

## 12장 판별분석

### 1. 서론

판별분석 - 분류기법, 로지스틱 회귀분석과 같이 분류와 프로파일링에 쓰이는 전통적 통계기법  
앞에서본 다른 분류기법들과 유사

### 2. 클래스로부터 관측자에 이르는 거리

클래스 - 관측자 간의 거리측정이 필요 - 가장 가까운 클래스에 배정

기초적인 방법: 유클리드 거리 (Euclidean distance)

$$D_{\text{Euclidean}}(X, \bar{X}) = \sqrt{(x_1 - \bar{x}_1)^2 + \dots + (x_p - \bar{x}_p)^2} \rightarrow \text{단점 존재}$$

① 거리측정이 임의변수에 대한 특정단위 선택에 따라 달라짐

↳ 연간 소득을 천 달러 기준으로 바꾸면 거리도 달라짐

② 변수의 변동성을 반영하지 못함

↳ 각 변수들의 변동성 (표준편차, 분산)을 반영하지 못함

⇒ 표준화된 값 (Z-score 등)을 이용해야 함

③ 변수들 간의 상관관계를 무시

변수가 많을 경우 특히 중요 - 상관관계에 의해서 변수 중요도가 바뀔 수 있음

⇒ 이러한 단점 해결책으로 통계적 거리 (Statistical Distance) 혹은 마할라노비스 거리 사용 (Mahalanobis)

↳ p개의 변수들 사이의 공분산 행렬 S

$$\begin{aligned} D_{\text{sta}}(X, \bar{X}) &= [X - \bar{X}]' S^{-1} [X - \bar{X}] \\ &= [(x_1 - \bar{x}_1), (x_2 - \bar{x}_2), \dots, (x_p - \bar{x}_p)] S^{-1} \begin{bmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \\ \vdots \\ x_p - \bar{x}_p \end{bmatrix} \end{aligned}$$

$[X - \bar{X}]'$  → 열벡터를 행벡터로 변환  $S^{-1}$  → S의 역행렬

공분산 행렬을 계산에 포함

↳ 통계적 거리는 예측변수들의 평균값 뿐만 아니라 예측변수 값들의 분산성과 상관관계도 고려

↳ 예측변수의 평균(중심점)과 변수들 간의 공분산 계산 필요

⇒ 판별분석은 통계적 거리를 기반으로 분리선 (예측변수가 세 개 이상인 경우에는 분리초평면)을 찾는다.

↳ 모든 클래스의 평균들로부터 동일 거리에 위치

### 3. 다수의 선형 분류함수

Linear Classification Problem - k개의 가변수와 m개의 클래스를 분류하는 문제



그 특  $\log(p_{ij})$ 를 분류함수에 상수로 추가

## 6. 서로 다른 오분류 비용

클래스 간의 오분류 비용이 현저하게 다를 경우 오분류의 기대 비용을 최소화해야 한다.

↳ 클래스 1에 속하는 멤버를 오분류하는 비용을  $q_1$ , 클래스 2를 오분류하는 비용은  $q_2$  라 할 때,

$\log(q_1)$ ,  $\log(q_2)$ 를 각각의 분류함수에 상수로 더해주면 된다.

↳ 사선택률과 오분류 비용을 다 고려하기 위해서  $\log(p_1 q_1)$ 을 더해주면 됨.

↳ 실제로 오분류 비용을 계산하는 것보다 그 비례를  $(q_2/q_1)$ 을 계산하는 것이 더 쉽다.

⇒ 클래스 2의 분류함수에  $\log(q_2/q_1)$ 만 더해주면 됨.

## 8. 판별분석의 장단점

주로 통계적 분류기법으로 여겨짐

판별분석의 사용과 성능은 다중선형 회귀분석과 유사 → 장단점도 유사

↳ 판별분석은 예측 변수의 최적 가능치를 찾는 과정을 포함

선형 회귀분석에서의 가절치는 반응(종속) 변수와 관련  $\hookrightarrow$  판별분석에서는 클래스들을 분리시키는 것과 관련

두 분석 모두 최소제곱법 사용하는 추정방법 이용, 결과로 얻은 추정치들은 극소 최적값에 영향을 받지 않고 강건.

↳ 기본과정은 정규성(normality) 판별분석에서는 예측변수들이 다변량 정규분포를 따른다고 가정

현실적으로 지켜지지 않는 경우가 많지만, 이런 문제에 강건

너무 한쪽으로 치우친 경우에는 로그 변환 등을 시행 → 성능향상이 가능

↳ 분류 방법으로서의 장점: 단일 예측 변수의 기여도에 대한 추정치를 제공

계산과정이 간단. 간결하고 데이터세트가 작을 때 유용