

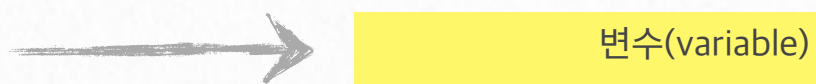
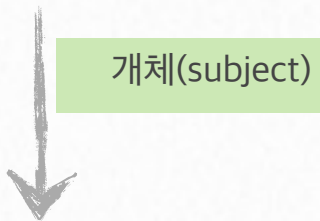
1

다변량분석 개념

1. 정의 Definition

2개 이상 (양적) 변수 데이터

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$



- 인과 관계(casual relationship)를 규명, 분석하거나 (회귀 분석: regression, 다변량 분산 분석: Multivariate ANOVA)
- 변수들의 상관관계(correlation)를 이용하여 변수의 차원을 축약(reduction)하거나
- 개체 간의 유사성을 측정하여 개체 분류(classification)에 관련된 분석 방법

일반적으로 (협의의) 다변량 분석은 후자 2개

2. 변수(variable)

1) 정의

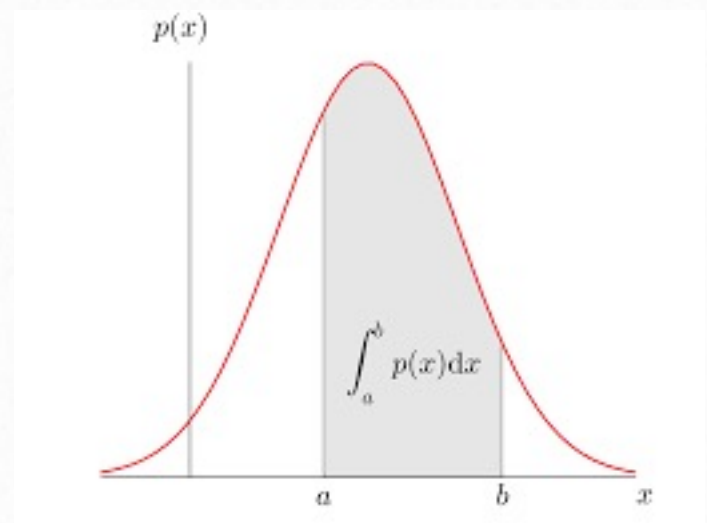
확률실험의 결과(원소 w)와 실수 값을 대응한 함수

$$X(w) = x$$

- 데이터에서 확률변수는 개체의 관심 특성의 총칭이며 변수의 관측값이 데이터

2) 확률분포함수

확률변수가 가지는 값과 그에 대응하는 확률을 표, 수식, 그래프로 나타낸 함수 ($f(x)$)



3) 변수 종류 : 분석적 측면

1) Metric (측정형 변수, measurable) : 실험 개체의 측정 가능한 특성을 측정한 변수로 키, 몸무게, 평점, IQ, 교통량, 사망자 수가 그 예이다. 연속형 변수는 모두 측정형 변수이고 이산형 변수 중 측정형 변수가 있을 수 있다. 예) 교통량

2) Non-metric (분류형 변수, classified, 범주형 categorical) : 개체를 분류하기 위해 측정된 변수를 의미하며 성별, 결혼여부 등이 그 예이다.

- 명목형 (nominal): 분류만 → 성별, 결혼여부, 소득 (단위: 만원)
- 순서형 (ordinal): 순서를 가진 분류 → 성적(A, B, ..) 소득수준(상, 중, 하), 5점 척도

대부분의 다변량 분석은 변수들의 분포가 다변량 정규 분포(multivariate normal distribution)를 가정하기 때문에 고전적 다변량 분석은 metric 변수들에 분석이나 non-metric 변수 중 순서형 변수에 대한 분석 방법도 제안되고 있다.

4) 변수종류 : 시간에 의해

자료가 시간적 순서를 가지면 이를 시계열 자료(time series)라 하고 그렇지 않은 경우를 횡단면 자료(Cross-section: 일정 시간에 한꺼번에 조사)라 한다.

경제 지표(환율, 수출량)나 기업의 연차별 자료(연도별 매출액)가 시계열 자료에 해당된다. 시계열 다변량 자료에 대한 분석 방법으로는 시계열 자료 분석, 계량 경제(econometrics) 방법을 이용하면 된다.

5) 변수종류 : 인과 관계

인과 관계(casual relationship)에서 원인이 되는(영향을 주는) 변수를 설명 변수(exploratory var.) 혹은 독립 변수(independent var.)라 하고 결과나 영향을 받는 변수를 종속 변수(dependent var.) 혹은 반응 변수(response var.)라 한다. 종속 변수는 Y, 설명 변수는 X로 표시한다. 분산 분석에서는 설명 변수를 처리 효과, 요인으로 불리어진다.

인과 관계는 이론적, 경험적 타당성에 근거하여 연구 목적에 설정되는 것이지 자료 분석 후 인과 관계가 설정되는 것은 아니다.

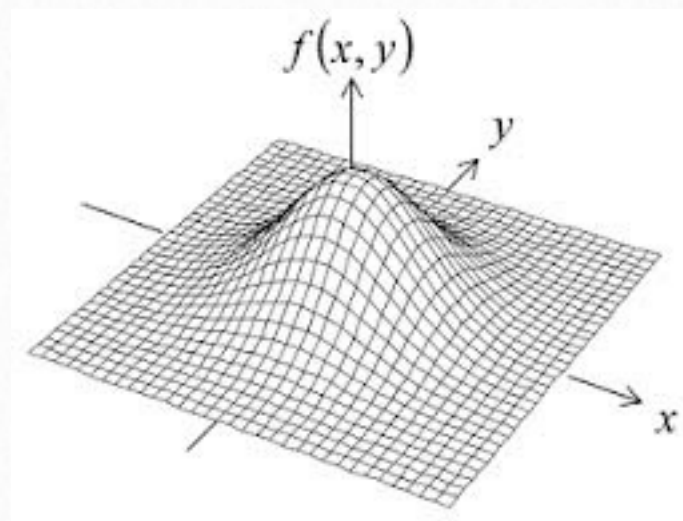
3. 다변량분포

확률변수가 $p (\geq 2)$ 개인 결합확률분포함수, 다변량분석에서는 변수들이 좌우 대칭인 분포를 따르는 것이 적합하다. 그렇지 않으면 적절한 정규변환이 필요함

정규변환 강의

1) 이변량 정규분포

확률변수 (x, y) 가 이변량 정규분포를 따른다고 하면 결합확률분포함수 그래프는 다음과 같다.

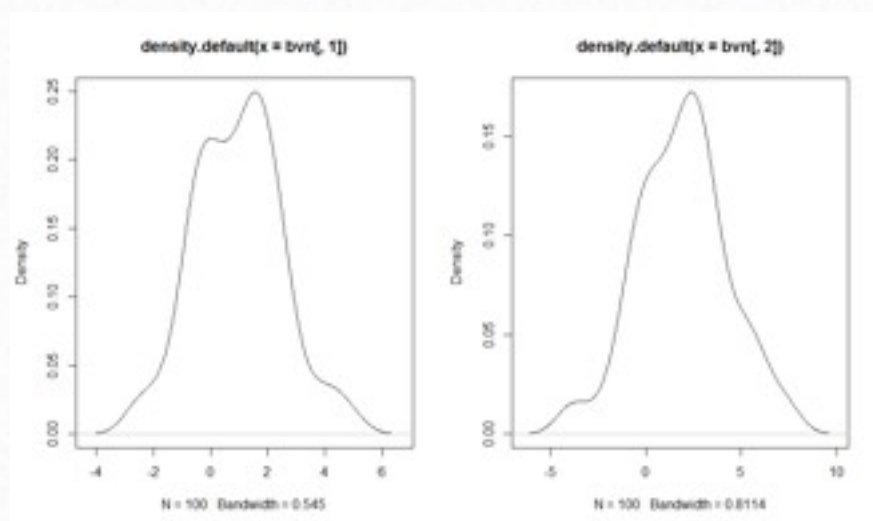


(정리) 만약 $Z_1, Z_2 \sim iid n(0, 1)$ 이면,

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \cdot \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_X}{\sigma_X}\right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 - 2\rho\frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y}\right]\right\}$$

$$X = \sigma_X Z_1 + \mu_X$$

$$Y = \sigma_Y[\rho Z_1 + \sqrt{1-\rho^2} Z_2] + \mu_Y \sim \text{BVN}$$



```
library(MASS) #mvnorm function package
Sigma <- matrix(c(2,3,3,5),2,2) #covariance matrix
bvn<-as.matrix(mvrnorm(n = 100, c(1, 2), Sigma))
par(mfrow=c(1,2)) #plot 2 in 1 row
plot(density(bvn[,1])) #marginal of X
plot(density(bvn[,2])) #marginal of Y
```

아변량 정규분포의 주변확률분포는 정규분포를 따른다.

2) 다변량 정규분포

$$\underline{x}_p = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

차수 p인 확률변수 벡터

$\underline{x}_p \sim N_p(\underline{\mu}_p, \Sigma_{p \times p})$, 평균이 $\underline{\mu}$, 공분산행렬이 Σ 인 다변량 정규분포의 확률분포함수는 다음과 같다.

$$f_{\underline{x}}(\underline{x}; \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp[-1/2(\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu})]$$

4. 상관계수와 산점도 [iBook강의]

다변량에서의 활용

similarity of variable : 상관계수가 높은 변수들은 유사하다는 의미임

- 두 변수가 유사하다는 것은 관심 개체에 대한 정보가 겹친다는 의미임
- (몸무게와 키) 두 변수의 상관계수가 매우 크므로 두 변수 중 하나의 관측값만 알아도 그 사람의 다른 관측값을 쉽게 예측할 수 있다. 사람들의 신체적 정보를 얻기 위하여 두 변수 모두 필요하는 것은 아니므로 변수의 차원을 줄일 수 있음
- (수학, 영어) 성적의 상관관계는 다소 낮음, 그러므로 수학성적만으로 영어성적을 예상하기는 쉽지 않

으므로 학생(개체)의 능력을 알려면 두 변수 모두 조사할 필요가 있음

5. 고유값 eigen value [iBook강의]

- 데이터 변수의 변동의 크기, 데이터의 공분산으로부터 구한 고유값은 원 변수들의 총변동의 크기를 반영함
- 변수의 변동(분산)은 관심 개체를 구별하는 능력 - 변수의 변동이 크다는 것은 개체를 구별할 수 있는 능력이 크다는 것임
- 예를 들어 H대학 통계학과 학생들의 고등학교 수학 성적과 영어성적을 조사한 자료이다.
- 수학평균=70점, 분산=20점 - 영어평균=80, 분산=5점
- 어느 과목을 이용하여 학생들의 능력을 구별하는 것이 좋을까? 당연히 변동(분산)이 큰 수학성적으로 학생들을 구별해야 용이하게 구별할 수 있음

6. 다변량 분석 요약

	주성분	요인	판별	군집	정준상관
변수 관계 탐색	S	D	N	N	S
자료 탐색	D	S	N	S	N
새 변수 만들기	Yes	Yes	No	No	Yes
개체 분류	No	No	Yes	Yes	No
변수 그룹	P	P	N	N	D
차원 줄이기	D	P	N	N	N

Sometimes, Definitely, Never, Possible, Rarely

다변량 분석 개념도



미국 59개 도시들의 기후지표, 사회 지표, 환경지표에 대한 데이터

	<- 기후지표->					<- 사회지표->					<- 환경지표->							
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	도시이름	1월기온	7월기온	상대습도	강수량	사망률지	교육기간	인구밀도	비백인비	사무직종	총인구	가구당인	가계소득	HCPot	NOxPot	S02Pot	NOx	남부여부
2	Akron, OH	27	71	59	36	921.87	11.4	3243	8.8	42.6	660328	3.34	29560	21	15	59	15	NO
3	Albany-S	23	72	57	35	997.87	11	4281	3.5	50.7	835880	3.14	31458	8	10	39	10	NO
4	Allentown	29	74	54	44	962.35	9.8	4260	0.8	39.4	635481	3.21	31856	6	6	33	6	NO
5	Atlanta, G	45	79	56	47	982.29	11.1	3125	27.1	50.2	2E+06	3.41	32452	18	8	24	8	YES
6	Baltimore	35	77	55	43	1071.3	9.6	6441	24.4	43.7	2E+06	3.44	32368	43	38	206	38	YES
7	Birmingham	45	80	54	52	1030.4	10.2	3225	28.5	43.1	883018	3.15	27835	30	32	72	32	YES

- 분석목적** (1) 59개 도시를 측정된 지표들을 이용하여 특성에 따른 분류
 (2) 측정된 지표들을 이용하여 “남부여부”를 판별하는 규칙(패턴) 탐색

- 1) 분석 (1) & (2)는 개체(도시)에 대한 분석임 -> 이를 개체 관련 기법 subject directed tech.
- 2) 특성 지표 3개분야, 총 15개 변수를 이용하여 도시의 특성을 파악하기는 변수의 수가 너무 많다. - 변수가 1개이면 기초통계량, 2개면 산점도, 3개면 버블산점도로 특성 파악 가능, 그러나 4개 이상이면 변수의 차원을 줄여야 도시 특성 파악이 가능하다.

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{21} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{n1} \end{bmatrix}$$

개체분류

- 군집분석 clustering analysis : 유사한 개체들을 동일 군집 분류
- 판별분석 discriminant analysis : 새로운 개체를 가장 적절한 특정 집단으로 분류

차원축약

- 주성분분석 principle component analysis : 변수들의 상관관계를 활용하여 변수 개수 축약 - 새로운 변수, 주성분 변수
- 요인분석 factor analysis : 유사한 변수들을 군집으로 분류