

1. 개념 정리

$\left\{ \begin{array}{l} \text{변수 선택} : \text{특정 변수를 선택하는 것 (하위 집합)} \\ \text{변수 추출} : \text{변수로 새로운 변수 추출 (SVD, PCA 등)} \end{array} \right.$

Supervised Feature Selection : IG, Stepwise, U12, ...

Feature Extraction : Partial Least Square

Unsupervised Feature Selection : PCA loading

Feature Extraction : PCA, Wavelets, AutoEncoder

2. PCA

n 개 관측치와 p 개 변수로 구성된 데이터 \rightarrow 주요한 k 개 변수로 구성된 데이터로 요약
 (주요한 변수의 선형 조합)

원래 데이터 분포를 최대한 보존하는 축을 찾아 데이터 사영.

\rightarrow 차원 축소, 시각화, 해석 \sim 초기 단계 사용.

$$Z_p = \alpha_p^T X = \alpha_{p1} X_1 + \alpha_{p2} X_2 + \dots + \alpha_{pp} X_p \quad \left\{ \begin{array}{l} X : \text{원래 변수} \\ \alpha_i : \text{변환 Loading} \\ Z_p : \text{사영 변환 후 변수} \end{array} \right.$$

3. 수식 배경

Covariance $Cov(X) = \frac{1}{n} (X - \bar{X})(X - \bar{X})^T$

공분산 행렬 대칭성 = p 변수의 분산 & 비모양 = 대응 단위들의 cov

데이터의 총 분산은 공분산 행렬 대칭성의 원 $\rightarrow [Cov(X)] = Cov_{11} + Cov_{22} + \dots + Cov_{pp}$

4. Projection

$$\vec{x} = p\vec{a} \quad p = \frac{\vec{b}^T \vec{a}}{\vec{a}^T \vec{a}}$$

Eigen Value & Eigenvector

$$Ax = \lambda x \quad \text{일 때 } \lambda : \text{고유값}$$

$$(A - \lambda I)x = 0 \quad \left\{ \begin{array}{l} x : \text{고유벡터} \sim \text{선형 변환에 의해 방향이 변하지 않는 벡터} \end{array} \right.$$

4. Algorithm ($\vec{v}_2 = 0$ 인 데이터가 있음)

to find α that produces the largest variance of Z \rightarrow 가능한 최대 분산

$$\begin{aligned} \text{Max } \text{Var}(Z) &= \text{Var}(X^T \alpha) \\ &= \alpha^T \text{Var}(X) \alpha = \alpha^T \underbrace{\text{Cov. Matrix}}_{= X^T X} \alpha \end{aligned}$$

$$= \alpha^T E \Lambda E^T \alpha$$

Eigenvalue Matrix

$$p = E^T \alpha \quad \|p\| = 1$$

$$\text{Max } \beta^T \Lambda \beta = \text{Max } \lambda_1 \beta_1^2 + \lambda_2 \beta_2^2 + \dots + \lambda_m \beta_m^2$$

$$\text{s.t. } \beta_1^2 + \beta_2^2 + \dots + \beta_m^2 = 1, \lambda_1 > \lambda_2 > \dots > \lambda_m$$

$$\Rightarrow \beta_1 = 1 \text{ 이 때 } \text{Max, optimal value} = \lambda_1, \alpha = e_1$$

$$[1 \ 0 \ 0 \ \dots \ 0]$$

① 데이터 Covariance Matrix (변분 공분산 행렬)

② 공분산 cov matrix 찾기.

③ Cov matrix \rightarrow eigenvalue & vector 찾기

④ Eigenvalue 및 해당 Eigenvector 순서대로 내림.

⑤ 정렬된 Eigenvector 기반 원본 변수 변환.

$\Rightarrow \beta_1 = 1$ 이라 하면 Max, optimal value = λ , $\lambda = 0$,
 $[1 \ 0 \ 0 \ \dots \ 0]$

5. PCA

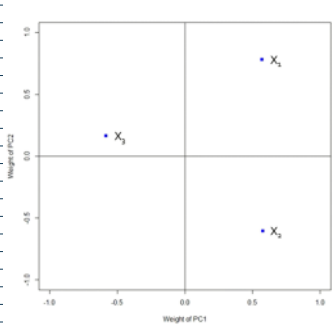
공분산행렬의 Eigen Value = 주성분의 분산

ex) $Cov(z) = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ $Var(z_1) = 2 = \lambda_3$ $\lambda_1, \lambda_2, \lambda_3$ 중 가장 큰 값을 λ_1 로 놓는다. $\lambda_1 = 2$
 $Var(z_2) = 1 = \lambda_2$ $\lambda_1, \lambda_2, \lambda_3$ 중 가장 큰 값을 λ_2 로 놓는다. $\lambda_2 = 1$
 $Var(z_3) = 1 = \lambda_3$ $\lambda_1, \lambda_2, \lambda_3$ 중 가장 큰 값을 λ_3 로 놓는다. $\lambda_3 = 1$

주성분 2개 선택

Elbow point / 설명 비율을 보거나

6. PCA Loading Plot



: 주성분 2개에 대해 가장 큰 영향을 미치는 변수를 나타낸다.

X_1 : PC1, PC2 모두 양의 영향을 받는다.

X_2 : PC1은 음의 영향을 받는다, PC2는 양의 영향을 받는다.

X_3 : PC1은 음의 영향을 받는다, PC2는 음의 영향을 받는다.

7. t-SNE

Non Gaussian / 다른 Gaussian 차이를 잘 표현해준다.

\rightarrow 커널 PCA, LE (Locally Linear Embedding)

비선형 관계를 고려, 비선형 관계를 잘 표현해준다.