

3.5 PCA

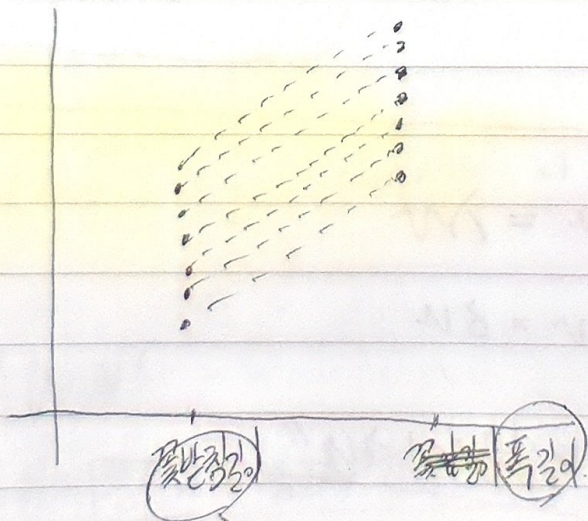
N개의 N차원 데이터가 있으면 보통 그 데이터들은 서로 다른 값을 가진다. 하지만 이러한 데이터간의 변이 (variation)는 무작위가 아니라 특정한 규칙에 따라 만들어지는 경우가 있다. 예를 들어, 꽃잎의 꽃받침 길이는 꽃다라 다르지만, 꽃받침 길이가 약 2배 커지면 꽃받침 폭도 약 2배 커지는 것이 일반적이다. 이러한 데이터간의 변이 **규칙을 찾아낼 때 PCA를 활용할 수 있다.**

PCA (Principal Component Analysis)

= 주성분분석

= 고차원 집합이 주어졌을 때 원래의 고차원 데이터와 가장 비슷하면서 더 낮은 차원의 데이터를 찾아내는 방법 (=차원 축소)

"더 낮은 차원의 데이터가 더 높은 차원의 데이터를 설명할 수 있다"
→ 몇 가지 원인으로 데이터 변이를 설명할 수 있다.



변이와 상관관계 문제 (규칙)

- "꽃의 크기"라는 근본적인 데이터가 "꽃받침 길이", "꽃잎 폭"으로 표현된 것
- 측정되지는 않지만 측정 데이터 자체에 숨겨져 측정 데이터를 결정짓는 데이터
→ 잠재변수 (latent variable)

PCA의 가정

- 장재변수와 측정 데이터가 선형적인 관계를 연결되어 있다고 가정함.
- 1번째 품목의 측정 데이터 벡터 x_i 의 각 원소를 선형조합하여
 2 뒤에 숨은 1번째 품목의 장재변수 u_i 의 값을 계산할 수 있다고 가정한다.

$$(u_i = w^T x_i)$$

- 이 식에서 w 는 측정 데이터 벡터의 각 원소를 조합한 가중치 벡터다.
 맛의 예에서는 골분침임과 골분침 폭을 선형조합하여
 맛의 크기를 나타낼 때의 값을 찾는 것이라 생각할 수 있다.

$$u_i = w_1 x_{i,1} + w_2 x_{i,2}$$

(latent)

(ex) 미식축구 쿼터백의 패시 레이트 (passer rate)

→ 쿼터백의 "실력"이라는 장재변수

$$\text{passing rate} = 5 \cdot \frac{\text{completions}}{\text{attempts}} + 0.25 \frac{\text{passing yards}}{\text{attempts}} + 20 \cdot \frac{\text{touchdowns}}{\text{attempts}} - 25 \frac{\text{interceptions}}{\text{attempts}} + 0.125$$

⇒ 차원축소 관점에서 보면 4차원 데이터를 1차원으로 줄인 것.

⇒ 2차원 데이터의 latent variable을 찾아 간단하게 만들자
 = PCA.

차원축소의 특징

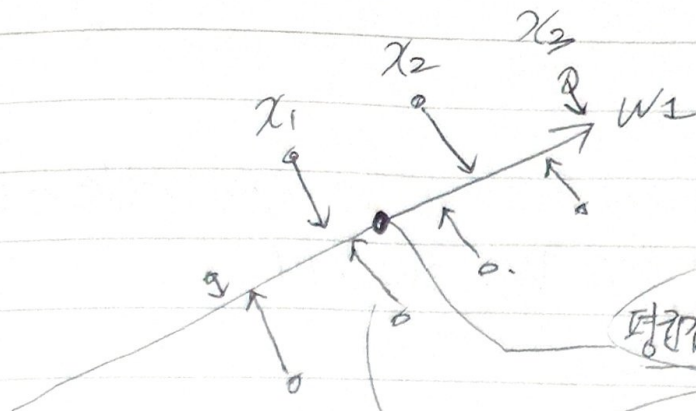
차원축소 문제는 다차원 벡터를 더 낮은 차원의 벡터공간에 투영하는 문제로 생각하여 풀 수 있다. 즉, 특정 차원 이하에서 실행된 근사 (low-rank approximation) 문제가 된다. 이 문제를 다음과 같이 서술할 수 있다.

N 개의 M 차원 데이터 벡터 x_1, x_2, \dots, x_N 을 정제집근인 기저벡터 w_1, w_2, \dots, w_k 로 구성된 k 차원 벡터공간으로 투영하여 가장 비슷한 N 개의 k 차원 벡터 $x_1^{||w}, x_2^{||w}, \dots, x_N^{||w}$ 를 만들기 위한 정제집근 기저벡터 w_1, w_2, \dots, w_k 를 찾는다.

다만, 원래의 근사 근사 문제와 달리, 근사 성능을 높이기 위해 직선이 원점을 지나야 한다는 제한조건을 있어야 한다. 따라서 문제는 다음과 같이 바뀐다.

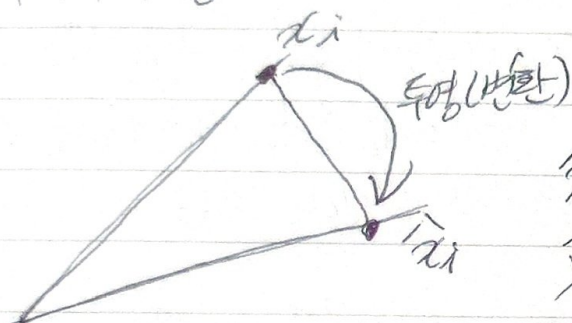
N 개의 M 차원 벡터 x_1, x_2, \dots, x_N 에 대해 어떤 상수 벡터 x_0 를 뺀 데이터 벡터 $x_1 - x_0, x_2 - x_0, \dots, x_N - x_0$ 을 정제집근인 기저벡터 w_1, w_2, \dots, w_k 로 구성된 k 차원 벡터공간으로 투영하여 가장 비슷한 N 개의 k 차원 벡터 $x_1^{||w}, x_2^{||w}, \dots, x_N^{||w}$ 를 만들기 위한 정제집근 기저벡터 w_1, w_2, \dots, w_k 와 상수 x_0 를 찾는다.

(단) x_0 은 데이터 벡터 x_1, x_2, \dots, x_N 의 평균벡터이고, w_1, w_2, \dots, w_k 는 가장 큰 k 개의 특잇값에 대응하는 오른쪽 특이벡터 v_1, v_2, \dots, v_k 이다.



$x_i \perp w_1$ 를 최소화하는 w_1 을 찾는다.

PCA의 수학적 설명



$$\hat{x}_i = W x_i$$

$$\hat{X} = X W^T$$

$$\begin{bmatrix} \hat{x}_1 \\ \vdots \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \end{bmatrix} \begin{bmatrix} w_1 \end{bmatrix}$$

- 데이터가 원점을 중심으로 존재하는 경우를 벡터에 변환행렬을 곱하는 연산을 두영(변환)을 계산할 수 있다. (설명 편위를 위해 데이터가 원점에 모여있다고 가정)
- ~~데이터가 원점을 중심으로 존재하는 경우를~~ 다음처럼 데이터 x_i 에 변환행렬 $W \in \mathbb{R}^{k \times M}$ 을 곱해서 새로운 데이터 \hat{x}_i 를 구하는 연산을 생각하자.

$$\hat{x}_i = W x_i$$

$$x \in \mathbb{R}^M, W \in \mathbb{R}^{k \times M}, \hat{x} \in \mathbb{R}^k$$

모든 데이터 x_i ($i=1, \dots, N$)에 대해 변환을 하면 벡터가 아닌 행렬로 표현 가능하다.

$$\hat{X} = X W^T$$

$$X \in \mathbb{R}^{N \times M}, \hat{X} \in \mathbb{R}^{N \times k}, W^T \in \mathbb{R}^{M \times k}$$

< PCA 목표 >

차원 축소 벡터 \hat{x} 가 원래 벡터 x 가 지녔던 정보와 가장 유사하게 만드는 행렬 W 를 찾는 것이다.

또한 \hat{x} 는,

$$\hat{x} = Ux$$

$$x \in \mathbb{R}^K, U \in \mathbb{R}^{M \times K}, \hat{x} \in \mathbb{R}^M$$

이렇게 변환과 역변환을 통해 원래의 차원으로 되돌린 벡터 $U\hat{x}$ 은 원래의 벡터 x 와 비슷할수록 장점이 갈수록 많다. 다만 이 값을 다시 한번 차원 축소 변환하면 또 \hat{x} 가 된다. 즉,

$$W\hat{x} = WUx = x$$

따라서 W 와 U 는 다음 관계가 있다.

$$WU = I$$

변환 행렬 U 를 알고 있다고 가정하고 변환했을 때, 원래 벡터 x 와 가장 비슷해지는 차원 축소 벡터 \hat{x} 를 다음과 같이 최적화를 이용하여 찾는다.

$$\arg \min_x \|x - U\hat{x}\|^2$$

목적함수는 다음과 같이 바꿀 수 있다.

$$\begin{aligned}\|x - U\hat{x}\|^2 &= (x - U\hat{x})^T (x - U\hat{x}) \\ &= x^T x - \hat{x}^T U^T x - x^T U \hat{x} + \hat{x}^T U^T U \hat{x} \\ &= \boxed{x^T x - 2x^T U \hat{x} + \hat{x}^T \hat{x}}\end{aligned}$$

이 목적함수를 최소화하려면 줄로 미분 식이 영벡터가 되는 값을 찾아야 한다.
이 부분은 행렬의 미분과 선형대수 부분에서 배워게 된다. 위 목적함수를 미분한
식은 다음과 같다.

$$\boxed{-2U^T x + 2\hat{x} = 0}$$

이 식을 정리하면

$$\hat{x} = U^T x$$

가 된다. 원래의 변환식

$$\hat{x} = Wx$$

라 비교하면,

$$U = W^T \quad (\because U^T x = Wx)$$

임을 알 수 있다. 따라서 다음식이 성립한다.

$$WW^T = I \quad (\because WU = I)$$

남은 문제는 행렬의 변환 행렬 W 를 찾는 것이다. 이 경우의 최적화 문제는
다음과 같이 된다.

$$\boxed{\arg \min_W \sum_{i=1}^N \|x_i - W^T W x_i\|^2}$$

모든 데이터에 대해 적용하면 목적함수는 다음처럼 된다.

$$\arg \min_x \|X - XW^TW\|^2$$

→ 이 문제는 rank-k 근사문제이고,
W는 행은 K개의 특징값에 대응하는
고유값 특이벡터로 만들어진 행렬이다.

사이킷런의 PCA 기능

사이킷런의 decomposition 시브패키지는 PCA 분석을 위한
PCA 클래스를 제공한다.

input parameter

- n_components (정수)

method

- fit_transform() : 특징행렬을 낮은 차원의 근사행렬로 변환
- inverse_transform() : 변환된 근사행렬을 원래차원으로 환원

attribute

- mean_ : 평균벡터
- components_ : 주성분벡터