

I.11 Norms of Vectors and Functions and Matrices

■ Norm이란?

□ 벡터의 norm: 벡터의 크기

- $\|\mathbf{x}\|_2 = (|x_1|^2 + \dots + |x_n|^2)^{1/2} \rightarrow$ Euclid norm 또는 l^2 -norm
- $\|\mathbf{x}\|_\infty = \max |x_i| \rightarrow l^\infty$ -norm (max norm)
- $\|\mathbf{x}\|_1 = |x_1| + \dots + |x_n| \rightarrow l^1$ -norm
- Preference: $\|\mathbf{x}\|_2, \|\mathbf{x}\|_1, \|\mathbf{x}\|_\infty$
- $\|\mathbf{x}\|_2$ 의 문제점: 작은 요소가 너무 작아지는 문제가 있음

□ 행렬의 norm

- $\|A\|_F = (\sigma_1^2 + \dots + \sigma_r^2)^{1/2} \rightarrow$ Frobenius norm
- $\|A\| = \max \sigma_i$
- $\|A\| = \sigma_1 + \dots + \sigma_r$

I.11 Norms of Vectors and Functions and Matrices

- Orthogonal invariance (대각 행렬에 대한 불변성)
 - 직교 행렬 Q 로 기저 변환 (또는 좌표축 변환)을 했을 때 변하지 않는 값들
 - 벡터의 길이: $(Q\mathbf{x})^T Q\mathbf{x} = \mathbf{x}^T \mathbf{x}$
 - 행렬 $A = U\Sigma V^T$ 의 singular value σ : $Q_1 A Q_2 = Q_1 U \Sigma V^T Q_2$
 - 직교 행렬의 곱은 직교 행렬임
 - Spectral norm $\|A\|_2 = \max \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \sigma_1$
 - Nuclear norm $\|A\|_N = \sum \sigma_i$
 - Frobenius norm $\|A\|_F = (\sigma_1^2 + \dots + \sigma_r^2)^{1/2} = \sqrt{\text{tr}(A^T A)}$
 - 단위 행렬(I)의 경우
 - $\|I\|_2 = 1$
 - $\|I\|_N = n$
 - $\|I\|_F = \sqrt{n}$

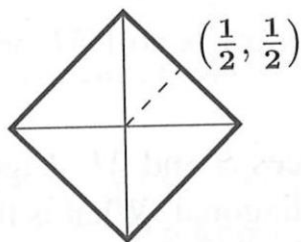
I.11 Norms of Vectors and Functions and Matrices

■ 벡터 norm의 성질

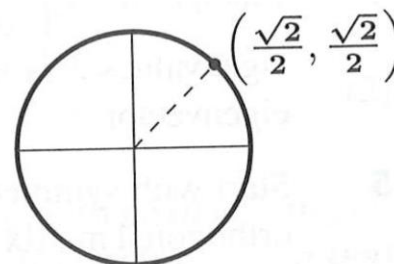
- 아래 3가지 종류의 벡터 norm에 적용되는 공통적인 성질은?
 - $\|\mathbf{x}\|_2 = (|x_1|^2 + \cdots + |x_n|^2)^{1/2} \rightarrow$ Euclid norm 또는 l^2 -norm
 - $\|\mathbf{x}\|_\infty = \max |x_i| \rightarrow l^\infty$ -norm (max norm)
 - $\|\mathbf{x}\|_1 = |x_1| + \cdots + |x_n| \rightarrow l^1$ -norm
- **0**-벡터를 제외한 모든 벡터 \mathbf{x} 에 대해 $\|\mathbf{x}\| > 0$
- $\|c\mathbf{x}\| = |c|\|\mathbf{x}\|$
- $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$: triangular inequality
- $\|\mathbf{x}\|_{1/2} = \left(\sqrt{|x_1|} + \sqrt{|x_2|} \right)^2$ 은 벡터 norm인가?
- 좋은 norm의 기준
 - 평면내 $\|\mathbf{x}\| \leq 1$ 를 만족하는 영역을 찾았을 때 영역의 모양이 Convex인 경우

I.11 Norms of Vectors and Functions and Matrices

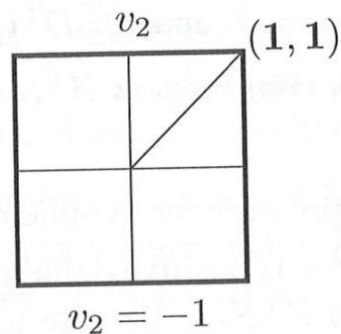
ℓ^1 norm
 $|v_1| + |v_2| \leq 1$
 diamond



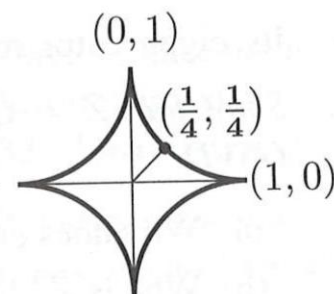
ℓ^2 norm
 $v_1^2 + v_2^2 \leq 1$
 circle



ℓ^∞ norm
 $|v_1| \leq 1, |v_2| \leq 1$
 square



$\ell^{1/2}$ norm
 $\sqrt{|v_1|} + \sqrt{|v_2|} \leq 1$
 not convex



I.9 PCA and the Best Low Rank Matrix

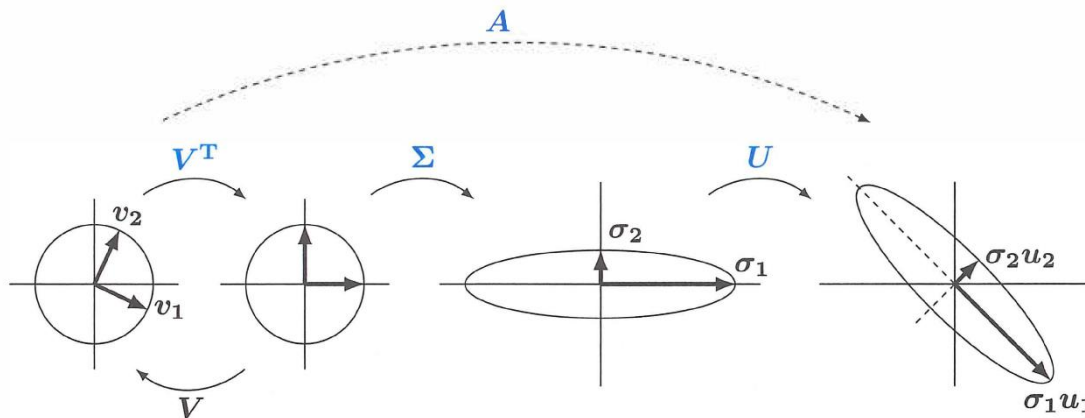
- 행렬 A 의 A_k
 - 행렬 A 가 주어졌을 때 A_k 는 A 를 SVD를 이용하여 rank-1 행렬들의 합으로 나타냈을 때, σ_1 부터 σ_k 까지 해당하는 항들만의 합을 나타냄

$$A_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

- A_k 의 rank는 당연히 k 임
- Eckart-Young 정리
 - 만약 행렬 B 의 rank가 k 이면, $\|A - B\| \geq \|A - A_k\|$ 이 항상 성립한다.
 - Rank가 k 인 행렬 중 행렬 A 에 가장 가까운 행렬은 A_k 이다.

I.8 Singular Values and Singular Vectors in SVD

- The singular vectors v_i
- $\frac{\|Ax\|_2}{\|x\|_2}$ 의 최대값
 - $\frac{\|Ax\|_2}{\|x\|_2}$ 은 벡터 $x = v_1$ 일 때 최대가 되고, 이 때 $\frac{\|Ax\|_2}{\|x\|_2} = \sigma_1$ 값을 가진다.
 - SVD를 아래 그림과 같이 입력 벡터의 회전 \rightarrow scaling \rightarrow 회전의 관점으로 보면 특이값(singular value)이 최대일 때, 비율이 최대가 되는 것은 당연해 보임.



I.8 Singular Values and Singular Vectors in SVD

- (증명) $\frac{\|Ax\|_2}{\|x\|_2}$ 은 벡터 $x = v_1$ 일 때 최대가 되고, 이 때 $\frac{\|Ax\|_2}{\|x\|_2} = \sigma_1$ 값을 가진다.
- 어떤 함수의 최대값 또는 최소값은 그 함수의 기울기가 모든 변수들에 대해 0이 되는 지점에서 일어남
- $\frac{\|Ax\|}{\|x\|} = \frac{x^T A^T A x}{x^T x} = \frac{x^T S x}{x^T x} \rightarrow$ 대칭행렬의 경우 이런 형태의 비율을 Rayleigh quotient라고 부름
- 위 함수는 x_1, \dots, x_n 의 함수이므로, 각각의 변수들에 대한 미분을 계산.

$$\frac{\partial}{\partial x_i} (x^T x) = \frac{\partial}{\partial x_i} (x_1^2 + \dots + x_n^2) = 2(x)_i$$

$$\frac{\partial}{\partial x_i} (x^T S x) = \frac{\partial}{\partial x_i} \left(\sum_j \sum_k S_{jk} x_j x_k \right) = 2 \sum_k S_{ik} x_k = 2(Sx)_i$$

- $\frac{\partial}{\partial x_i} \left(\frac{x^T S x}{x^T x} \right)$ 의 분모는 $x \neq 0$ 인 경우 0보다 크므로, 분자가 0이 되는 경우만 고려 $\rightarrow \left\{ \frac{\partial}{\partial x_i} (x^T S x) \right\} (x^T x) - (x^T S x) \frac{\partial}{\partial x_i} (x^T x) = 2(Sx)_i (x^T x) - (x^T S x) 2(x)_i = 0$
- $\{(x^T x)S\}x - (x^T S x)x = 0 \rightarrow Sx = \frac{(x^T S x)}{(x^T x)} x \rightarrow$ 벡터 x 가 eigenvector일 때 기울기가 0이 됨.
- 벡터 x 가 eigenvector v_i 일 때는 $\frac{\|Ax\|}{\|x\|} = \sqrt{\lambda_i} = \sigma_i$ 이므로 특이값이 최대인 σ_1 이 $\frac{\|Ax\|}{\|x\|}$ 의 최대값이 됨.

I.9 PCA and the Best Low Rank Matrix

- Spectral norm: $\|A\|_2 = \max \frac{\|Ax\|_2}{\|x\|_2} = \sigma_1$ (또는 l^2 -norm으로도 불림)
- Eckart-Young 정리의 예 (l^2 -norm의 경우)
 - 주어진 행렬 A 가 대각 행렬인 경우, rank-2를 가진 행렬 B 중 $\|A - B\|_2$ 를 최소로 하려면 A_2 가 당연함

- $A = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ 과 $A_2 = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$ 의 비교

- l^2 -norm의 값은 임의의 대각 행렬 Q_1, Q_2 에 의한 변환에 불변하므로, 위의 예는 4,3,2,1의 특이값을 가지는 임의의 4x4행렬에 대해서 적용됨

I.9 PCA and the Best Low Rank Matrix

- Eckart-Young 정리 (l^2 -norm의 경우)

- $\text{rank}(B) \leq k$ 인 임의의 행렬 B 가 주어지면, $\|A - B\|_2 = \max \frac{\|(A-B)x\|_2}{\|x\|_2} \geq \sigma_{k+1}$ 이다.
- (증명) $\|A - A_k\|_2 = \sigma_{k+1}$ 임을 이용. 증명의 핵심은 $x \neq \mathbf{0}$ 이면서 $Bx = \mathbf{0}$ 와 $x = \sum_{i=1}^{k+1} c_i v_i$ 를 동시에 만족하는 벡터 x 를 찾는 것임.
- 이러한 조건을 만족하는 벡터 x 는 다음의 관계가 성립

$$\begin{aligned}\|(A - B)x\|_2^2 &= \|Ax\|_2^2 = \left\| \sum_{i=1}^{k+1} c_i \sigma_i u_i \right\|_2^2 = \sum_{i=1}^{k+1} c_i^2 \sigma_i^2 \\ \sum_{i=1}^{k+1} c_i^2 \sigma_i^2 &\geq \left(\sum_{i=1}^{k+1} c_i^2 \right) \sigma_{k+1}^2 = \|x\|_2^2 \sigma_{k+1}^2\end{aligned}$$

- 따라서 $\|(A - B)x\|_2^2 \geq \|x\|_2^2 \sigma_{k+1}^2$ 이므로, $\frac{\|(A-B)x\|_2}{\|x\|_2} \geq \sigma_{k+1}$ 이 성립
- 위의 조건을 만족하는 벡터 x 가 항상 존재한다고 말할 수 있는가?

I.9 PCA and the Best Low Rank Matrix

- Eckart-Young 정리 (l^2 -norm의 경우)
 - (증명 계속) $x \neq \mathbf{0}$ 이면서 $Bx = \mathbf{0}$ 와 $x = \sum_{i=1}^{k+1} c_i v_i$ 를 동시에 만족하는 벡터 x 가 항상 존재한다고 말할 수 있는가?
 - 행렬 B 의 rank는 최대 k 임. 따라서 행렬 B 의 nullspace의 차원은 적어도 $n - k$ 임.
 - $x = \sum_{i=1}^{k+1} c_i v_i$ 를 만들 때 사용된 v_i 가 생성하는 공간의 차원은 $k + 1$ 임. 이 벡터 x 가 행렬 B 의 nullspace에 포함되지 않는다면 벡터 x 내에 nullspace의 성분이 없어야 함.
 - 하지만, nullspace의 차원과 v_i 가 생성하는 공간의 차원의 산술적인 합은 적어도 $(n - k) + (k + 1) = n + 1$ 임.
 - 두 부분공간이 공통으로 가지고 있는 차원이 있어야 하고, 벡터 x 를 그 차원내에 선택하면 위의 두 조건을 동시에 만족가능함.
 - 따라서 $\|A - B\|_2 = \max \frac{\|(A-B)x\|_2}{\|x\|_2} \geq \sigma_{k+1}$ 은 성립해야 함

I.9 PCA and the Best Low Rank Matrix

- Frobenius Norm의 몇가지 다른 형태

- 1. $\|A\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2$

- 2. $\|A\|_F^2 = \text{tr}(A^T A) = (A^T A)_{11} + \dots + (A^T A)_{nn}$

- 3. $\|A\|_F^2 = |a_{11}|^2 + |a_{12}|^2 + \dots + |a_{mn}|^2$ (모든 $|a_{ij}|^2$ 성분들의 합)

- 2의 정의로부터 3의 유도

- $(A^T A)_{11}$ 는 A 의 첫번째 열의 성분들의 제곱의 합 $|a_{11}|^2 + |a_{21}|^2 + \dots + |a_{m1}|^2$ 임
 - 마찬가지로 $A^T A$ 의 대각성분 $(A^T A)_{ii}$ 는 A 의 i 번째 열의 성분들의 제곱의 합 $|a_{1i}|^2 + |a_{2i}|^2 + \dots + |a_{mi}|^2$ 이므로, 모든 대각 성분의 합은 행렬 A 의 모든 성분들의 제곱의 합이 됨.

- Frobenius norm에 대해서도 Eckart-Young 정리가 성립함을 보일 수 있음 → 교재 74페이지 참고.

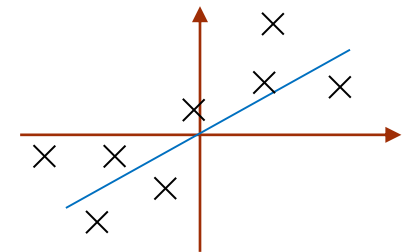


I.9 PCA and the Best Low Rank Matrix

- Principal component analysis (PCA)
 - n 명의 사람을 측정하는데 매번 측정마다 (나이, 키, 몸무게)같은 m 개의 측정값을 얻었다고 가정
 - 이 경우 행렬 A_0 는 각 사람의 데이터를 열벡터로 가진 $m \times n$ 행렬로 가정

평균 대비	사람1	사람2	사람3	사람4	사람5	사람6
나이 차이	3	-4	7	1	-4	-3
키 차이	7	-6	8	-1	-1	7

- 우선 행렬 A_0 의 각 행마다의 평균값을 구한 다음, 해당하는 행에서 빼 행렬을 A 라고 함
- 행렬 A 의 각 열은 \mathbf{R}^m 공간내 n 개의 점에 해당
- 행렬 A 의 각 행의 평균은 0
- 행렬 A 의 모든 열의 합과 평균도 0임
- n 개의 점들이 원점을 중심으로 분포되어 있는데, 특정한 선이나 평면같은 \mathbf{R}^m 공간의 낮은 차원을 가진 부분공간일 가능성이 있음
- 원점을 지나는 직선은 행렬 A 의 \mathbf{u}_1 벡터임



I.9 PCA and the Best Low Rank Matrix

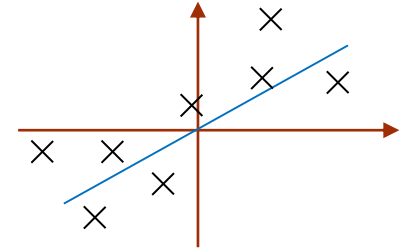
- Principal component analysis (PCA)의 통계적 관점
 - 행렬 A_0 의 각각의 행에서 평균을 제거하였으므로 행렬 A 의 모든 값은 평균으로부터 떨어진 거리임
 - 분산
 - 평균으로부터 떨어진 거리의 제곱의 합
 - 주어진 행렬 A 의 AA^T 의 대각 원소들의 합임
 - 공분산
 - 측정 데이터들간의 상관관계를 나타냄
 - 주어진 행렬 A 의 AA^T 의 대각이 아닌 원소들과 관련됨
 - ➔ (A 의 i 번째 행)과 (A 의 j 번째 행)간의 내적임
 - 샘플들의 공분산 행렬은 샘플의 개수가 n 일 때, 다음과 같이 나타남
 - $$S = \frac{AA^T}{n-1}$$

I.9 PCA and the Best Low Rank Matrix

- Principal component analysis (PCA)의 통계적 관점
 - 공분산 행렬 $S = \frac{AA^T}{n-1}$
 - $A = \begin{bmatrix} 3 & -4 & 7 & 1 & -4 & -3 \\ 7 & -6 & 8 & -1 & -1 & 7 \end{bmatrix}$
 - $S = \frac{AA^T}{6-1} = \begin{bmatrix} 20 & 25 \\ 25 & 40 \end{bmatrix}$
 - eigenvalue: 57, 3
 - eigenvalue 57에 해당하는 eigenvector \mathbf{u}_1 는 $\mathbf{u}_1 \approx (0.6, 0.8)$ 임



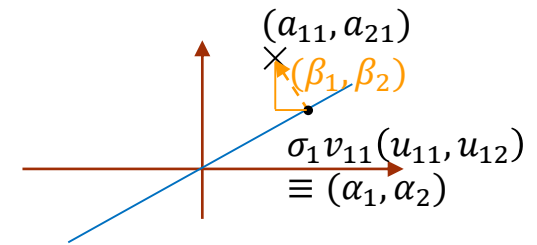
I.9 PCA and the Best Low Rank Matrix



- PCA의 기하학적 의미
 - ▣ 벡터 \mathbf{u}_1 이 나타내는 선과 점들간에는 perpendicular least square의 관계가 존재함.
 - ▣ 행렬 A 가 $A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$ 와 같이 principal component(PC)들의 합으로 나타났을 때, 주어진 예와 같이 행렬 A 의 크기가 $2 \times n$ 일 때 첫번째 PC는 아래와 같은 형태를 띠م
 - ▣ $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T = \sigma_1 \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} [v_{11} \ \dots \ v_{1n}] = \sigma_1 \begin{bmatrix} v_{11} \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} & \dots & v_{1n} \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} \end{bmatrix}$
 - ▣ 따라서 행렬 A 에서 첫번째 PC를 제거하는 것은 각 열에 있는 데이터에서 벡터 $\mathbf{u}_1^T = [u_{11} \ u_{12}]$ 와 평행한 성분을 제거하는 것에 해당한다. → Why?
 - $A - \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T = \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$ 는 다른 PC들로만 이루어져 있는데, 다른 PC들의 모든 열벡터들은 \mathbf{u}_1 와 직교하므로, 이들의 1차결합으로 이루어진 $A - \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$ 의 모든 열벡터들도 \mathbf{u}_1 와 직교함
 - → Gram-Schmidt직교화와 유사
 - ▣ 행렬 A 의 첫번째 열에 저장된 첫번째 샘플 데이터 $\begin{bmatrix} a_{11} \\ a_{21} \end{bmatrix}$ 에 끼치는 영향을 생각해 보면,

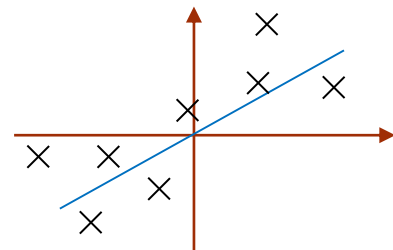
$$(a_{11}, a_{21}) = (\alpha_1, \alpha_2) + (\beta_1, \beta_2) = (\alpha_1 + \beta_1, \alpha_2 + \beta_2)$$
 - ▣ 원점에서 거리를 생각해 보면

$$a_{11}^2 + a_{21}^2 = (\alpha_1 + \beta_1)^2 + (\alpha_2 + \beta_2)^2 = \alpha_1^2 + \alpha_2^2 + \beta_1^2 + \beta_2^2$$
 - ▣ 위에서 (α_1, α_2) 과 (β_1, β_2) 는 직교하므로 $\alpha_1 \beta_1 + \alpha_2 \beta_2 = 0$ 임
 - ▣ 따라서 첫번째 PC를 제거하는 것은 각각의 샘플 데이터에서 원점까지의 거리의 제곱에서 \mathbf{u}_1 라인상의 샘플 데이터와 가장 점까지의 원점에서의 거리의 제곱을 뺀 것과 같고, 이는 데이터로부터 \mathbf{u}_1 라인까지의 거리의 제곱이라고 말할 수 있음.





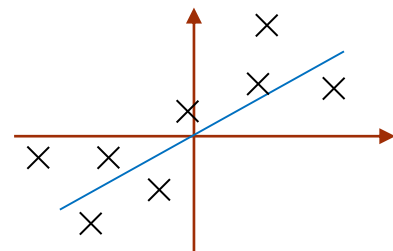
I.9 PCA and the Best Low Rank Matrix



- PCA의 기하학적 의미
 - ▣ 벡터 \mathbf{u}_1 이 나타내는 선과 점들간에는 perpendicular least square의 관계가 존재함.
 - ▣ 일반적인 최소 자승법(least square)
 - 다른 종류의 해석임
 - 는 $A\mathbf{x} = \mathbf{b}$ 를 푸는 해의 근사값을 얻기 위해 $\|A\mathbf{x} - \mathbf{b}\|^2$ 를 최소화하는 \mathbf{x} 를 구하는 방법임
 - $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$ 형태로 해를 구함
 - ▣ PCA는 특이값과 해당하는 특이벡터로 해를 구함
 - ▣ 비슷한 결과를 얻는 경우도 많으나, 최적화하는 목표가 다름
 - ▣ \mathbf{u}_1 이 나타내는 선으로부터 거리의 제곱의 합이 최소가 됨



I.9 PCA and the Best Low Rank Matrix



- PCA의 기하학적 의미
 - ▣ \mathbf{u}_1 이 나타내는 선으로부터 거리의 제곱의 합이 최소가 됨
 - 행렬 A 의 j 번째 열의 열벡터 \mathbf{a}_j 를 직교하는 두 기저 $\mathbf{u}_1, \mathbf{u}_2$ 와 평행한 성분으로 나눠서 생각해보면,

$$\sum_{j=1}^n \|\mathbf{a}_j\|^2 = \sum_{j=1}^n |\mathbf{a}_j^T \mathbf{u}_1|^2 + \sum_{j=1}^n |\mathbf{a}_j^T \mathbf{u}_2|^2$$
 - 등식의 오른쪽 첫번째 항은 $\mathbf{u}_1^T \mathbf{a}_j \mathbf{a}_j^T \mathbf{u}_1$ 의 합이므로, $\mathbf{u}_1^T (A A^T) \mathbf{u}_1$ 에 해당
 - 등식의 왼쪽 항이 샘플 데이터에 의해 고정된 값일 경우에는, 오른쪽 첫번째 항이 최대가 됨에 따라 나머지 거리가 최소로 됨
- 선형 대수적 관점
 - ▣ 모든 데이터의 분산의 합: $T = \|A\|_F^2 / (n - 1) = (\|\mathbf{a}_1\|^2 + \dots + \|\mathbf{a}_n\|^2) / (n - 1)$
 - ▣ $T = (\sigma_1^2 + \dots + \sigma_r^2) / (n - 1)$: 각각의 principal component가 전체 데이터의 분산에 끼치는 영향을 볼 수 있음