

머신러닝 차원 축소

선형대수학 응용 예제

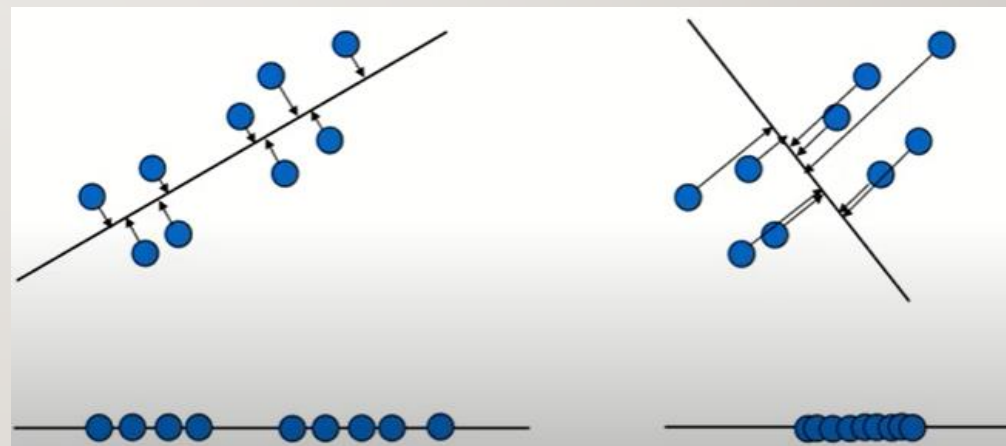
202021224 주민찬



주성분분석(PCA)

- Principal Component Analysis

1. 가장 기본적인 차원 축소 기법
2. 데이터를 축에 사영했을 때 가장 높은 분산(정보 손실 최소화)을 가지는 데이터의 축(주성분)을 찾아 그 축으로 차원을 축소하는 것



데이터 정규화

1	2	3
4	5	6
7	8	9

$$\begin{array}{lll} m_1 = 4 & m_2 = 5 & m_3 = 6 \\ Sd_1 = \sqrt{6} & Sd_2 = \sqrt{6} & Sd_3 = \sqrt{6} \end{array}$$

$X =$

$\frac{1 - 4}{\sqrt{6}}$	$\frac{2 - 5}{\sqrt{6}}$	$\frac{3 - 6}{\sqrt{6}}$
$\frac{4 - 4}{\sqrt{6}}$	$\frac{5 - 5}{\sqrt{6}}$	$\frac{6 - 6}{\sqrt{6}}$
$\frac{7 - 4}{\sqrt{6}}$	$\frac{8 - 5}{\sqrt{6}}$	$\frac{9 - 6}{\sqrt{6}}$

공분산 행렬

$$X = \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array}$$

$$X^T = \begin{array}{|c|c|c|} \hline 1 & 4 & 7 \\ \hline 2 & 5 & 8 \\ \hline 3 & 6 & 9 \\ \hline \end{array}$$

$$X^T X = \begin{array}{|c|c|c|} \hline 1 \cdot 1 + 4 \cdot 4 + 7 \cdot 7 & 1 \cdot 2 + 4 \cdot 5 + 7 \cdot 8 & 1 \cdot 3 + 4 \cdot 6 + 7 \cdot 9 \\ \hline 2 \cdot 1 + 5 \cdot 4 + 8 \cdot 7 & 2 \cdot 2 + 5 \cdot 5 + 8 \cdot 8 & 2 \cdot 3 + 5 \cdot 6 + 8 \cdot 9 \\ \hline 3 \cdot 1 + 6 \cdot 4 + 9 \cdot 7 & 3 \cdot 2 + 6 \cdot 5 + 9 \cdot 8 & 3 \cdot 3 + 6 \cdot 6 + 9 \cdot 9 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|} \hline 66 & 78 & 90 \\ \hline 78 & 93 & 108 \\ \hline 90 & 108 & 126 \\ \hline \end{array}$$

e.value = $\lambda_1, \lambda_2, \lambda_3$ ($\lambda_1 > \lambda_2 > \lambda_3$)

e.vector v_1, v_2, v_3 corresponding to $\lambda_1, \lambda_2, \lambda_3$

주성분 선택 및 주성분으로 투영

원본 데이터

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

e.vector

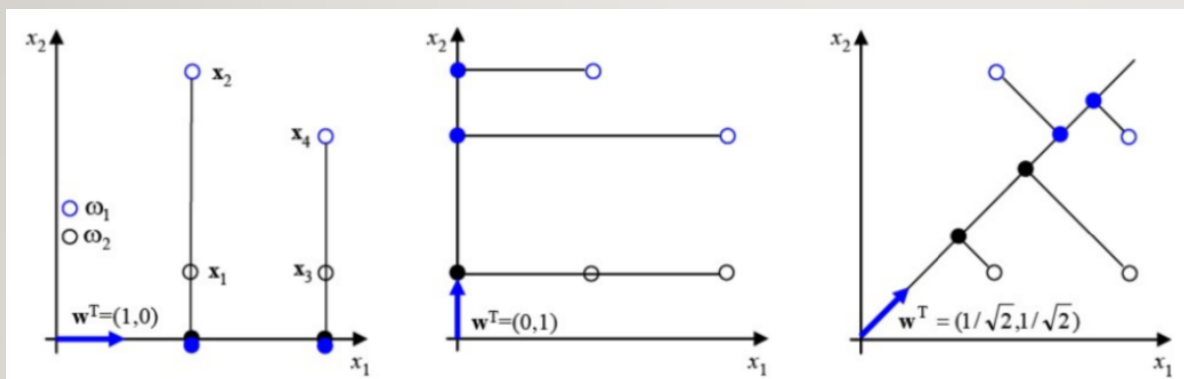
$$v_1, v_2$$

2차원으로 변환

$$\begin{bmatrix} \text{proj}_{v_1} x_1 & \text{proj}_{v_2} x_1 \\ \text{proj}_{v_1} x_2 & \text{proj}_{v_2} x_2 \\ \text{proj}_{v_1} x_3 & \text{proj}_{v_2} x_3 \end{bmatrix}$$

PCA의 한계점

- PCA의 주요 한계점으로 최대의 분산의 각 축이 반드시 클래스 간의 구별을 잘하는 좋은 피처를 뽑아준다는 보장이 없다는 점



- 그 차원을 축소하는 데 있어 클래스 간의 차별성을 최대화할 수 있는 방향으로 수행하는 것이 LDA입니다.