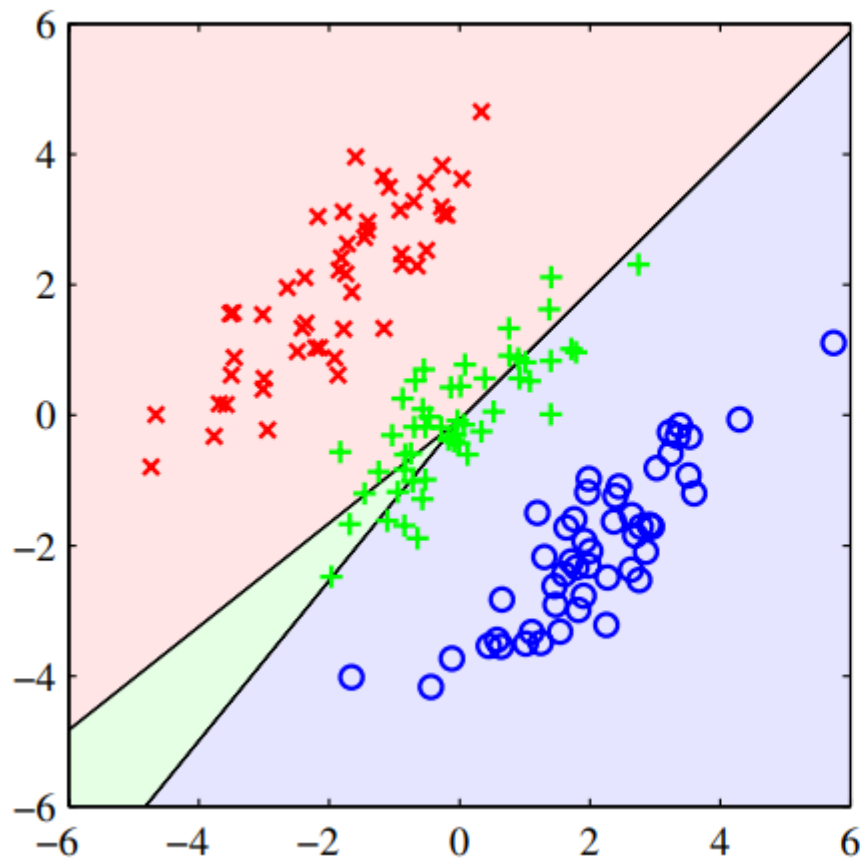
## 4.1.3 분류를 위한 최소 제곱법



이상값에 민감하다! → 강건성이 부족

## 4.1.3 분류를 위한 최소 제곱법

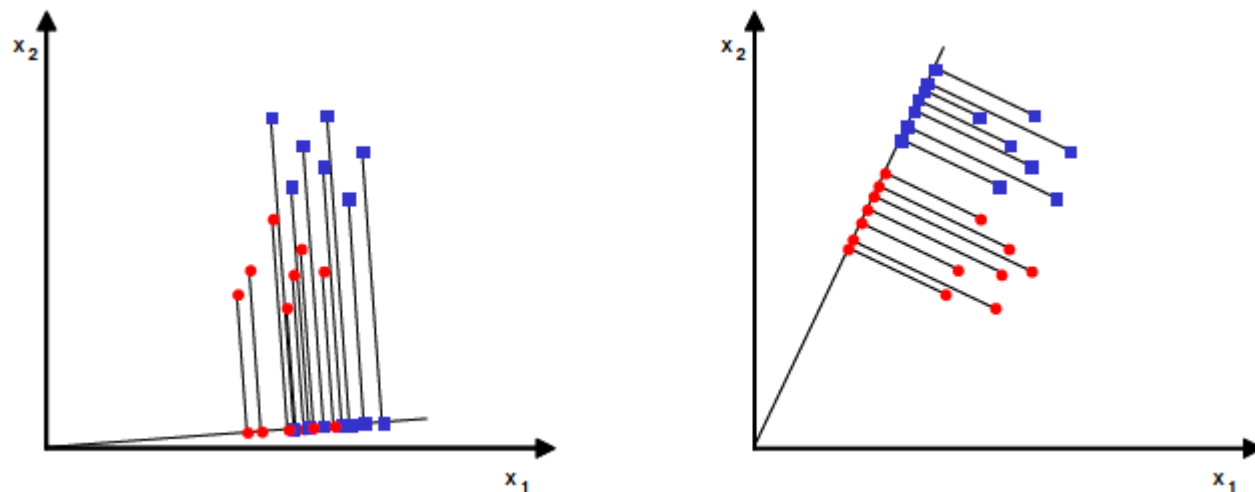
각 클래스마다 결정 경계 존재



최소 제곱법은 가우시안 조건부 분포를 가정하기에 이런 일이 발생함



## 4.1.4 피셔의 선형 판별



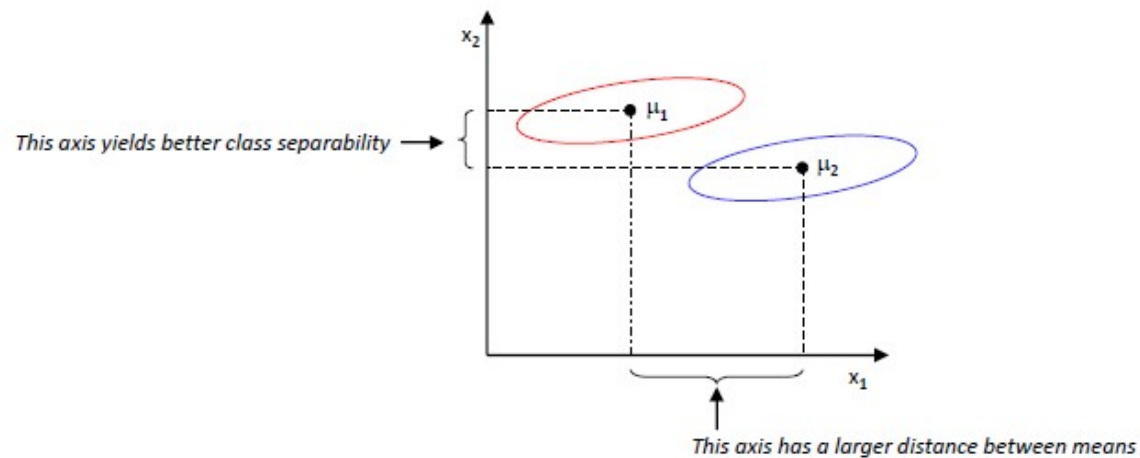
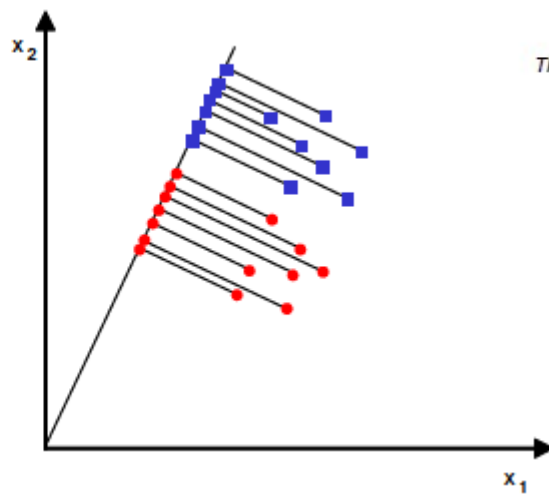
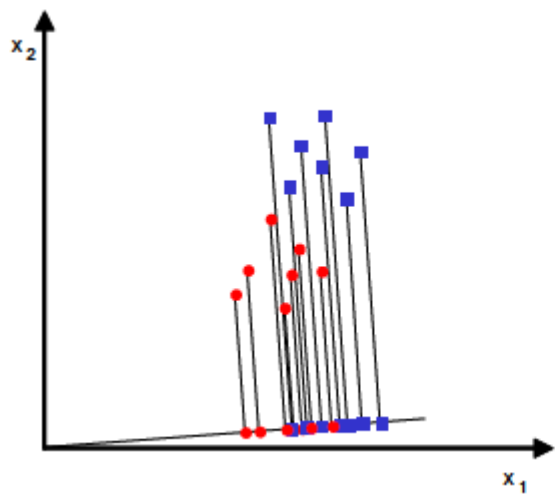
표본의 두 집단을 가장 잘 분리시키는 선에 정사영 하는 방법

1차원 투영 시 상당한 양의 정보 손실 &  
원래 잘 분리되었던 분리가 잘 이루어지지 X



가중 벡터  $w$ 의 성분들을 잘 조절하여  
클래스 간의 분리를 최대화 해야함.

## 4.1.4 피셔의 선형 판별



분리 척도 (measure of separation)

1. 정사영된 곳에서 각 집합의 평균 차이 : 좋지 X

각각 집단의 표본 분산을 고려하지 않았서 안 좋음

$x, y$  공간들의 평균 벡터

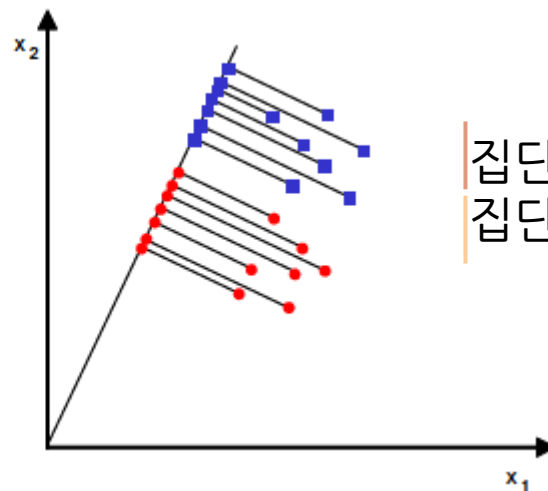
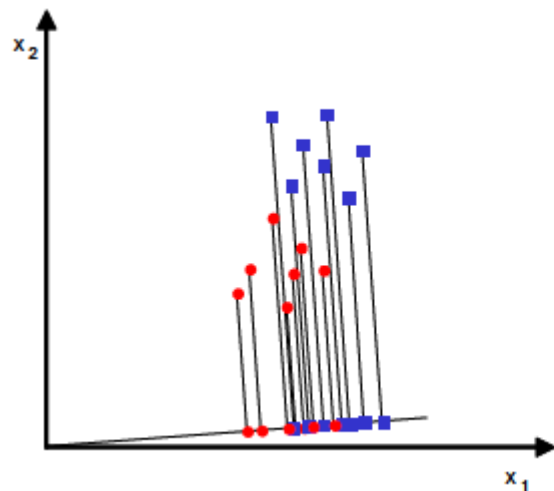
$$\mu_i = \frac{1}{N_i} \sum_{x \in \omega_i} x \text{ and } \tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in \omega_i} y = \frac{1}{N_i} \sum_{x \in \omega_i} w^T x = w^T \mu_i$$

정사영된 평균의 차

$$J(w) = |\tilde{\mu}_1 - \tilde{\mu}_2| = |w^T (\mu_1 - \mu_2)|$$

## 4.1.4 피셔의 선형 판별

클래스 내 분산  $s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$



집단(클래스) 간 분산 값 ↑  
집단(클래스) 내 분산 값 ↓

분리 척도 (measure of separation)

### 2. 피셔의 해법

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

투영된 클래스 평균 사이의 분리 정도 ↑  
각 클래스 내의 분산 ↓

클래스 간 중복 최소화

## 4.1.4 피셔의 선형 판별

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}.$$

$J(\mathbf{w})$ 는 위의 경우에서 극대화됨

$$y = \mathbf{w}^T \mathbf{x}. \quad \text{선형 함수 정의}$$

$$m_k = \mathbf{w}^T \mathbf{m}_k$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$$

$$\longrightarrow J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

$\mathbf{w}$ 에 대해 미분

스칼라만 보기

클래스 간 공분산 행렬(between class)

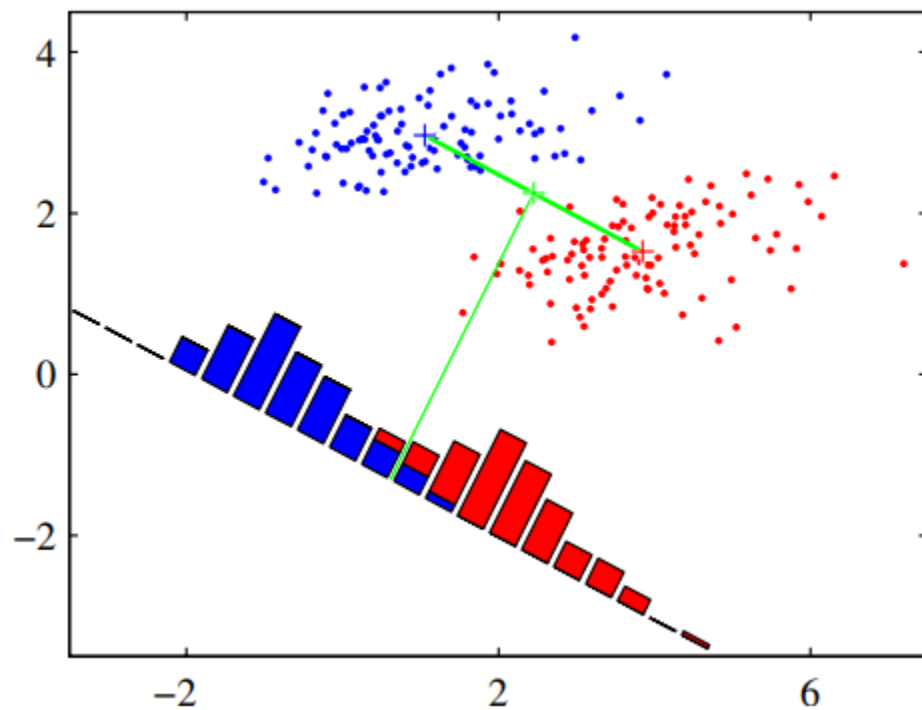
$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

클래스 내 공분산 행렬(within class)

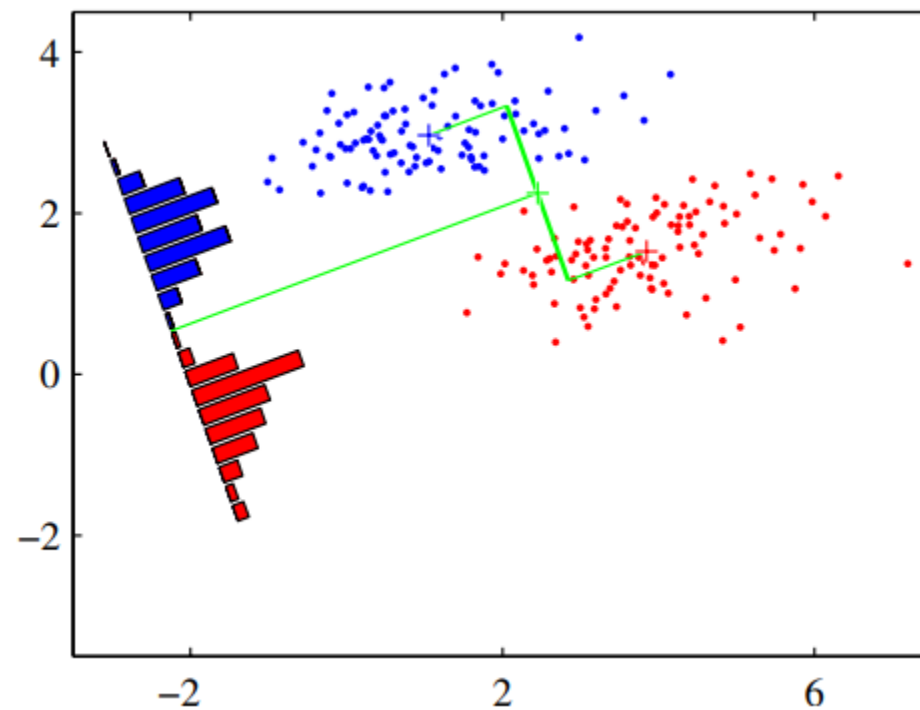
$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T.$$

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1).$$

## 4.1.4 피셔의 선형 판별

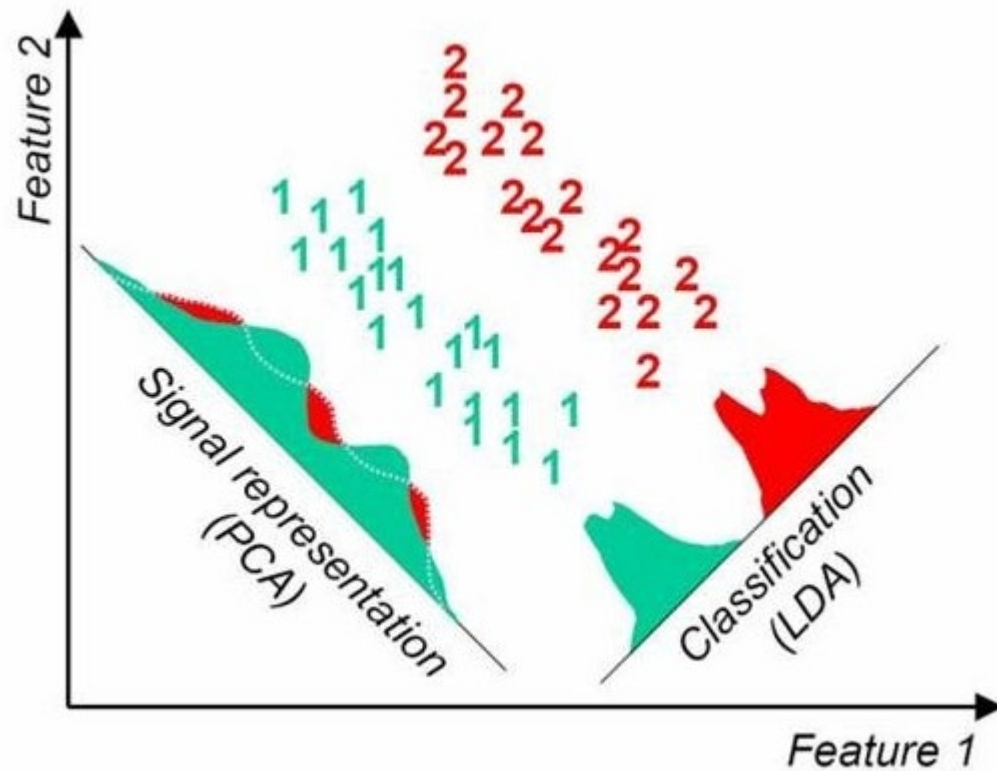


분산 고려 X



클래스 내 분산 고려

## 4.1.4 피셔의 선형 판별



## 4.1.5 최소 제곱법과의 관계

선형 판별 함수의 ‘매개변수’ 학습

출력 공간 상에서 클래스 간의 분리를  
최대화 하기 위해

모델의 예측값을 가능한 한 표적값  
에 가깝게 하기 위해

4.1.3 분류를 위한 최소 제곱법

4.1.4 피셔의 선형 판별

(4.1.5 두 클래스 문제에서 피셔는 최소 제곱법의 특별 케이스)

4.1.6 다중 클래스에 대한 피셔 판별식

4.1.7 퍼셉트론 알고리즘

피셔 기준은 최소 제곱법의 특별 케이스

## 4.1.5 최소 제곱법과의 관계

원 핫 인코딩을 표적값에 적용



다른 부호화 적용  
→ 최소 제곱 해는 곧 피셔 해

$C_1$ 의 표적값:  $N / N_1$

$C_2$ 의 표적값:  $-N / N_2$

$N$ : 전체 패턴들의 숫자

$N_1$ : 클래스  $C_1$ 에 있는 패턴들의 숫자

$N_2$ : 클래스  $C_2$ 에 있는 패턴들의 숫자



## 4.1.5 최소 제곱법과의 관계

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

$\mathcal{C}_1$  if  $y(\mathbf{x}) = \mathbf{w}^T(\mathbf{x} - \mathbf{m}) > 0$  and class  $\mathcal{C}_2$  otherwise.

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2.$$

↓  $w_0$  와  $\mathbf{w}$ 에 대해 미분, 각 미분값을 0으로 설정

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0$$

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0.$$

표적값 부호화를  $t_n$ 에 대해 적용

$$w_0 = -\mathbf{w}^T \mathbf{m}$$

$t_n$  식

$$\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} = 0$$

$\mathbf{m}$ 은 전체 데이터 집합의 평균

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2).$$

정리한 후  $t_n$ 에 대해 부호화 적용

$$\left( \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2)$$

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T.$$

## 4.1.6 다중 클래스에 대한 피셔 판별식

$$D > K$$

클래스 K개 ( $K > 2$ )

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}. \quad \text{특징값 벡터 } \mathbf{y}$$

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad \text{클래스 내 공분산 행렬}$$

전체 공분산 행렬

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T$$

$\mathbf{m}$ 은 데이터 집합 평균

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k$$

$$\left| \begin{array}{l} \mathbf{S}_k = \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \\ \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n \end{array} \right|$$

## 4.1.6 다중 클래스에 대한 피셔 판별식

$D > K$

$N = \sum_k N_k$  전체 데이터 포인트들의 개수

클래스 K개 ( $K > 2$ )

클래스 내 공분산 행렬  
 $\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$   
추가적인 행렬

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T.$$

클래스 간 공분산

원 x 공간상에 정의

$$\mathbf{S}_W = \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} (\mathbf{y}_n - \boldsymbol{\mu}_k)(\mathbf{y}_n - \boldsymbol{\mu}_k)^T$$

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{y}_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k.$$

비슷한 행렬을 투영된 D'차원의 y 공간상에 정의

## 4.1.6 다중 클래스에 대한 피셔 판별식

$D > K$

클래스 K개 ( $K > 2$ )

클래스 간 공분산  $\uparrow$   
클래스 내 공분산  $\downarrow$

$$J(\mathbf{W}) = \text{Tr} \{ \mathbf{S}_W^{-1} \mathbf{S}_B \} .$$

$\downarrow$  투영 행렬  $\mathbf{W}$ 에 대한 명시적인 함수

$$J(\mathbf{w}) = \text{Tr} \{ (\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W}^T) \} .$$

가중치 값은  $\mathbf{S}_W^{-1} \mathbf{S}_B$ 의  $D$ 개의 고유값  
에 해당하는 고유 벡터들에 의해 결정

## 4.1.6 다중 클래스에 대한 피셔 판별식

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$$

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T.$$

---

K개의 행렬

각각의 행렬은 두 벡터의 외적으로 계수(rank)가 1

→ (K-1)개의 행렬들만 독립적

최대 (K-1)개의 행렬 계수(=0이 아닌 고윳값) 가짐

여기서 고유 벡터들에 의해 (K-1) 차원의 부분 공간에 투영  
→ J(W)의 값을 바꾸지 않음, 이를 이용해서 (K-1)개  
보다 많은 선형 ‘특징’ 찾는 것은 불가능

## 4.1.7 퍼셉트론 알고리즘

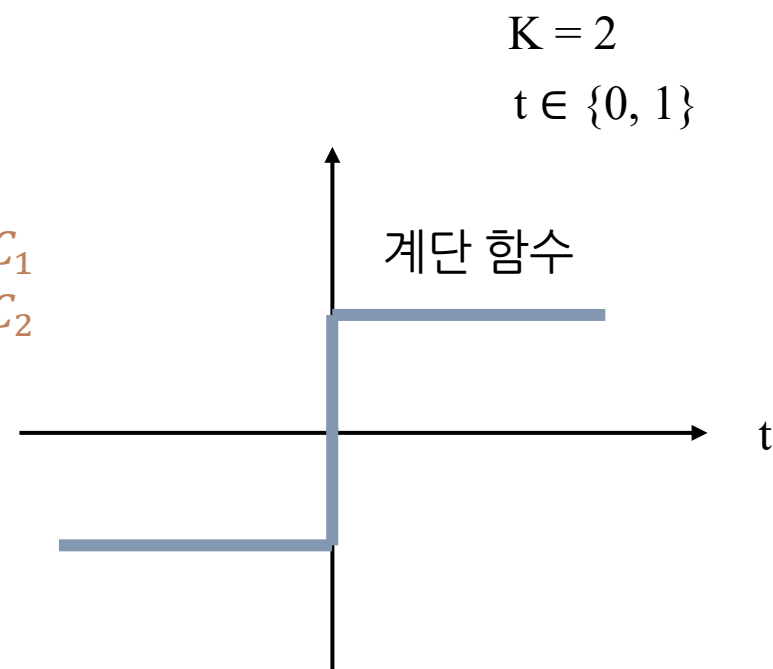
비선형 변환

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

입력 벡터  $\mathbf{x}$       특징 벡터

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases} \begin{matrix} C_1 \\ C_2 \end{matrix}$$

비선형 활성화 함수



매개변수  $\mathbf{w}$ 를 구하자

오류 함수를 최소화하자

~~오류 함수 : 오분류된 패턴의 총 숫자~~

학습 알고리즘 복잡(w에 대해 조각별 상수 함수)  
→ 불연속성, 기울기 이용 못 함.

오류 함수 : 퍼셉트론 기준(perceptron criterion)

## 4.1.7 퍼셉트론 알고리즘

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0. \end{cases}$  → 이 식을 만족하는 가중치 벡터  $\mathbf{w}$  찾기

→  $t \in \{-1, +1\}$

$$\mathbf{w}^T \phi(\mathbf{x}_n) t_n > 0.$$

오류 함수 : 퍼셉트론 기준(perceptron criterion)

올바르게 분류된 패턴 - 0 배정

오분류된 패턴  $\mathbf{x}_n$ 에 대해서  $\mathbf{w}^T \phi(\mathbf{x}_n) t_n$  최소화

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$

오분류된 패턴들의 전체 집합

## 4.1.6 다중 클래스에 대한 피셔 판별식

알고리즘의 단계에 대한 지표(정수)

확률적 경사 하강법 적용  $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$   
학습률

참고 : 훈련 중  $w$ 가 변함  $\rightarrow$  오분류가 되는 패턴 집합도 변함.

올바르게 분류  $\rightarrow w$  값 안 변함.

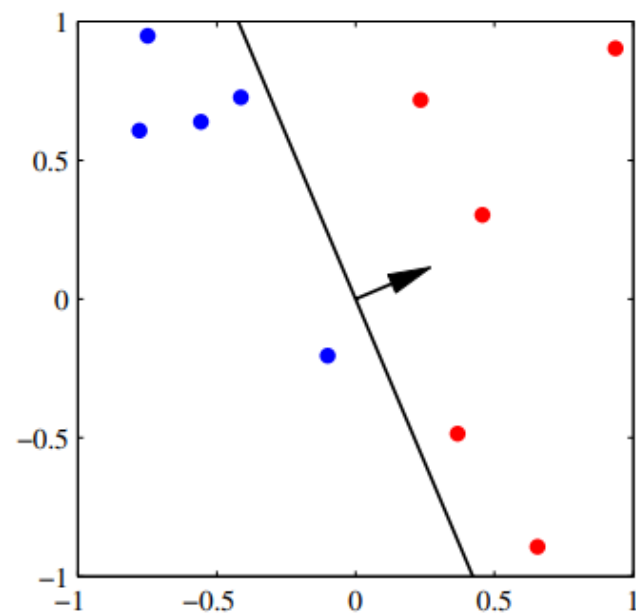
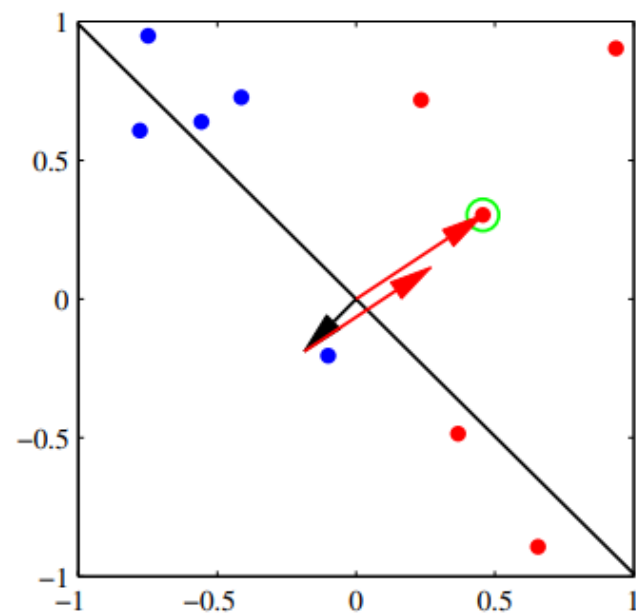
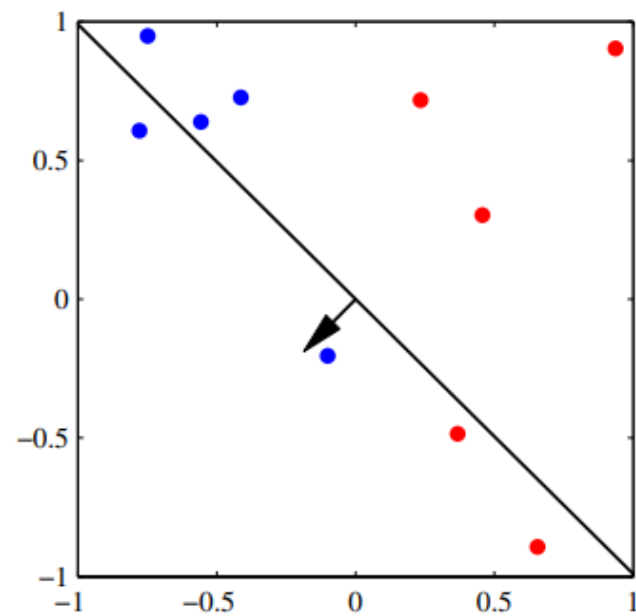
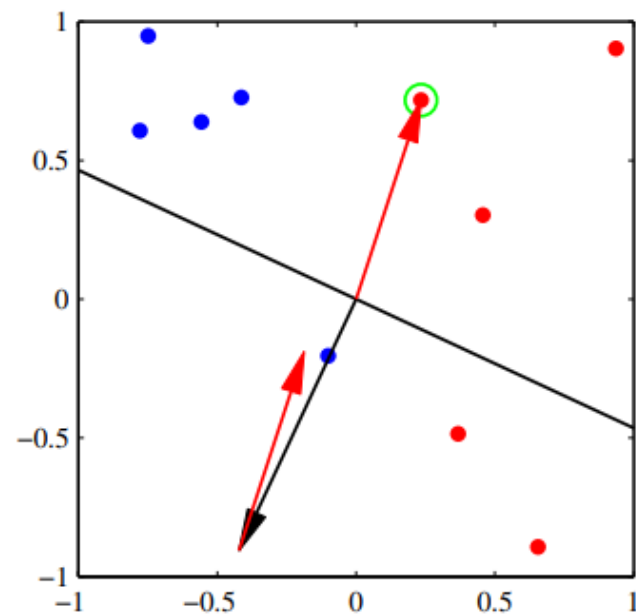
틀리게 분류  $\rightarrow C_1$  벡터  $\phi(\mathbf{x}_n) + w$ 의 예측치

$C_2$  벡터  $\phi(\mathbf{x}_n) - w$ 의 예측치

$$\left. \begin{array}{l} \eta = 1 \\ (\phi_n t_n)^2 > 0 \end{array} \right| -\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n$$

각 단계가 전체 오류 함수의 값을 줄인다고 보장 X





4.1.7 퍼셉트론 알고리즘

## 4.1.7 퍼셉트론 알고리즘

퍼셉트론 수렴 정리(perceptron convergence theorem)

If 문제가 정확한 해를 가지고 있기만 한다면  
(=훈련 집합이 선형적으로 분리가 가능하다면)

분리 가능 여부 모름  
서로 다른 여러 해 존재 가능성

정확한 해를 유한한 단계 안에 확실히 구할 수 있다.  
(단, 수렴을 위해 필요한 단계의 수는 알 수 없음)

시간 문제

한계점

학습 알고리즘 자체가 어려움  
확률적인 출력값 내지 않음  
K>2 클래스 문제에 대한 일반화X  
고정된 기저 함수들의 선형 결합 → 가장 큰 한계점 발생

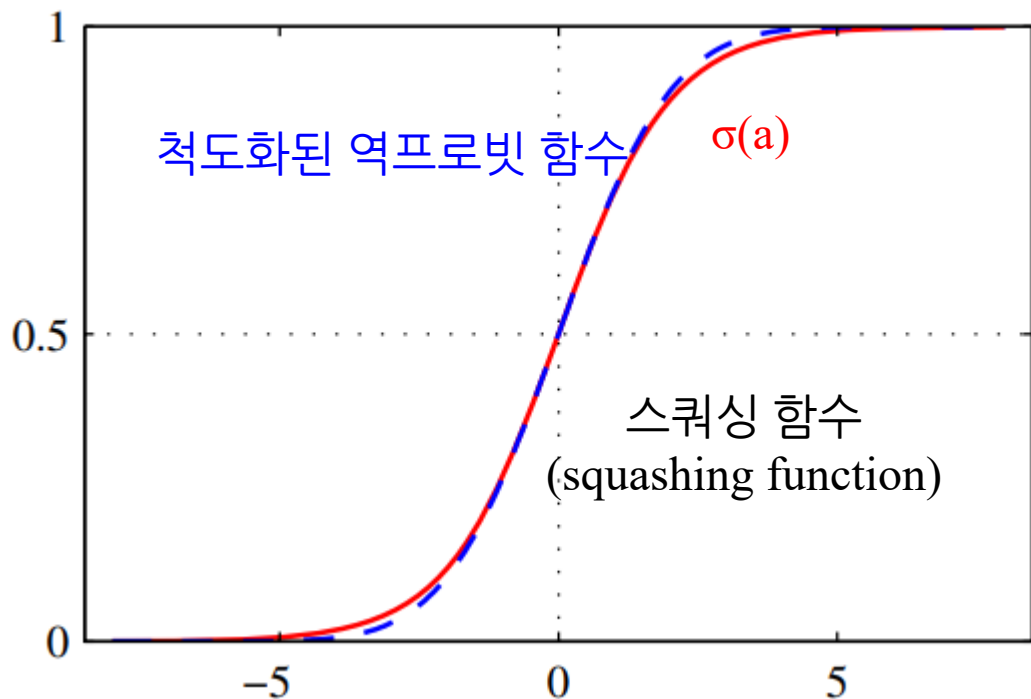
매개변수 초기화 방법  
데이터 포인트들 입력 순서

## 4.2 확률적 생성 모델

$p(\mathbf{x}|\mathcal{C}_k)$ : 클래스별 조건부 밀도

$p(\mathcal{C}_k)$ : 클래스 사전 분포

→  $p(\mathcal{C}_k|\mathbf{x})$ : 사후 확률



$$\sigma(-a) = 1 - \sigma(a) \quad \text{역 : } a = \ln \left( \frac{\sigma}{1 - \sigma} \right)$$

Sigmoid 성질                      로짓 함수

$$\ln [p(\mathcal{C}_1|\mathbf{x})/p(\mathcal{C}_2|\mathbf{x})] \quad \text{log odds}$$

클래스 2개

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

$$\underline{\sigma(a)} = \frac{1}{1 + \exp(-a)}$$

$\sigma(a)$ : Logistic Sigmoid

## 4.2 확률적 생성 모델

클래스 K개 ( $K > 2$ )

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

→ 정규화된 지수 함수(normalized exponential function) = softmax function

logistic sigmoid를 여러 클래스에 대해 일반화한 형태

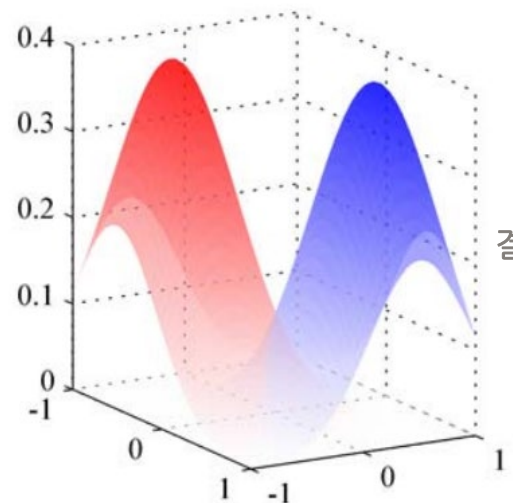
$$a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

$$j \neq k, \quad a_k \gg a_j$$

$$p(\mathcal{C}_k|\mathbf{x}) \simeq 1, \text{ and } p(\mathcal{C}_j|\mathbf{x}) \simeq 0.$$

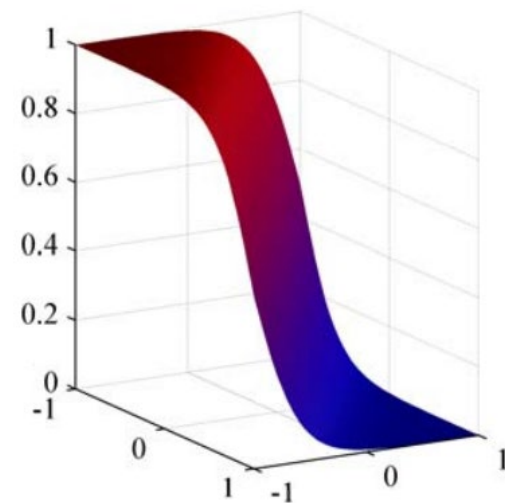
## 4.2.1 연속 입력

클래스별 조건부 밀도 → 가우시안  
모든 클래스가 같은 공분산 행렬 공유



결정경계는 사후 확률  
 $p(C_k|\mathbf{x})$ 이 상수  
 $\mathbf{x}$ 의 선형 함수

두 클래스들의 클래스 조건부 밀도



$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) \right\}.$$

$p(C_1|\mathbf{x})$ : 사후 확률

$p(C_2|\mathbf{x})$ : 사후 확률

$p(C_2|\mathbf{x}) = 1 - p(C_1|\mathbf{x})$

$p(C_k)$ 는  $w_0$  통해서만 연관 → 결정 경계 평행 이동

클래스 2개

$$p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

정의

$$\begin{cases} \mathbf{w} &= \Sigma^{-1}(\mu_1 - \mu_2) \\ w_0 &= -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)} \end{cases}$$

## 4.2.1 연속 입력

클래스 K개 ( $K > 2$ )

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

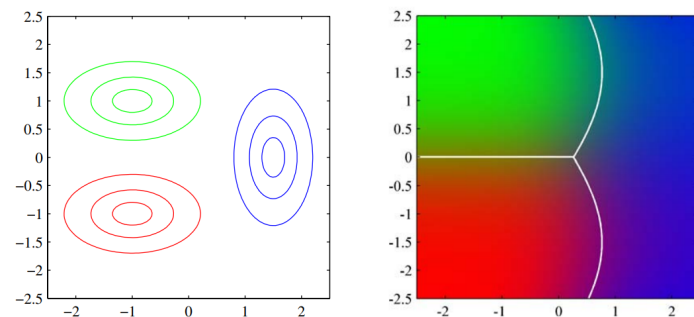
정의

$$\left\{ \begin{array}{l} \mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k \\ w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k). \end{array} \right.$$

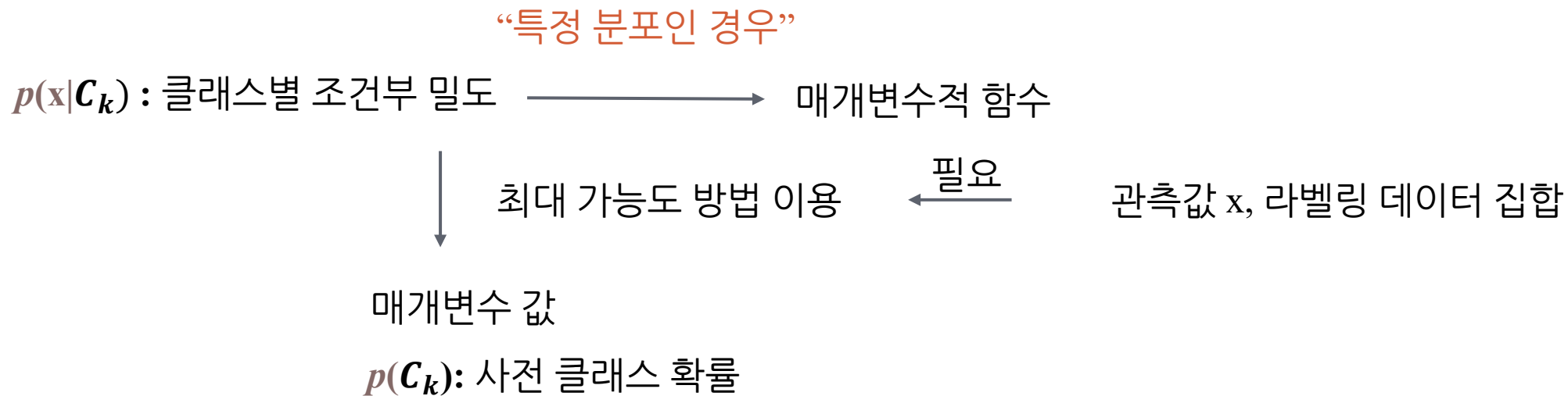
클래스별 조건부 밀도  $\rightarrow$  가우시안  
~~모든 클래스가 같은 공분산 행렬 공유~~

$p(\mathbf{x}|\mathcal{C}_k)$ : 클래스별 조건부 밀도  $\longrightarrow$  각각 공분산 행렬  $\Sigma_k$  가진다

이차항 안 사라짐  $\rightarrow$   $\mathbf{x}$ 의 이차 함수 : 이차 판별식(quadratic discriminant)



## 4.2.2 최대 가능도 해

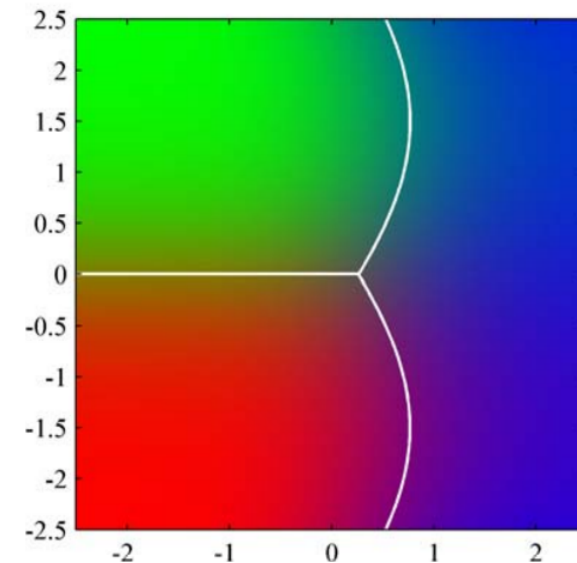
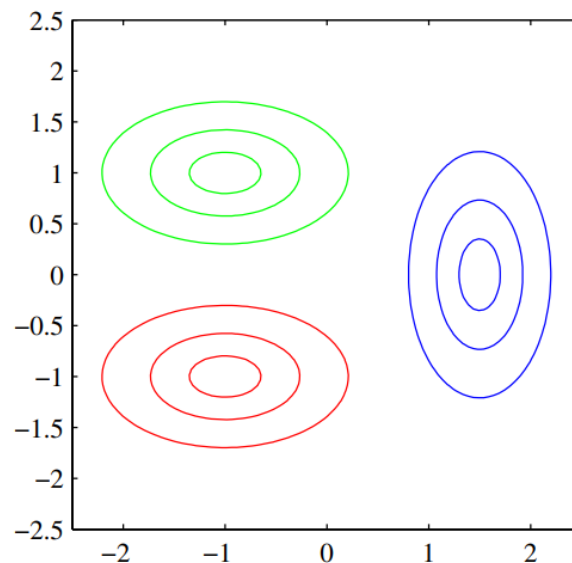


## 4.2.2 최대 가능도 해

클래스별 조건부 밀도  $\rightarrow$  가우시안  
모든 클래스가 같은 공분산 행렬 공유

클래스 2개

$p(\mathcal{C}_1): \pi, p(\mathcal{C}_2): 1 - \pi$



$$\mathcal{C}_1 \ t_n = 0 \quad p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n|\mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}).$$

$$\mathcal{C}_2 \ t_n = 1 \quad p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n|\mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

$$\text{가능도 함수 } p(\mathbf{t}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

$$\mathbf{t} = (t_1, \dots, t_N)^T$$



## 4.2.2 최대 가능도 해

$C_1$  데이터 포인트 수 :  $N_1$   
 $C_2$  데이터 포인트 수 :  $N_2$

클래스 2개

$\pi$ 에 대해 종속적인 로그 가능도 함수 항

$$\sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi)\}$$

↓  $\pi$ 에 대해 미분, 각 미분값을 0으로 설정

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

$\mu$ 들에 대해 같은 방식으로 해서 구하기

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

$\mu$ 에 대해 종속적인 로그 가능도 함수 항

$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \mu_1)^T \Sigma^{-1} (\mathbf{x}_n - \mu_1) + \text{const.}$$

## 4.2.2 최대 가능도 해

공유된 공분산 행렬 $\Sigma$ 에 대한 최대 가능도 해 고려

$$\begin{aligned} & -\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \\ & -\frac{1}{2} \sum_{n=1}^N (1 - t_n) \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \\ & = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr} \{ \Sigma^{-1} \mathbf{S} \} \end{aligned}$$

$$\left| \begin{aligned} \mathbf{S} &= \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2 \\ \mathbf{S}_1 &= \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \\ \mathbf{S}_2 &= \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T. \end{aligned} \right.$$

정의 이용

$$\Sigma = \mathbf{S}$$

두 클래스 쌍의 해당하는 공분산 행렬들의 가중 평균

클래스 K개 ( $K > 2$ )

확장 가능, 하지만 이상점이 포함되어 있는 경우는 강건X

## 4.2.3 이산 특징

$$\mathbf{x}_i \in \{0, 1\}$$

나이브 베이즈 해 가정

특징의 수 제한을 위한  
제한된 표현 방법 필요

각각의 값들이 클래스  $\mathcal{C}_k$  에 대해 조  
건부일 때 서로 독립적으로 취급

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \quad \leftarrow D\text{개의 독립변수 가짐}$$

$$a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(\mathcal{C}_k) \quad \longrightarrow \quad x_i \text{에 대한 선형 함수}$$

입력 :  $D$

일반적인 분포 :  $2^D$

( $D$ 차원의 이진 입력을 고려할 때  
가능한 모든 조합에 대한 클래스 확률)

합산 제약 조건  $\rightarrow$  이 분포는  $2^D - 1$ 개의 독립 변수 가짐

클래스 2개

로지스틱 시그모이드  
함수 고려 가능

## 4.2.4 지수족

각각의 클래스들이 각자의 매개변수 벡터  $\lambda_k$  를 가지지만, 척도 매개변수  $s$  는 공유한다고 가정

모든 지수족 분포엔 고유 파라미터가 있음, 그래서 클래스별로 달라야 해서  $\lambda_k$  이용, 모든 클래스가 동일한 스케일 파라미터 가정

클래스  $K = 2$  로지스틱 시그모이드 함수

클래스  $K$  개 ( $K \geq 2$ ) 소프트맥스 함수

→ 더 일반화

$p(\mathbf{x}|\lambda_k) = h(\mathbf{x})g(\lambda_k) \exp \{ \lambda_k^T \mathbf{u}(\mathbf{x}) \}$  지수족 성질을 이용하여  $x$ 에 대한 분포 작성

$\mathbf{u}(\mathbf{x}) = \mathbf{x}$ 인 부분 집합들에 대해서만 고려

$$p(\mathbf{x}|\lambda_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right) g(\lambda_k) \exp \left\{ \frac{1}{s} \lambda_k^T \mathbf{x} \right\}$$

$$p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$$

← 척도 매개변수  $s$  도입

## 4.2.4 지수족

클래스  $K = 2$   $\longrightarrow$   $a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$  식 4.58 대입

$$a(\mathbf{x}) = (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(\mathcal{C}_1) - \ln p(\mathcal{C}_2).$$

사후 클래스 확률이 선형 함수  $a(\mathbf{x})$ 에 대한 로지스틱 시그모이드 함수로 주어짐

클래스  $K$ 개 ( $K \geq 2$ )  $\longrightarrow$   $a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ . 식 4.63 대입

$$a_k(\mathbf{x}) = \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(\mathcal{C}_k)$$

$\mathbf{x}$ 에 대한 선형 함수

## 4.3 확률적 판별 모델

생성적 모델링

4.2 내용

$p(\mathbf{x}|\mathbf{C}_k)$ : 클래스별 조건부 밀도

$p(\mathbf{C}_k)$ : 사전 클래스 확률



$p(\mathbf{C}_k|\mathbf{x})$ : 사후 확률 결정(간접적으로 파라미터 결정)

판별적 모델링

<반복 재가중 최소 제곱법 = IRLS 알고리즘>

<iterative reweighted least squares>

일반화된 선형 모델 함수 형태를 명시적으로 사용

최대 가능도



직접적으로 매개변수 값 얻기  
(생성적 모델보에 비해 구해야 할 적응  
매개변수 숫자가 적음.)

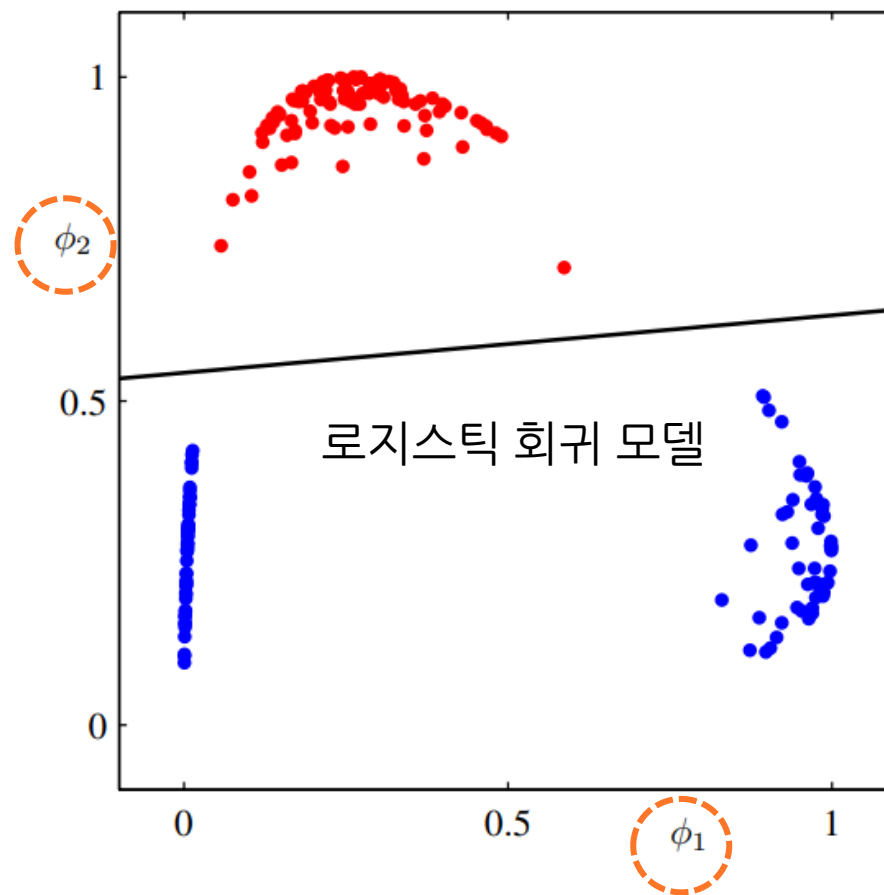
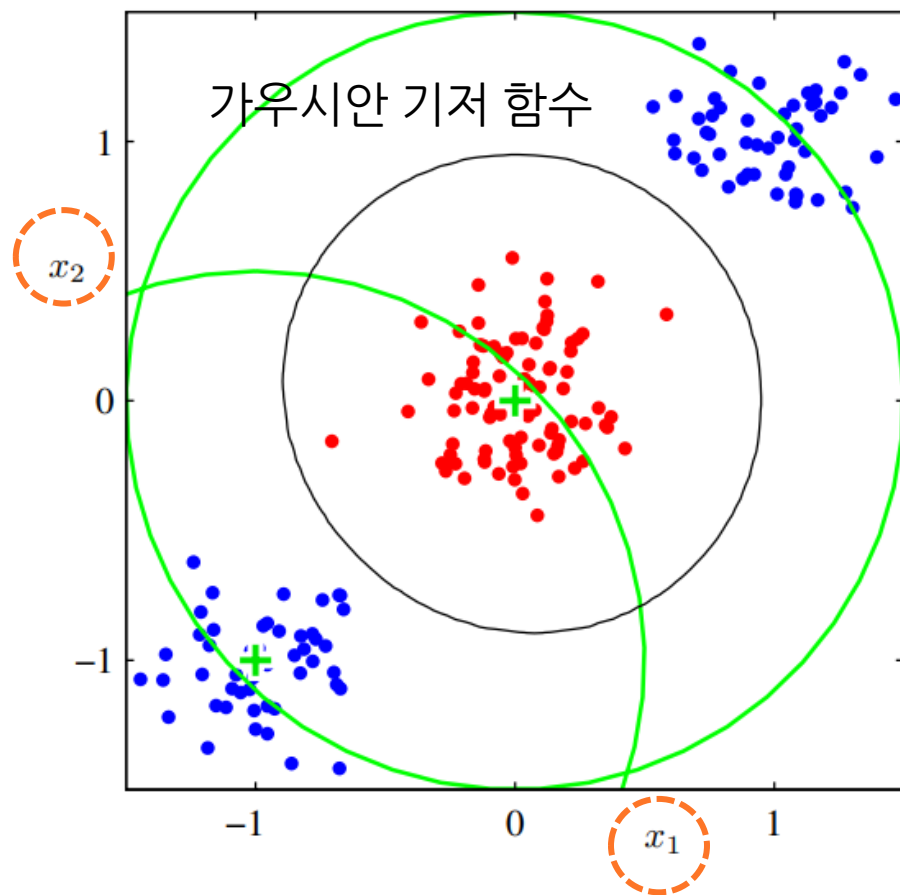


$p(\mathbf{C}_k|\mathbf{x})$  정의

MLE 활용

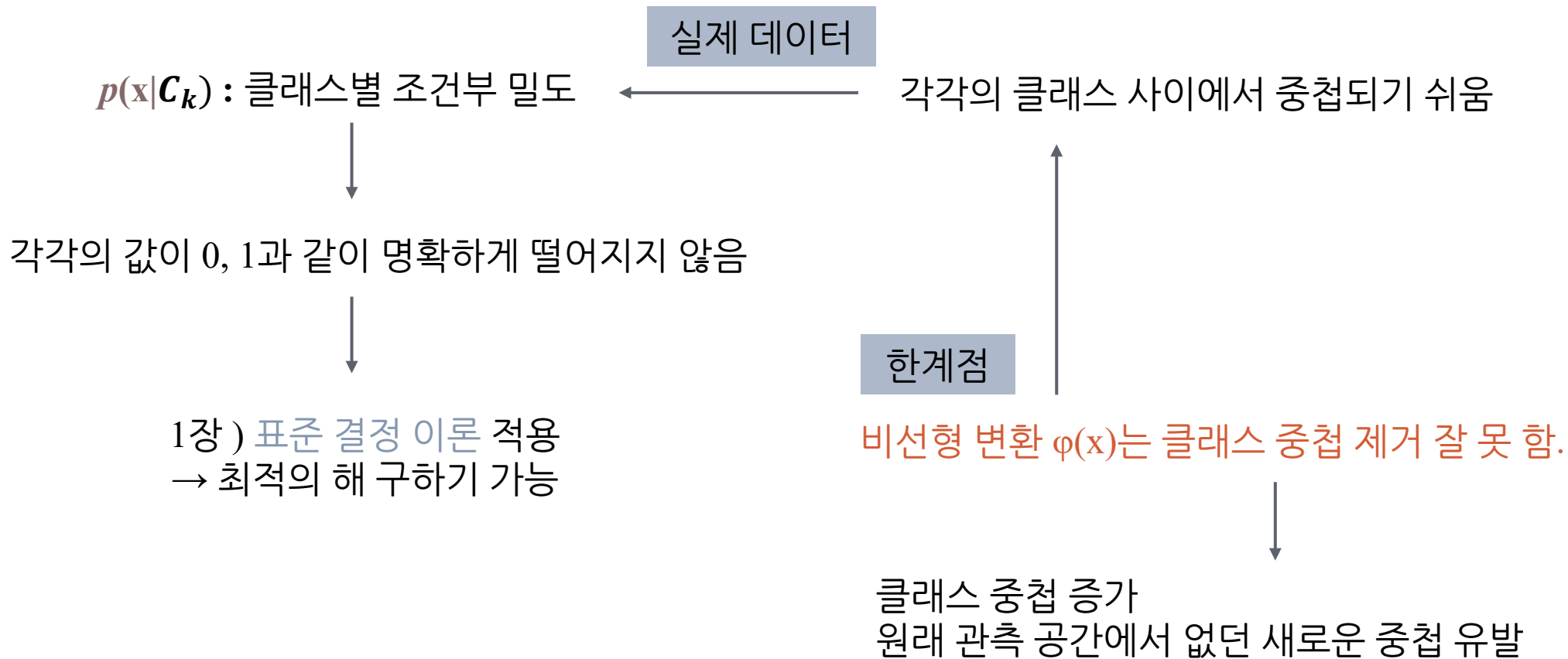
생성적 모델링이 잘 안 맞을 때 쓰면 좋은 결과 얻기 가능

## 4.3.1 고정된 기저 함수



입력 벡터를  $\phi_x$  로 변환하는 상태 이용할 것

## 4.3.1 고정된 기저 함수



그럼에도 불구하고 적절한 비선형성 선택 시 사후 확률 모델링 과정 쉬워짐



## 4.3.1 고정된 기저 함수

### 한계점

기저 함수 자체가 데이터에 대해 적응되도록 만들면 해결 가능

1. 비선형성과 복잡한 패턴의 제한 (데이터가 복잡한 비선형 패턴 잘 포착X)
2. 기저 함수 선택의 어려움 (데이터에 대한 깊은 이해도 필요)
3. 고차원 데이터에서의 한계 (차원의 저주, 효율성 저하)
4. 유연성의 부족 (적응성 결여, 기저 함수 동적 조절 불가)
5. 과적합 문제 (기저 함수 너무 복잡하다면)

## 4.3.2 로지스틱 회귀

클래스  $K = 2$  일반화된 선형 모델

4.2 내용

클래스  $K = 2$

$C_1$ 에 대한 사후 확률 :  
특정 벡터의 선형 함수에 대한  
로지스틱 시그모이드 함수

M차원 특정 공간  $\phi$

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad \text{“로지스틱 회귀”}$$

“가우시안 클래스 조건부 밀도”

$$p(C_2|\phi) = 1 - p(C_1|\phi)$$

매개변수 : M개

매개변수(평균값) : 2M개

매개변수(공분산 행렬) :  $M(M+5)/2 + 1$ 개

매개변수 개수 M이 큰 경우에는 로지스틱 회귀 모델을 다루는 것 이 더 유리할 수 있음

## 4.3.2 로지스틱 회귀

로지스틱 회귀 모델 매개변수 구하기 ← 최대 가능도 이용

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma).$$

a data set  $\{\phi_n, t_n\}$

$$\phi_n = \phi(\mathbf{x}_n)$$

$$t_n \in \{0, 1\}, n = 1, \dots, N$$

가능도 함수

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

음의 로그값



교차 엔트로피  
(cross entropy)

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

$$y_n = p(\mathcal{C}_1|\phi_n)$$

## 4.3.2 로지스틱 회귀

교차 엔트로피  
(cross entropy)

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma).$$

이용

$\mathbf{w}$ 에 대하여 오류 함수 계산

$$\left| \begin{array}{l} y_n = \sigma(a_n) \\ a_n = \mathbf{w}^T \phi_n \end{array} \right.$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \underbrace{(y_n - t_n)}_{\text{오류}} \phi_n$$

$$\left| \begin{array}{l} \nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T. \end{array} \right. \text{선형 회귀 모델의 제곱합 오류 함수의 기울기}$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

가중치 벡터 업데이트 식

## 4.3.2 로지스틱 회귀

선형 분리 가능한 데이터 집합에 최대 가능도 방법 사용

→ 심각한 과적합 문제

헤비사이드 계단 함수, 각각의 클래스  $K$ 에서 온 모든 훈련 포인트들이 사후 확률 1을 가지게 됨

최대 가능도 방법을 통해서 어떤 하나의 해를 다른 해보다 선호하게 할 수 없음. 해를 찾는 건 최적화 알고리즘과 매개변수 초기화에 달림

→ 사전 확률 포함,  $w$ 에 대해 MAP해를 찾는 방식, 오류함수에 정규화항 추가하는 방식으로 해결 가능

## 4.3.3 반복 재가중 최소 제공법

가우시안 노이즈 모델 – 최대 가능도 해 : 닫힌 형태

why? 로그 가능도 함수가 매개변수 벡터  $\mathbf{w}$ 에 대해 이차 종속성 가짐

로지스틱 회귀 모델

닫힌 해X

로지스틱 시그모이드 함수의 비선형성

이차식 형태 → 오류 함수 **Convex** → 유일한 최솟값



<뉴턴 라프슨(Newton-Raphson) >

로그 가능도 함수에 대한 지역적인 이차식 근삿값을 구하는 방식

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

헤시안 행렬

최소화 대상

$\mathbf{H} = \nabla \nabla E(\mathbf{w})$  각 원소 :  $E(\mathbf{w})$ 를  $\mathbf{w}$ 의 각 성분으로 이차 미분한 값

## 4.3.3 반복 재가중 최소 제곱법

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (3.3)$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 \quad (3.12)$$

<뉴턴 라프슨 적용>

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{t}$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

$\boldsymbol{\Phi}$  is the  $N \times M$  design matrix

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \{ \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w}^{(\text{old})} - \boldsymbol{\Phi}^T \mathbf{t} \} \\ &= (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

이 결괏값은 표준 최소 제곱 해

오류 함수 : 이차식

뉴턴 라프슨 공식 한 단계 만에 정확한 해 구해냄

## 4.3.3 반복 재가중 최소 제곱법

$$\frac{d\sigma}{da} = \sigma(1 - \sigma).$$

이용

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \text{ 교차 엔트로피 함수}$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t})$$

오류 함수 기울기

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

헤시안

$$R_{nn} = y_n (1 - y_n).$$

$N \times N$ : 대각행렬  $\mathbf{R}$

시그모이드 출력 범위  $0 < y_n < 1$

$$\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$$

헤시안 - 상수  $\times$   
가중 행렬  $\mathbf{R}$ 을 통해  $\mathbf{w}$ 에 종속성 가짐  
 $\rightarrow$  오류 함수가 quadratic 아님

$\mathbf{H}$ 는 양의 정부호 행렬

$\rightarrow$  오류 함수  $\mathbf{w}$ 에 대한 convex, 유일한 최솟값 가짐



## 4.3.3 반복 재가중 최소 제공법

로지스틱 회귀 모델에 대한 뉴턴 라프슨 업데이트

$$\begin{aligned}\mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}\end{aligned}$$

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}).$$

$\mathbf{z}$ 는  $N$ 차원 벡터

→ 가중된 최소 제공 문제의 정규 방정식의 집합 형태

행렬  $\mathbf{R}$ 이 상수가 아니고 벡터  $\mathbf{w}$ 에 종속적 이기에  
정규 방정식들을 ‘반복적’으로 적용해야 함.

매 반복마다 새 가중 벡터  $\mathbf{w}$ 를 이용하여 수정된 가중 행렬  $\mathbf{R}$ 을 구해야 함.

## 4.3.3 반복 재가중 최소 제곱법 = IRLS 알고리즘

대각 가중 행렬  $R$ 의 원소  $\rightarrow$  분산으로 해석

$\downarrow$  (로지스틱 회귀 모델에서의  $t$ 의 평균과 분산)

$$\begin{aligned}\mathbb{E}[t] &= \sigma(\mathbf{x}) = y \\ \text{var}[t] &= \mathbb{E}[t^2] - \mathbb{E}[t]^2 = \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2 = y(1 - y)\end{aligned}$$

$$\begin{aligned}t &\in \{0, 1\} \\ t &= t^2\end{aligned}$$

IRLS를  $a = \bar{\mathbf{w}}^T \phi$ .

공간상에서 선형된 문제의 해로 해석

$$\begin{aligned}a_n(\mathbf{w}) &\simeq a_n(\mathbf{w}^{(\text{old})}) + \left. \frac{da_n}{dy_n} \right|_{\mathbf{w}^{(\text{old})}} (t_n - y_n) \\ &= \phi_n^T \mathbf{w}^{(\text{old})} - \frac{(y_n - t_n)}{y_n(1 - y_n)} = z_n.\end{aligned}$$

## 4.3.4 다중 클래스 로지스틱 회귀

생성적 모델링

4.2 내용

$p(\mathbf{x}|\mathbf{C}_k)$ : 클래스별 조건부 밀도

$p(\mathbf{C}_k)$ : 사전 클래스 확률



$p(\mathbf{C}_k|\mathbf{x})$ : 사후 확률 결정(간접적으로 파라미터 결정)

---

$$p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad \Bigg| \quad a_k = \mathbf{w}_k^T \phi.$$

$\mathbf{w}_k \leftarrow$  최대 가능도 방법 활용하여 직업 구하기

$$\frac{\partial y_k}{\partial a_j} = y_k(\underline{I_{kj}} - y_j)$$

항등 행렬의 원소

## 4.3.4 다중 클래스 로지스틱 회귀

가능도 함수(원 핫 인코딩 사용해서 표현)

$$p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(\mathcal{C}_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad \left| \quad y_{nk} = y_k(\phi_n)\right.$$

음의 로그값

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

N × K 행렬, 원소  $t_{nk}$

다중 클래스 분류 문제상에서의 교차 엔트로피(cross entropy) 오류 함수

## 4.3.4 다중 클래스 로지스틱 회귀

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \underbrace{(y_{nj} - t_{nj})}_{\text{오류값}} \underbrace{\phi_n}_{\text{기저함수}} \quad \leftarrow \text{이용}$$

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

$$\sum_k t_{nk} = 1$$

소프트맥스의 특징

선형 모델에 제곱합 오류 함수를 사용했을 경우의 기울기 함수 형태가 같음.

→ 순차적 알고리즘 만들 수 있음.

교차 엔트로피 오류 함수  $\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$

소프트맥스 활성화 함수, 다중 클래스 식 모두 같은 형태

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T. \quad \text{선형 회귀 모델의 제곱합 오류 함수의 기울기}$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

가중치 벡터 업데이트 식

## 4.3.4 다중 클래스 로지스틱 회귀

Batch 알고리즘 찾기

→ 뉴턴 라프슨 업데이트 적용  
다중 클래스상에서 IRLS 알고리즘 구하기

계산 
$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T.$$

다중 클래스 로지스틱 회귀 모델(이중 클래스와 마찬가지로)

H는 양의 정부호 행렬

→ 오류 함수 w에 대한 convex, 유일한 최솟값 가짐

## 4.3.5 프로빗 회귀

### 4.2.4 지수족

클래스  $K = 2$   $\longrightarrow a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$  식 4.58 대입

$$a(\mathbf{x}) = (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(\mathcal{C}_1) - \ln p(\mathcal{C}_2).$$

사후 클래스 확률이 선형 함수  $a(\mathbf{x})$ 에 대한 로지스틱 시그모이드 함수로 주어짐

클래스  $K$ 개 ( $K \geq 2$ )  $\longrightarrow a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ . 식 4.63 대입

$$a_k(\mathbf{x}) = \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(\mathcal{C}_k)$$

$\mathbf{x}$ 에 대한 선형 함수

이런 단순한 형태의 사후 확률이 모든 종류의 클래스 조건부 밀도 분포에 대해서 결과로 나오는 것은 아님.

## 4.3.5 프로빗 회귀

$p(t = 1|a) = \underline{f(a)}$  일반화된 선형 모델의 틀

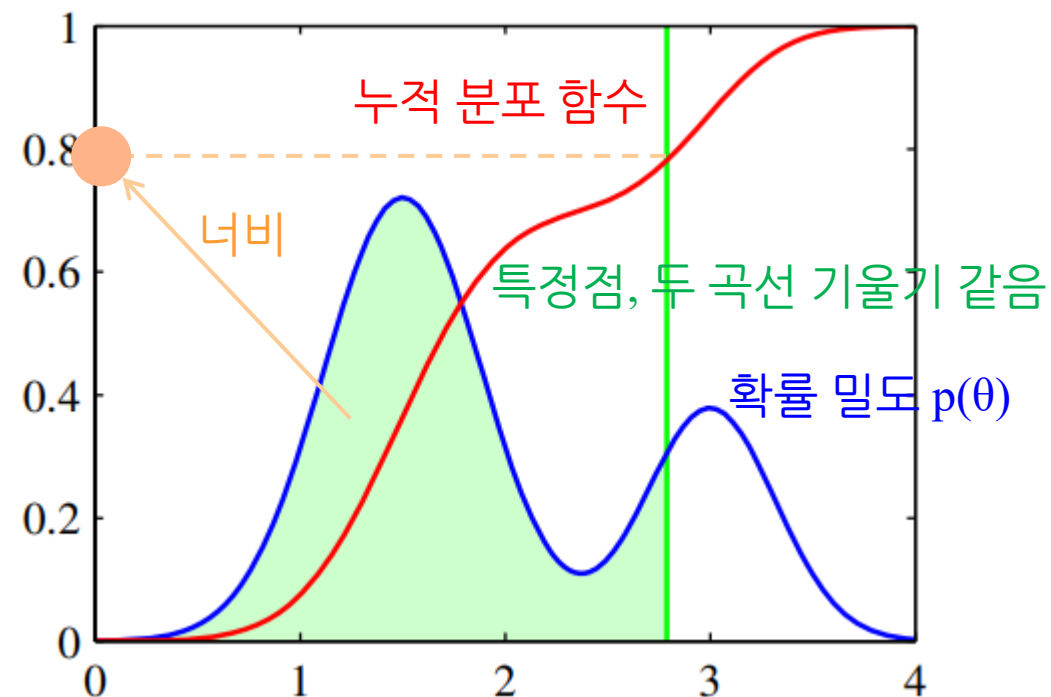
$a = \mathbf{w}^T \boldsymbol{\phi}$   $f$ 는 활성 함수

대안적인 연결 함수 : 노이즈 임계값 모델

$$\begin{cases} t_n = 1 & \text{if } a_n \geq \theta \\ t_n = 0 & \text{otherwise.} \end{cases}$$

### 대안적인 연결 함수

일반화 선형 모델에서 사용되며, 선형 예측값을 종속 변수의 기대값과 연결합니다. 주로 확률적 성격을 띠는 모델에서 사용되며, 특정 확률 분포와 관련이 있습니다.



$$f(a) = \int_{-\infty}^a \underline{p(\theta)} d\theta$$

확률 밀도



## 4.3.5 프로빗 회귀

$$f(a) = \int_{-\infty}^a \underline{p(\theta)} d\theta$$

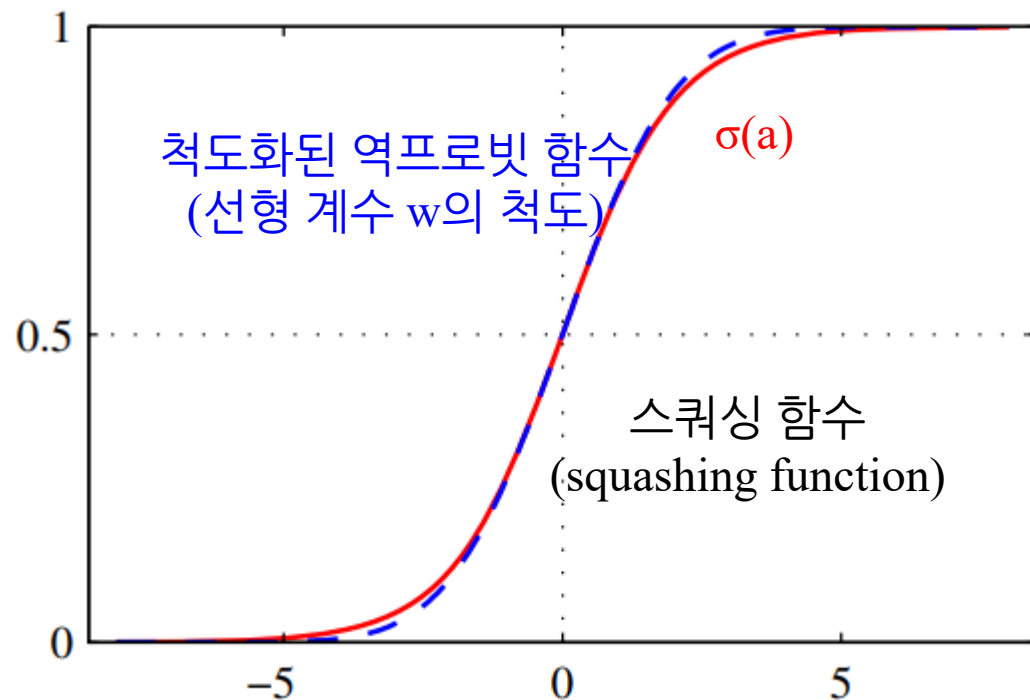
확률 밀도

가우시안 분포 가정: 평균 0, 분산 1

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0, 1) d\theta \quad \text{역프로빗(inverse probit)}$$

오차 함수(오류 함수 아님)

$$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-\theta^2/2) d\theta$$



역프로빗과 오차 함수 관계

$$\Phi(a) = \frac{1}{2} \left\{ 1 + \frac{1}{\sqrt{2}} \text{erf}(a) \right\}$$

역프로빗 활성화 함수를 바탕으로 한 일  
반화된 선형 모델 : 프로빗 회귀

## 4.3.5 프로빗 회귀

실제 응용 사례 이슈

“이상값(outlier)”

입력 벡터  $x$  측정 시 오류  
표적 벡터  $t$  잘못 라벨링

이상적인 결정 경계면을 기준으로 잘못된 방향에 놓임  
→ 분류기를 심각하게 왜곡할 수 있음.

$x \rightarrow \infty$

로지스틱 회귀 모델

점근적으로  $\exp(-x)$ 와 같이 감소

반응

프로빗 회귀 모델

점근적으로  $\exp(-x^2)$ 와 같이 감소

더 예민하게 반응

## 4.3.5 프로빗 회귀

데이터가 올바르게 라벨링 되었다고 가정

표적값  $t$ 가 잘못된 값을 부여받았을 확률  $\epsilon$

$$\begin{aligned} p(t|\mathbf{x}) &= (1 - \epsilon)\sigma(\mathbf{x}) + \epsilon(1 - \sigma(\mathbf{x})) \quad \longleftarrow \text{잘못된 라벨링의 효과를 확률적 모델에 적용} \\ &= \epsilon + (1 - 2\epsilon)\sigma(\mathbf{x}) \end{aligned}$$

$\sigma(\mathbf{x})$  활성화 함수

$\epsilon$ : 미리 정해두거나 데이터로부터 유추할 수 있는 초매개변수처럼 다루기 가능

## 4.3.6 정준 연결 함수

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \underbrace{(y_{nj} - t_{nj})}_{\text{오류값}} \underbrace{\phi_n}_{\text{기저함수}} \leftarrow \text{이용}$$

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

$$\sum_k t_{nk} = 1$$

소프트맥스의 특징

선형 모델에 제곱합 오류 함수를 사용했을 경우의 기울기 함수 형태가 같음.

→ 순차적 알고리즘 만들 수 있음.

교차 엔트로피 오류 함수

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

소프트맥스 활성화 함수, 다중 클래스 식 모두 같은 형태

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T.$$

선형 회귀 모델의 제곱합 오류 함수의 기울기

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

가중치 벡터 업데이트 식

타깃 변수의 조건부 분포가 지수족에 포함되는 경우에 일반적으로 얻을 수 있다는 증명 부분

이때의 활성화 함수

정준 연결 함수 (canonical link function)

## 4.3.6 정준 연결 함수

$p(\mathbf{x}|\boldsymbol{\lambda}_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k) \exp \{ \boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x}) \}$  지수족 성질을 이용하여  $\mathbf{x}$ 에 대한 분포 작성

$\mathbf{u}(\mathbf{x}) = \mathbf{x}$ 인 부분 집합들에 대해서만 고려

$$p(\mathbf{x}|\boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp \left\{ \frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x} \right\} \longleftarrow \begin{array}{l} p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right) \\ \text{척도 매개변수 } s \text{ 도입} \end{array}$$

이번에는 타깃 변수  $t$ 를 지수족 분포로 가정

$$p(t|\eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp \left\{ \frac{\eta t}{s} \right\}.$$

$y$ 로 표현되는  $t$ 의 조건부 평균



$$y \equiv \mathbb{E}[t|\eta] = -s \frac{d}{d\eta} \ln g(\eta).$$

$y$ 와  $\eta$ 는 반드시 연관성이 있어야 함

$$\eta = \psi(y).$$

머신 러닝 문헌 :  $f$ 는 활성화 함수  
통계학 문헌 :  $f$ 의 역함수는 연결 함수

가정 : 모든 관측값은 동일한 척도 매개변수 공유  
(가우시안 분포는 노이즈의 분산)  $\rightarrow s$ 는  $n$ 에 대해 독립적

## 4.3.6 정준 연결 함수

Nelder and Wedderburn에 따르면

일반화된 선형 모델은  $y$ 가 입력(또는 특징)의 선형 결합을 비선형 함수에 넣은 결과로 표현되는 모델이라 정의

$$y = f(\mathbf{w}^T \boldsymbol{\phi})$$

$\eta$ 의 함수로 표현되는 이 모델의 로그 기능도 함수

$$\ln p(\mathbf{t}|\eta, s) = \sum_{n=1}^N \ln p(t_n|\eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const}$$



모델 매개변수  $\mathbf{w}$ 에 대해 미분

$$\begin{aligned} \nabla_{\mathbf{w}} \ln p(\mathbf{t}|\eta, s) &= \sum_{n=1}^N \left\{ \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right\} \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n \\ &= \sum_{n=1}^N \frac{1}{s} \{t_n - y_n\} \psi'(y_n) f'(a_n) \boldsymbol{\phi}_n \end{aligned}$$

$$a_n = \mathbf{w}^T \boldsymbol{\phi}_n$$

$$y \equiv \mathbb{E}[t|\eta] = -s \frac{d}{d\eta} \ln g(\eta).$$

$$y_n = f(a_n)$$

$$f^{-1}(y) = \psi(y)$$

$$f(\psi(y)) = y$$

$$f'(\psi) \psi'(y) = 1$$

$$\nabla \ln E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \boldsymbol{\phi}_n.$$

오류 함수의 기울기

## 4.4 라플라스 근사

목적 : 연속 변수의 집합에 대해 정의된 확률 밀도의 가우시안 근사치 찾기

$$p(z) = \frac{1}{Z} f(z) \quad z \text{는 단일 연속 변수}$$

$$\text{정규화 계수 } Z = \int f(z) dz$$

“분포  $p(z)$ 의 최빈값을 중심으로 한 가우시안 근사  $q(z)$  찾기”

1)  $p(z)$ 의 최빈값 구하기

$p'(z_0) = 0$ 이 되는  $z_0$ 값 찾기

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0.$$

## 4.4 라플라스 근사

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2}A(z - z_0)^2$$

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}.$$

$z_0$ 가 분포의 최댓값  $\rightarrow$  테일러 전개 일차항 안 나타남

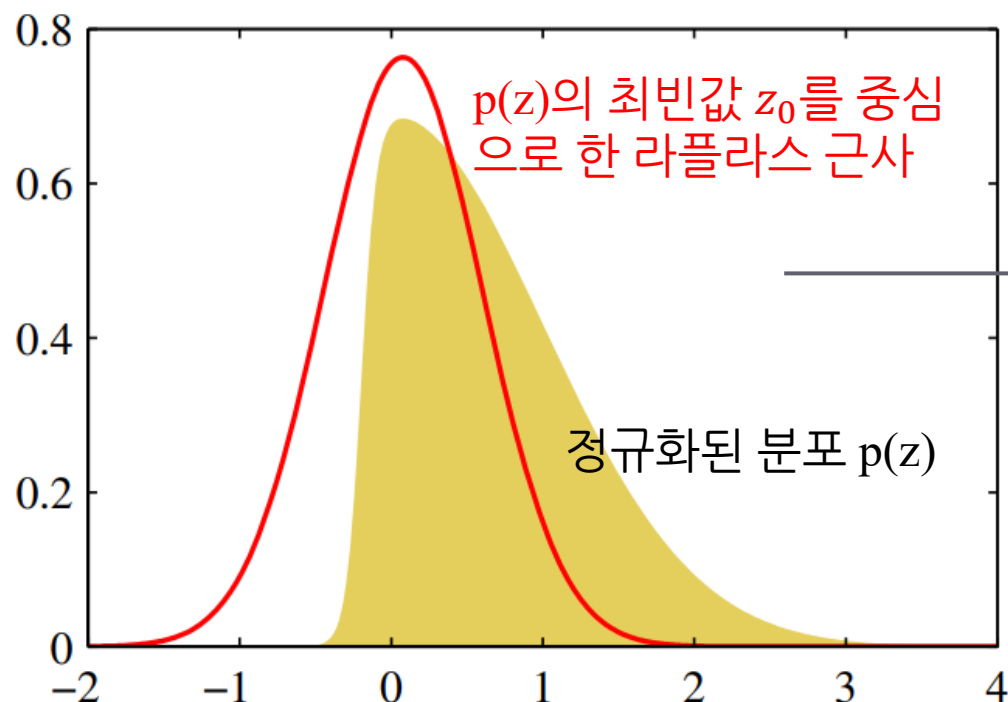
$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}.$$

정규화된 가우시안 분포에 대한 표준 방법 이용  
 $\rightarrow$  정규화된 분포  $q(z)$ 를 구할 수 있음.

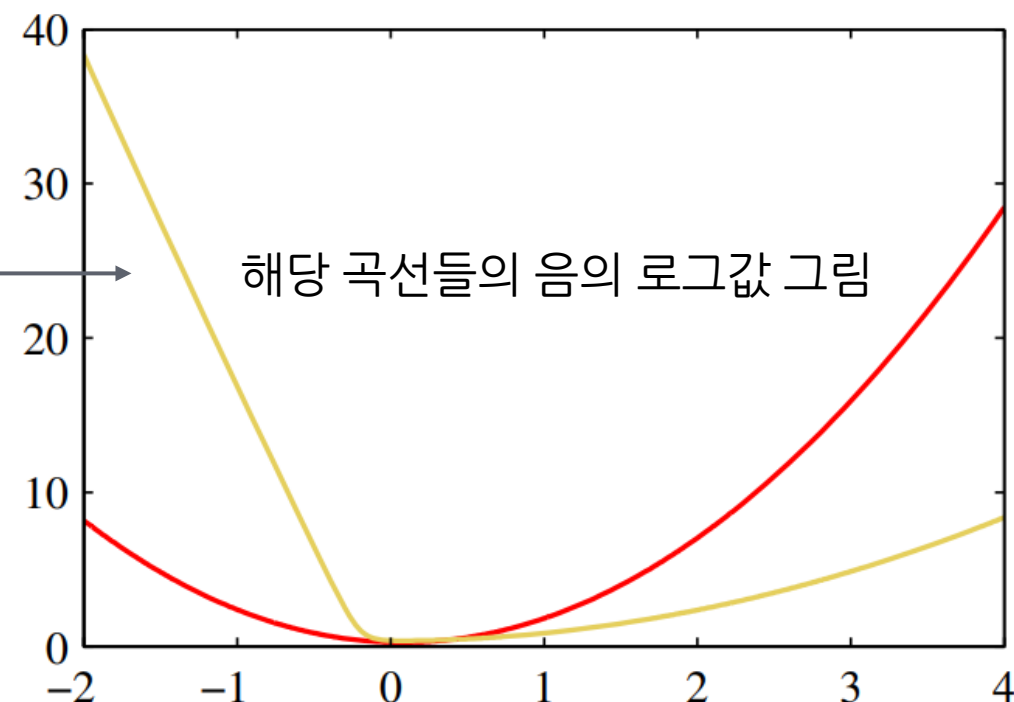
$$q(z) = \left( \frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}.$$



## 4.4 라플라스 근사



$$p(z) \propto \exp(-z^2/2)\sigma(20z + 4)$$



분포의 정밀도  $A$ 가 0보다 큰 경우에만 잘 정의  
(임계점  $z_0$ 가 지역적 최댓값,  $f(z)$ 의  $z_0$ 에서의 2차 미분값이 음수여야함)

## 4.4 라플라스 근사

m차원 공간  $\mathbf{z}$ 에 대해 정의된 분포  
 $p(\mathbf{z}) = f(\mathbf{z}) / Z$ 에 대해서 라플라스 근사 확장

$$\nabla f(\mathbf{z})$$

임계점에서 기울기 0

→ 임계점 근처에 대해서 전개 시행

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

$$\left| \mathbf{A} = -\nabla \nabla \ln f(\mathbf{z}) \right|_{\mathbf{z}=\mathbf{z}_0} \quad M \times M \text{ 헤시안 행렬 } \mathbf{A} \text{의 정의}$$

$\nabla$ 는 기울기 연산자, 양쪽 변에 대해 지수 함수 취하기

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}.$$

분포  $p(\mathbf{z})$ 는  $f(\mathbf{z})$ 에 비례, 적절한 정규화 계수는 다변량 가우시안에 대한 표준 결과를 이용해서 구하기 가능

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$$

## 4.4.1 모델 비교와 베이esian 정보 기준

$P(z)$  말고  $Z$ 도 근사치 구할 수 있음

$$\left. \begin{aligned} Z &= \int f(\mathbf{z}) \, d\mathbf{z} \\ &\simeq f(\mathbf{z}_0) \int \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} \, d\mathbf{z} \\ &= f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned} \right| \begin{aligned} p(\mathcal{D}) &= \int p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \, d\boldsymbol{\theta}. \\ f(\boldsymbol{\theta}) &= p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \text{ and } Z = p(\mathcal{D}) \end{aligned}$$

---

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \underbrace{\ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|}_{\text{Occam factor}}$$

$$\mathbf{A} = -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) p(\boldsymbol{\theta}_{\text{MAP}}) = -\nabla \nabla \ln p(\boldsymbol{\theta}_{\text{MAP}}|\mathcal{D}).$$

## 4.4.1 모델 비교와 베이지안 정보 기준

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \underbrace{\ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|}_{\text{Occam factor}}$$

↓ 매개변수에 대한 가우시안 사전 분포가 넓게 퍼졌고,  
헤시안 행렬이 최대 계수를 가진다 가정

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2} M \ln N$$

베이지안 정보 기준(BIC) = 슈바르츠 기준

아카이케 정보량 기준(AIC)과 비교했을 때  
베이지안 정보 기준은 모델의 복잡도에 대해 더 큰 불이익을 줌.

## 4.4.1 모델 비교와 베이지안 정보 기준

복잡도를 측정하는 단위는 측정이 쉽다는 장점이 있지만

잘못된 결과도 끌어낼 수 있음.

Ex) 헤시안 행렬이 완전 행렬 계수를 가진다는 가정은  
대부분 사실이 아님 ← 매개변수들이 잘 확정되지 않기 때문

## 4.5 베이지안 로지스틱 회귀

로지스틱 회귀의 정확한 베이지안 추론은 다루기 아주 어려움

사후 분포 계산  $\leftarrow$  정규화 (사전 분포 \* 가능도 함수)

각 데이터 포인트마다  
로지스틱 시그모이드  
함수를 모두 곱한 값

예측 분포 계산  $\leftarrow$  사후 분포와 비슷한 이유로 다루기 어려움.

→ 베이지안 로지스틱 회귀 문제에 라플라스 근사를 적용할 것

## 4.5.1 라플라스 근사

사후 분포의 가우시안 표현 찾기

$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$  가우시안 사전 분포

고정된 초매개변수

$p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w}) p(\mathbf{t} | \mathbf{w})$   $\mathbf{w}$ 에 대한 사후 분포

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

양변에 로그  
사전 분포 식  
기능도 함수  $\rightarrow$

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

$$\begin{aligned} \ln p(\mathbf{w} | \mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{const} \end{aligned}$$

$$y_n = \sigma(\mathbf{w}^T \phi_n).$$

## 4.5.1 라플라스 근사

사후 분포의 가우시안 표현 찾기

사후 분포를 최대화  $\rightarrow$  최대 사후 분포 해인  $\mathbf{w}_{map}$  을 구해야 함.

가우시안 분포의 평균값 정의 가능

공분산 : 음의 로그 가능도 함수의 이차 미분값 행렬의 역

$$\mathbf{S}_N = -\nabla\nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n)\phi_n\phi_n^T.$$

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{MAP}, \mathbf{S}_N).$$

사후 분포의 가우시안 근삿값



## 4.5.2 예측 분포

실제 예측을 위해 주변화 시행

새 특징 벡터 :  $\phi(x)$ , 클래스  $C_1$ 에 대한 예측 분포

사후 분포  $p(\mathbf{w}|\mathbf{t})$ 에 대해 주변화 시행

가우시안 분포  $q(\mathbf{w})$ 로 근사 가능

$$p(C_1|\phi, \mathbf{t}) = \int p(C_1|\phi, \mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$$

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad \leftarrow \begin{array}{l} \text{w에 대해 종속적인 부분} \\ \text{: w의 } \phi \text{에 대한 투영} \end{array} \quad a = \mathbf{w}^T \phi$$

클래스  $C_2$ 에 대한 확률  $p(C_2|\phi, \mathbf{t}) = 1 - p(C_1|\phi, \mathbf{t})$

## 4.5.2 예측 분포

실제 예측을 위해 주변화 시행

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da$$

$\delta$ 는 디랙(Dirac) 델타 함수이기에 다음을 구할 수 있음.

$$\int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da$$

$$p(a) = \int \delta(a - \mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}.$$

델타 함수로 인해  $\mathbf{w}$ 에 대해 선형 제약 조건이 걸림

## 4.5.2 예측 분포

$$\underline{p(a)} = \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w}.$$

결합 분포  $q(\mathbf{w})$ 를  $\phi$ 에 대해 직교하는  
모든 방향으로 적분하여 주변 분포 만들기

$q(\mathbf{w})$  : 가우시안 분포  $\rightarrow$  주변 분포도 가우시안

$$\text{평균 } \mu_a = \mathbb{E}[a] = \int p(a) a da = \int q(\mathbf{w}) \mathbf{w}^T \boldsymbol{\phi} d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}$$

$$\begin{aligned} \text{공분산 } \sigma_a^2 &= \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} da \\ &= \int q(\mathbf{w}) \{(\mathbf{w}^T \boldsymbol{\phi})^2 - (\mathbf{m}_N^T \boldsymbol{\phi})^2\} d\mathbf{w} = \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}. \end{aligned}$$

## 4.5.2 예측 분포

$$p(\mathcal{C}_1|\mathbf{t}) = \int \sigma(a)p(a) \, da = \int \underline{\sigma(a)\mathcal{N}(a|\mu_a, \sigma_a^2)} \, da.$$

a의 분포가 아래의 선형 회귀 모델에서의 예측 분포  
가 노이즈 분산이 0일 때와 같은 형태를 띈다

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T\phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

## 4.5.2 예측 분포

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

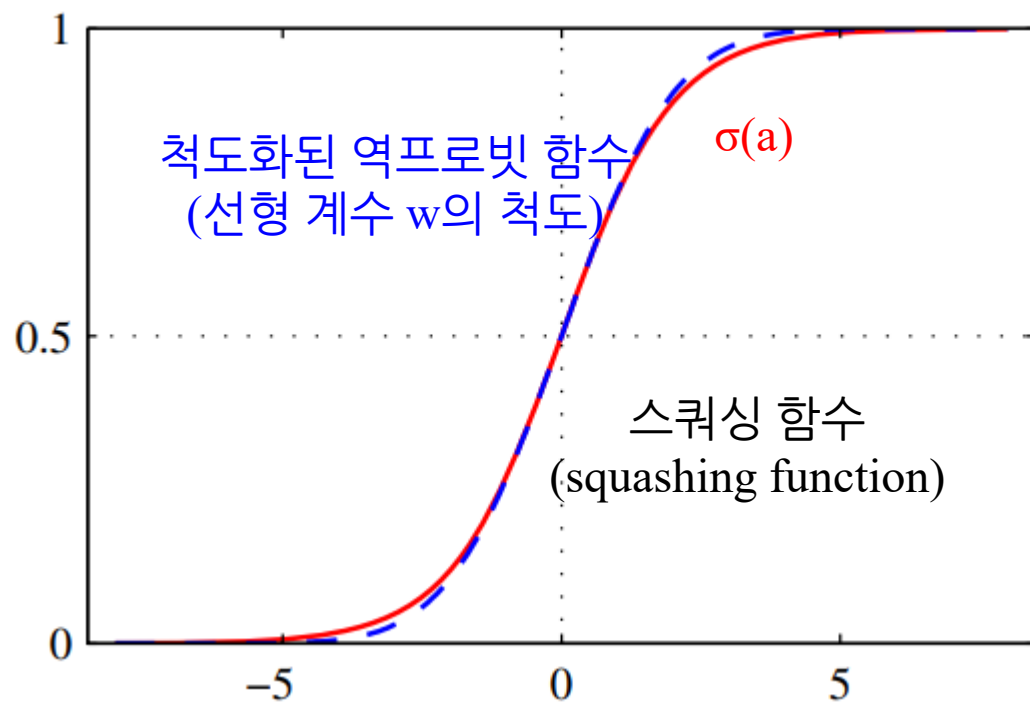
a에 대한 적분은 가우시안을 시그모이드와 컨볼루션 한 것

해석적으로 계산 불가

→ 로지스틱 시그모이드 함수와 역프로빗  
함수의 유사성 이용해서 유사 분석 가능

$$\lambda^2 = \pi/8.$$

로지스틱 함수의 최고 근사치를 얻기 위해 가로축의 척도 변경



## 4.5.2 예측 분포

역프로빗 함수와 가우시안 콘볼루션  
→ 다른 역프로빗 함수로 표현 가능

$$\int \Phi(\lambda a) \mathcal{N}(a|\mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right).$$

$$\sigma(a) \simeq \Phi(\lambda a)$$

양변에 대입

로지스틱 시그모이드와 가우시안의 콘볼루션에 대한 다음의 근사 얻기 가능

$$\int \sigma(a) \mathcal{N}(a|\mu, \sigma^2) da \simeq \sigma(\kappa(\sigma^2)\mu)$$

$$\kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2}.$$

여기에 적용

$$p(\mathcal{C}_1|\mathbf{t}) = \int \sigma(a)p(a) da = \int \sigma(a)\mathcal{N}(a|\mu_a, \sigma_a^2) da.$$

$$p(\mathcal{C}_1|\phi, \mathbf{t}) = \sigma(\kappa(\sigma_a^2)\mu_a)$$

예측 분포에 대한 근사치

## 4.5.2 예측 분포

Thank you