

Week 12. 주성분 분석

- 고차원의 데이터는 계산과 시각화가 어려우니 분석하기가 쉽지 않다.
- 따라서 원래 데이터의 분포를 가능한 유지하면서 데이터의 차원을 줄이는 것이 필요하다.
→ 이를 차원 축소 (dimensionality reduction)라 한다.

- 주성분 분석 (Principal Component Analysis)은 가장 널리 사용되는 차원 축소 기법 중 하나로, 원 데이터의 분포를 최대한 보존하면서 고차원 공간의 데이터를 저차원 공간 데이터로 변환한다.

12.1 차원 축소

1974년 Motor Trend Magazine, 자동차 특성 기술하고 경쟁력하는 16가지 변수.

	mpg	cyl	disp	hp	drat	wt	spc qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160	110	3.90	2.620	16.96	0	1	4	4
Mazda RX4 Wag	21.0	6	160	110	3.90	2.875	17.02	0	1	4	4
...

자동차 1대 → 11차원 벡터의 정보 고차원 데이터로 계산, 시각화 어렵다.

원 데이터의 분포를 가능한 유지하면서 데이터 차원을 줄이는 것이 필요하다. → 차원 축소

(?) 물론 일부 변수의 데이터만 뽑아서 사용할 수도 있지만 (feature selection)

변수들 사이에 어떠한 밀접한 관계가 있는지 미리 알기 어렵고, 원 데이터 분포를 훼손할 수 있다.

12.2 주성분 분석 (PCA)

Principal Component Analysis

원 분포를 최대한 보존하면서, 고차원 공간의 데이터를 저차원 공간으로 변환한다.

→ PCA는 기저 변수를 조합하여 서로 상관성이 없는 새로운 변수, 즉 주성분을 생성

→ 첫번째 주성분 PC1이 원데이터 분포를 가장 많이 보존하고, 두번째 PC2가 그 다음으로 원 데이터 분포를 많이 보존하는 식

→ 11차원 데이터를 PCA하여 1개의 주성분을 생성할 수 있으며,

PC1, PC2, PC3이 원데이터 분포(성질)의 약 90%를 보존한다면

10%의 손실을 감한해도 합리적 분석이 큰 의미가 있어 3차원으로 줄일 수 있다.

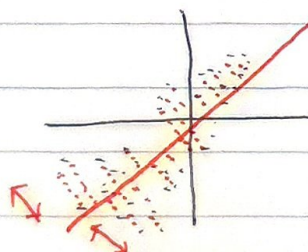
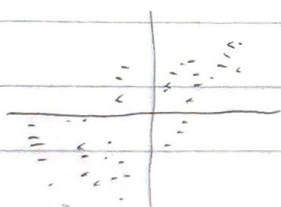
~ 계산/시각화 용이!

주성분을 구하기 위해 ⁶⁶ 새로운 축⁹⁹을 찾아야 함 → 주축 (principal axes)

→ 주방향 (principal directions)

ex) 2차원 데이터의 첫번째 주축을 찾아 정사영

→ 정사영된 데이터 (성분)들이 PC1을 이루게 된다.



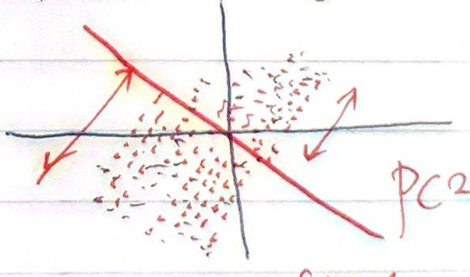
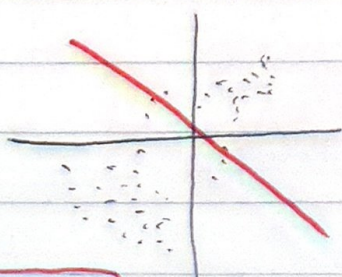
PC1 ∈ principal axis
(principal directions)

ex2) 두번째 축 — 원데이터를 정사영할 때, 세 데이터 (PC2) 분포가 PC1 다음으로 큰대.

* PC1과 PC2가 서로 관계가 없도록, **첫번째 축 ⊥ 두번째 축 (수직)**

→ 예상치 못한 변수 간 영향을 방지하기 위하여

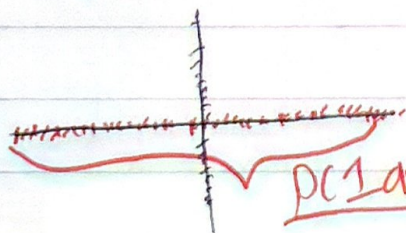
→ 서로 관계가 있는 변수끼리는 각 주성분이 모여도록 하되, 각 주성분들(축) 사이 관계성 제거



PC2

$$PC1 \text{ ax} = x \text{ 축}$$

$$PC2 \text{ ax} = y \text{ 축}$$



PC1 ax가 PC2 ax보다 분포가 넓다.

12.3 주성분 분석의 계산

$X = \text{data matrix}_{(n \times p)}$
 # data features. 4개는 확률변수 x_1, x_2, \dots, x_p 개수
 # sample 개수
 (동계학에서는 일반적으로 데이터 행렬 X 와 관련이 있어, 확률변수 x_1, x_2, \dots 대신에 소문자 x_1, x_2, \dots 를 쓴다.)

* 데이터 행렬 X 의 (i, j) 성분 x_{ij} 는 i 번째 표본의 j 번째 확률변수 x_j 에 대한 하나의 데이터를 의미하고, j 열 $X^{(j)}$ 는 확률변수 x_j 가 갖는 모든 데이터를 의미한다.

확률변수 x_1, x_2, \dots, x_p

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}_{n \times p}$$

$X^{(1)}$, x_1 의 모든 데이터

X 를 센터링 (centering, 확률변수의 평균을 0으로 조정) 하여, \tilde{X} 를 정의하자. 이를 위해

- ① X 의 각 열의 평균을 구한다.
 확률변수 x_j 의 평균은 \bar{x}_j 로 표기한다.

$$\bar{x}_1 = \frac{1}{n} (x_{11} + x_{21} + \dots + x_{n1})$$

- ② 각 열별로 데이터에서 열의 평균 \bar{x}_j 를 빼다. 이 행렬은 센터링된 행렬 \tilde{X} 라 한다.

$$\tilde{X} = X - \begin{bmatrix} \bar{x}_1 & \dots & \bar{x}_p \\ \vdots & \ddots & \vdots \\ \bar{x}_1 & \dots & \bar{x}_p \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} - \begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \end{bmatrix}$$

여기서 센터링된 행렬 \tilde{X} 의 각 열의 평균은 0이 된다. \rightarrow mean-centered matrix
 앞으로 주어진 행렬 X 는 '이미 센터링이 된' 행렬이라고 가정하자.

- ③ X 의 특이값 분해 (SVD)를 구한다.

$$X = [u_1 \ u_2 \ \dots \ u_k \ u_{k+1} \ \dots \ u_n] \begin{bmatrix} s_1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_k & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & s_{k+1} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & s_p \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_p^T \end{bmatrix}$$

U, V - 직교행렬 ($U^T U = I, V^T V = I$)
 $s_i =$ 크기순 배열된 특이값 $s_1 \geq s_2 \geq \dots \geq 0$
 을 주대각선 성분으로 하는 $p \times p$ 대각선행렬 $S = \text{diag}(s_1, \dots, s_p)$

$$= USV^T = U \begin{bmatrix} s_1 \\ 0 \end{bmatrix} V^T = \sum_{i=1}^p s_i u_i v_i^T$$

④ V 의 열벡터 $\begin{bmatrix} V_1^T \\ V_2^T \\ \vdots \\ V_p^T \end{bmatrix}$ 가 주축 (principal axes) 이 된다.
 → p 개의 주축이 생성 가능하다

⑤ $Z = US$ 의 열벡터들이 원 데이터를 주축에 정사영하여 얻어진 주성분 점수 (principal component score, PC score)가 된다.

⑥ 특이값 S_i 를 이용하여 계산한 $\frac{S_i^2}{n}$ 은 i 번째 PC의 분산 이고,
 $\frac{S_i^2}{S_1^2 + S_2^2 + \dots + S_p^2}$ 은 i 번째 PC가 원 데이터의 분포를 보존하는 비율이다.
 각 PC별로 원 데이터의 분포를 보존하는 비율을 계산하여, 몇차원으로 축소할지 결정한다.

⑦ 데이터를 p 차원까지 k ($k \leq p$)차원으로 줄이기 위하여, U 의 처음 k 개의 열벡터 (U_k)와 S 의 k 번째 선행 주 부분행렬 (leading principal submatrix) (S_k)을 택하면,
 $U_k S_k$ 는 처음 k 개의 PC를 포함하는 $n \times k$ 행렬이 된다.

$$Z_k = U_k S_k = [u_1, u_2, \dots, u_k] \begin{bmatrix} s_1 & s_2 & \dots & s_k \\ & & & \\ & & & \\ & & & s_k \end{bmatrix}$$

12.4 주성분 분석 사례: 4차원 데이터 → 2차원 데이터

※ 센터링 진행

for col in range(p):

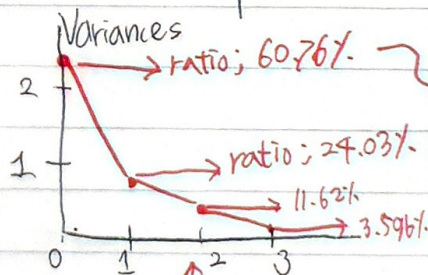
$v = np.array(X.column(col))$

$m = \text{mean}(v)$

$st = \text{std}(v)$

for row in range(n):

$X_ctr[row, col] = (X[row, col] - m) / st$



PC1, PC2
 분산비율합
 84.79%

||
 PC1, PC2
 원 분포
 보존율 84.79%

특이값분해

$U, S, V = X_ctr.SVD()$

대각성분 특이값 → 분산

$\text{Var-PC} = [(S[i, i]^2 / n), n (\text{digits}=4) \text{ for } i \text{ in range}(p)]$

i 번째 PC가 원 분포 보존 비율

$\text{Prop-Var} = [100 * \text{Var-PC}[i] / \text{sum}(\text{Var-PC}) \text{ for } i \text{ in range}(p)]$

p 개의 주성분 열 차 행렬

$\text{PC-Score} = U * S$

12.5 주성분 분석과 공분산행렬

주성분 분석은 원 데이터의 분포를 최대한 보존하면서 고차원 데이터 \rightarrow 저차원 변환.
 But 원 데이터 분포에 대한 정보는 공분산행렬이 담겨 있다. 왜냐하면 공분산행렬은
 주대각선 상의 성분은 각 확률변수가 얼마나 퍼져 있는지를 나타내는 **분산**과 주대각선
 이외의 성분이 **확률변수 간 상관관계**를 나타내는 **공분산**으로 되어 있기 때문이다.

데이터 행렬 X 가 선형 (mean-centered) 된 행렬이면, $p \times p$ 공분산행렬 Σ 은
 다음과 같이 계산한다.

$$\Sigma = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_p) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \dots & \text{Cov}(x_2, x_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_p, x_1) & \text{Cov}(x_p, x_2) & \dots & \text{Var}(x_p) \end{bmatrix}$$

$$= \frac{1}{n} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pp} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n x_{i1}^2 & \frac{1}{n} \sum_{i=1}^n x_{i1} x_{i2} & \dots \\ \frac{1}{n} \sum_{i=1}^n x_{i2} x_{i1} & \frac{1}{n} \sum_{i=1}^n x_{i2}^2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

비대각 성분 = 각 확률변수 간 공분산 (상관관계)

대각 성분 = 각 확률변수 x_{np} 의 분산

$$= \frac{1}{n} X^T X$$

주성분 성분의 목적은 '공분산행렬로부터 얻은 정보'를 최대한 보존하는 '더 적은 개수의 새 변수들'을 찾는 것이다.

$X = USV^T = U \begin{bmatrix} S_1 \\ 0 \end{bmatrix} V^T$ 로부터 다음 관계를 얻는다.

$$\text{공분산 행렬 } \Sigma = \frac{1}{n} X^T X = \frac{1}{n} V \begin{bmatrix} S_1^T & 0 \end{bmatrix} U^T U \begin{bmatrix} S_1 \\ 0 \end{bmatrix} V^T = V \left(\frac{S_1^2}{n} \right) V^T$$

따라서 $\Sigma V = V \left(\frac{S_1^2}{n} \right)$, 즉 $\Sigma V_i = \left(\frac{S_1^2}{n} \right) V_i \quad (i=1, 2, \dots, p)$ 관계가 성립한다.

$$= \lambda_i \cdot V_i = \text{고유값} \cdot \text{고유벡터}$$

1번째 PC의 분산

$\lambda_i = \frac{S_i^2}{n}$
 $\rightarrow \Sigma$ 의 고유값
 (eigenvalue)

$V_i \rightarrow \Sigma$ 의 고유벡터

1번째 주축
 (Principal Axis)

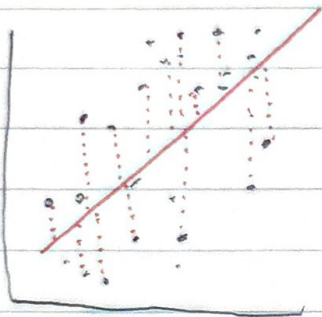
12.6 주성분 분석과 선형회귀

(그림 1) 선형회귀에서는 각 데이터 (x_i, y_i) 와 일차함수 $y = a + bx$ 에 의하여 계산된 점 (x_i, \hat{y}_i) 사이의 거리 (즉 y 축 상에서의 직선거리)가 최소가 되도록 a, b 를 결정

(그림 2) 주성분 분석에서는 각 데이터 (x_i, y_i) 에서 일차함수 $y = a + bx$ 에 이르는 거리가 최소 되도록 a 와 b 를 결정 ... 그래야 가장 많은 분포를 얻을 수 있다.

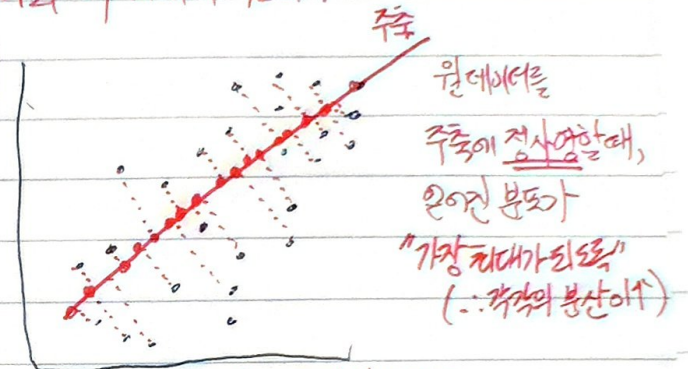
⇒ PCA로 얻어진 직선이 데이터와 더 가까이 있다.

[그림 1]



$\min (y - \hat{y})$

[그림 2]



* PCA는 data scale이 민감하므로
mean-centered, normalization (정규화)
를 반드시 가정해야 한다.