## 10.1 10.2 두 개 그룹의 판별

2019년 6월 16일 일요일 오전 3:11

지상단에 대한 정보로부터 지난을 귀열할 수 있는 판별할수 또는 판별규칙을 만들고 새로운 개체에 대해 이느 지난에 속하는지를 판별하며 불류하는 다면한 기법

四级 = 工意

工意花의 大的量 3711 湖平台工意 亚姆拉宁 子和双形花

吐燈站午是工意花의 731量本12H3 型31计至率 吐气 (超午量의 化耐湿放作) 冠

판년하는 구한 후에는 판년하는수를 이용해 기존의 개체들을 분류하다 오분류율을 기계사하는.

从五元 TH和1011 THAH41는 亚岭社产 이용하다 车하는 그동을 추정站 수 있는

#### 단얼라 불병적의 목적

- 1) 受对的 空初礼工意。圣学的工意的 특성을 나타내주고 귀절해주는 社수를 결정하点.
- 2) 理智知 吐燈花午童 이용하다 (H至은 판측치를 판煙하다 7H知를 변報法.

#### 10.2.1 Fisher = 1154

年7H21 工意 UII文计正 吐煙站 수 있도록 (地수들의 化剂交流)으로 판煙站수를 구하고자站.
다(地球, 비리는 일(地球, 地수로 (地土)하지 그 등을 판별하는 나는 나는 지, 건규성 가건은 필요하지 않음.

그룹을 판별하기위해 대, 대로 가능한 한 딸이지 구별되게 하는 신해판별하수 Y= 오'X 를 구하고자하.
두 그룹은 공통 공임사행렬을 가지.

对工意创创 证理论学的 对证的放政 总化 , 总化量工行社 车工意的 对证 为时, Fisher的 化剂证理论学 P259 无본에 证此 Fisher的 化剂证性理论学, 无怪对证例时, 无怪恐怕处对望。 就是没是化物理。 P260

가 그룹의 덩균점은  $\overline{y_1} = \hat{\mathcal{L}}' \overline{\chi_1}$  ,  $\overline{y_2} = \hat{\mathcal{L}}' \overline{\chi_2}$  이고 이들의 중간점은  $\hat{m} = \frac{1}{2}(\overline{y_1} + \overline{y_2})$  이 됨. 서로운 관측에서  $\chi_0$  에 대한 판별감수와 그 판별감수를 이용한 판별규칙은 p260

Fisher의 位部延慢站台间 의部 似色色 y字이 部级되어 두 그룹의 국化程度 기준으로 두 그룹이 나뉘어 분포된

## 10.2.2 다면냥 컴퓨션도를 따르며 두 그룹의 공원산행렬이 같은 경우

f(X|G1): G1 工贵二至早日 时代社 X의 李量型五社个

PION P2는 각각 X7+ G1, G2 그룹에서 발생할 사전학을 (P2 = 1 - P1)

7321 10.1 오년유 학률을 소1소학하는 소1각 변유규칙 P263 만야 P1 = P2 이번 소1각 분유규칙 P263 이는 첫대유도를 이용한 규칙이 됨.

특히 X 7+ 다년라 경구분도를 따를 거유. ( P264 )

\$童望在这个, 如对 世界不刻, 平对证 小家

시내로이 관측된 시계된 X7+ 속하는 그룹에 대한 판별규칙 => 신해 판별 규칙

만야 PI = PZ 이번, 정규성으로부터 유도된 선행판별규칙은 Fisher의 판별규칙과 끝아지.
그리고 선행판별규칙을 이용한 분류규칙은 근사적으로 최적인 판별규칙이 됨.

10.2.3 다녀라 정규본도를 따르따 두 그룹의 공본산행렬이 다른 경우

千至UI, OI大时经验中, 如对 坚和和是 P266

无坚竟可容益如的大工世经社会 无坚社会 P266

문본하수를 이용한 분류규칙은 근사적으로 소계각규칙은 아닌

실제 판별인건에서 판별하는 구하고자 할 때, 그룹들이 공통공원산행결을 가지다하며 귀우가 결을 기가하게 못하는 기상에는 전해된 발달하는 사용하고 귀우가 결을 기가하는데 될 기상에는 이것나만일하는 선택하는 수 있음.

# 10.3 세 개 이상 그룹의 판별

2019년 6월 16일 일요일 오전 3:11

工意的 37H 이사상인 기상우 97H 工意의 차이를 가지상 크게 하도록 하는 1년수들의 전혀기결합식을 찾는것.

보산보다법이는.

#### 10.3.1 Fisherel 11/41

水場 5元 Zi = ○'Xi 이 되

가 그룹덩균을 분리하내는 식의 기준을 9개 그룹에 대해 착장하기 위하며 분산분석에서의 그룹간 해결 B 와 그룹 내 해결 E 를 이용해 포현값. 판별값수  $Z_1 = \alpha_1'$  는 그룹덩균들의 차이를 가장 크게 해주는, 즉 그룹을 구별해주는 값수가 됨.

시내로운 7H처1에 다나한 분류규칙은 P270

## 10.3.2 공보산행렬이 모두 같은 73우

转量经验产 叶芒 7% 如何 贴惯形剂 P271

(吡야 그룹에 대한 사전학률이 모두 准의단위 판별규칙은 유도함수만을 이용하게 되므로 최대유도를 이용한 규칙이 됨)

x7+ C+1년강등 정규본포를 따를 거유 化部站产生 吐煙开礼은 P271 이 판별규칙은 군사적으로 최적인 분류규칙이 된 10.3.3 공발산행렬이 또두 다를 거우 X7+ 다년당 정규범포를 따를 거우 圣堂二至早日191 李行記号 이용한 이大·新沙社수, 판岭至于礼鲁 P272 이 판별규칙은 군사적으로 최적인 분류규칙이 된 그룹에 대한 사건학률이 모두 같으면 이것나에서라는수에서 100기나 사라지

# 10.4 오분류율 계산

2019년 6월 16일 일요일 오전 3:11

### 10.4.1 재대성 발뉴에 의한 오류별 7계산

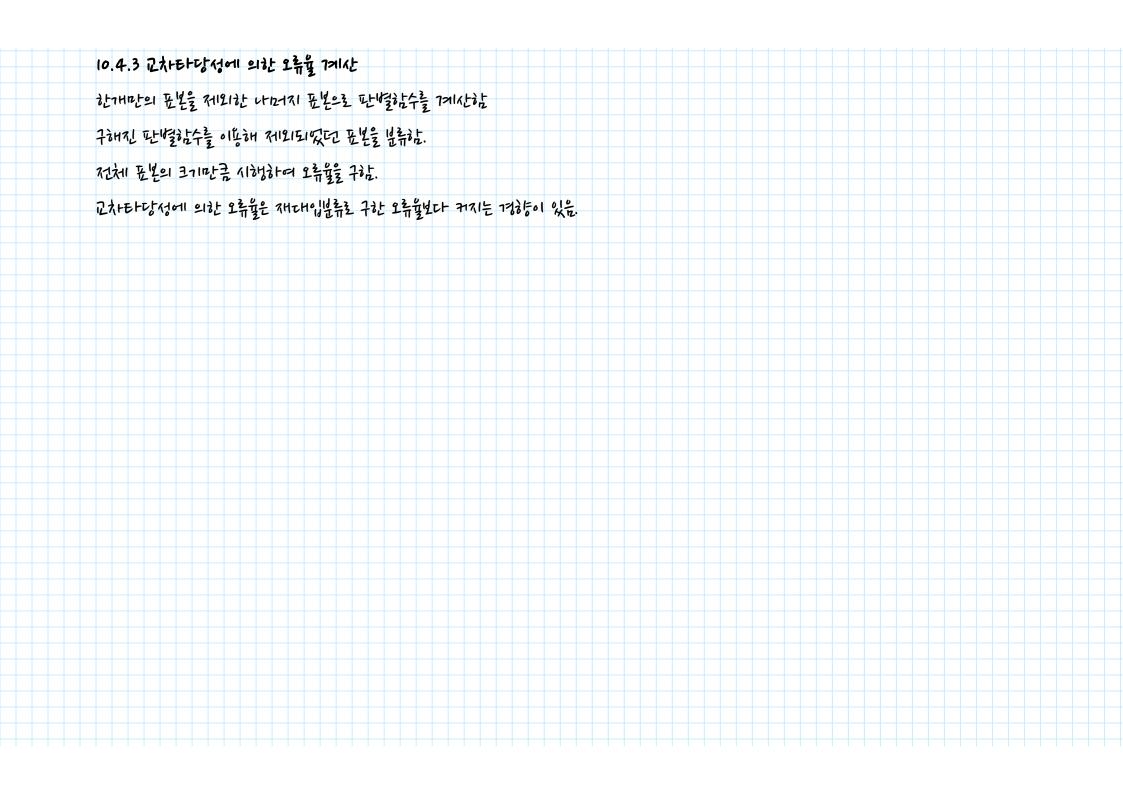
데이터 3부터 유도된 판별감수를 다시 데이터에 각당하는 제다비를 분류에 의해 오분류별을 기시산감 수 있는. 문을 만들어서 기시산간 데네식간 오류별, 정확간 분류별은 P274 제다비를 분류에 의간 오류별은 실제 bias 보다 각계 기시산될 수 있는.

#### 10.4.2 显性验明 의社 2清量 721化

포본을 두 부분으로 나누어 훈련포본은 판별학수를 만드는데 이용하고 타당성 검사 포본은 만든 판별학수를 이용해 분류한 후 판별학수의 판별 능력을 떨거나다. 추정된 오류율은 불편수정강등이 됨

단점1: 문본을 두 부분으로 나누어야 하나요 네고적 큰 문본의 크기가 오구

吐祖2: 설계逐 사袋鞋 吐燈站수에 대해서는 평가寺 수 때는



# 10.5 10.6 판별함수의 표준화 & 유의성 검정

2019년 6월 16일 일요일 오전 3:11

판년하는 에 기대하는 변수의 상대적 비혹은 판년하는 기계수에 나타보.
변수가 등의 단위가 다르면 동등한 비교 불가능
단위가 다를 때는 문문화면수를 이용하여 판년하는 구한 후 판년하는의 기계수를 비교하는.

두 개의 그룹이 있으면 공반산 행렬이 같다고 할 수 있는 13우는 P276 대의 개의 그룹이 있으면 공반산행렬이 같다고 할 수 있는 13우는 P277

Fisher 의 판별학수는 지난 또는 그룹간의 행관하다 첫대나 되도록 판별학수를 구하는 것이다고 그룹간의 행균 (벡터) 하이가 있다면 판별학수를 이용학 필요가 있게 된 => 판별학수의 유의성 경정
판별학수의 유의성 경쟁에는 경규성 가정이 필요함.

10.6.1 두 7H 그룹의 7경우 판년감수의 유의성 1년정 ( P278 )

두 그룹의 발리를 위해 두 그룹 간 거리가 최다니가 되도록 만든 판별하수의 7계수벡터 판별하수의 7계수벡터의 유의성에 대한 귀투가설 두 그룹의 덩균비교에 대한 가설건강은 Hotelling T-square 건강을 이용하.

10.6.2 미러 개 그룹의 거우 판별감수의 유의성 검정 ( P278 )

# 10.7 판별함수의 변수선택 : 단계별 변수 선택

2019년 6월 16일 일요일 오전 3:12

#### 10.7.1 地个位时比如

科学社 他个位时皇皇计 亚姆拉宁 核二时 亚姆林光色 圣堂 খ叶 아니라 生品量 快量 수 있는.

- 1) 전진격 선택 (forward selection): 그룹간의 거리를 참다비로 하는 한 개의 변수를 선택한 후 그 다음으로 그룹간의 거리를 참다비로 하는 다른 한 개의 (변수를 선택한. 1년수가 더 이사 선택되지 않을때마지 진행한.
- 3) 단기기각 시간 (Stepwise Selection): 전지각 시간 나방대라 후진각 시간 나방대의 혼하다. 변수가 시간되는 때 단기에 나다 시간 된 된 1년수가 기뻐하는 내를 건강하며 결정하는 가 단기에서 판별하는 를 기에 사하지는 이동이 설립적으로 때문기에 MANOVA 라강을 수하는 하는것. 1년수 시간 라강이 끝나면 시간되면 1년수를 이용하며 판별하는수를 구하는 이외 가는이 1년수를 시간되하며 판멸하는수를 따드는 각 기을 단기에 각 판별분석이라하는.

### 10.7.2 부분 F-통7계당을 이용한 단7개1/절 1년수 (전략 (P282)