

Mathematical basis for PCA:

Singular value decomposition (SVD)

Eckart & Young (39.)

If B has rank M

$$\text{then } \|A - B\| \geq \|A - A_M\|$$

$$A_M = \sum_{j=1}^M \sigma_j \vec{u}_j \vec{v}_j^T$$

where \vec{u}_i are the right eigenvectors of A
 \vec{v}_i are the left eigenvectors of A
 σ_i are the singular values of A

These definitions
follow from
the SVD

SVD

Analogue to
Diagonalization
of
square
matrices

$$A V = V \Lambda$$

$$V = (\vec{v}_1 \dots \vec{v}_n)$$

$$\Lambda = \text{diag}(\lambda_i)$$

$$A = V \Lambda V^{-1}$$

$A_{n \times n}$

↳ but for rectangular matrices.

SVD

Analogue for rectangular matrices:

Given $A_{m \times n} \in \mathbb{R}^{m \times n}$

$$A_{m \times n} V_{n \times n} = U_{m \times m} \Sigma_{m \times n}$$

$$V_{n \times n} = (\vec{v}_1 \dots \vec{v}_n)$$

$$U_{m \times m} = (\vec{u}_1 \dots \vec{u}_m)$$

$$\Sigma_{m \times n} = \left(\begin{array}{ccc|c} \sigma_1 & & 0 & 0 \\ & \ddots & & \\ 0 & & \sigma_r & 0 \\ \hline & 0 & & 0 \end{array} \right)$$

Singular values.

$$\min(m, n) \geq r$$

r is the rank of A (at most $\min(m, n)$)
but could be smaller.

Reduced form of the SVD:

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T = \sum_{j=1}^r \sigma_j \vec{u}_j \vec{v}_j^T = (U_r)_{m \times r} \Sigma_r (V_r^T)_{r \times n}$$

Truncations of U and V taking the first r vectors \vec{u}_j and \vec{v}_j .

What are these U and V ?

Definition of
SVD

$$\begin{aligned} AV &= U\Sigma & VV^T &= I \\ & & UU^T &= I \\ A &= U\Sigma V^T \end{aligned}$$

$$(A^T A)_{n \times n} = V \Sigma^T \underbrace{U^T U}_I \Sigma V^T = V(\Sigma^T \Sigma) V^T$$

$$(A A^T)_{m \times m} = U \Sigma \underbrace{V^T V}_I \Sigma^T U^T = U(\Sigma \Sigma^T) U^T$$

$\left\{ \begin{array}{l} V_{n \times n} \text{ contains the eigenvectors of } (A^T A) \\ U_{m \times m} \text{ contains the eigenvectors of } (A A^T) \end{array} \right.$
 $\left\{ \Sigma^T \Sigma \text{ and } \Sigma \Sigma^T \right.$ containing the non-zero eigenvalues of $A^T A$ and $A A^T$

Singular values

$\sigma_k = \sqrt{\text{eigenvalues of } (A^T A) \text{ and } (A A^T)}$

Left and right singular vectors come in pairs with the same σ_k^2

$$\left\{ \begin{array}{l} (A^T A) \vec{v}_k = \sigma_k^2 \vec{v}_k \\ (A A^T) \vec{u}_k = \sigma_k^2 \vec{u}_k \end{array} \right.$$

with the same σ_k^2

It follows from the equality that:

$$\vec{u}_k = \frac{A \vec{v}_k}{\sigma_k}$$

check: $(AA^T) \vec{u}_k = A A^T \frac{A \vec{v}_k}{\sigma_k} = \sigma_k^2 \frac{A \vec{v}_k}{\sigma_k} = \sigma_k^2 \vec{u}_k \checkmark$

Eckart-Young:

of rank at most k

"Given $A_{m \times n}$ find $B_{m \times n}$ that is close to A "

i.e., $\|A - B\|$ small.

Matrix norm.

Matrix Norms:

- Induced norms: e.g. ℓ_2 norm \equiv Spectral norm

$$\|A\|_2 = \max_{\vec{x}} \frac{\|A\vec{x}\|}{\|\vec{x}\|} = \sigma_1$$

largest singular value of A

- Element-wise norms:

e.g. Frobenius norm:

$$\|A\|_F^2 = \sum_{i,j} |a_{ij}|^2 = \text{Tr}(A^T A) =$$

$$= \text{Tr}(U \Sigma^2 U^T) = \text{Tr}(\Sigma^2) = \sum_{j=1}^r \sigma_j^2$$

Equivalent to the vector norm of $\|\text{vec}(A)\|$

These are vector norms

(1) E-Y for spectral norm:

Statement: If $\text{rank}(B) \leq k$ then $\|A-B\| = \max_{\vec{x}} \frac{\|(A-B)\vec{x}\|}{\|\vec{x}\|}$

$$\geq \frac{\|(A-A_k)\vec{x}\|}{\|\vec{x}\|}$$

where $A_k = U_k \Sigma_k V_k$

Proof: (i) $\|(A-A_k)\vec{x}\| \geq \sigma_{k+1}$ since $A_k = U_k \Sigma_k V_k$ where these are from columns

$\exists \vec{x} \neq 0_v$ / $B\vec{x} = 0_v$ and $\vec{x} = \sum_{j=1}^{k+1} c_j \vec{v}_j$

given that $\begin{cases} \dim[\text{Ker } B] \geq n-k \\ \Rightarrow \{\vec{v}_j\}_{j=1}^{k+1} \cap \text{Ker } B \neq \emptyset \end{cases}$

Then we have: $\|(A-B)\vec{x}\|^2 = \|A\vec{x}\|^2 = \sum_{j=1}^{k+1} c_j^2 \sigma_j^2$

$$\geq \left(\sum_{j=1}^{k+1} c_j^2 \right) \sigma_{k+1}^2 = \|\vec{x}\|^2 \sigma_{k+1}^2$$

(ii) $\|(A-B)\vec{x}\|^2 \geq \|\vec{x}\|^2 \sigma_{k+1}^2$ ✓

where the $\{\sigma_1, \sigma_2, \dots, \sigma_k, \sigma_{k+1}, \dots, \sigma_r\}$ are ordered in decreasing order.

Hence: from (i) and (ii) $\frac{\|(A-B)\vec{x}\|^2}{\|\vec{x}\|^2} \geq \sigma_{k+1}^2 = \|A-A_k\|^2$ ✓

(2) E-y for Frobenius norm:

Find B that minimises
 $\min_B \|A - B\|_F^2$

where B of rank k
 or less

Any B of rank k or less
 can be written as:

$$B = CR$$

$m \times k \quad k \times n$

where

$$C^T C = D$$

$$R R^T = I$$

Consider E and minimise over C, R :

① $E = \|A - CR\|_F^2$

This is from SVD
 (for example)

Matrix
 calculus

$$\frac{\partial E}{\partial C} = 2(CR - A) R^T \quad ①$$

$$\frac{\partial E}{\partial R} = 2(R^T C^T - A^T) C \quad ②$$

~~$(A^T A) R^T$~~

At the
 minimum \Rightarrow

$$C R R^T = A R^T \quad ①$$

$$R^T C^T C = A^T C \quad ②$$

SVD
 gives best
 CR
 ✓

$$\left\{ \begin{array}{l} (A^T A) R^T = R^T D \\ \underline{R^T = V_k} \\ (A A^T) C = C D \\ \underline{C = U_k} \end{array} \right\}$$

Combining
 ① & ②
 we get

And the error is then:

~~$E = \|A - U_k V_k^T\|_F^2$~~ $E = \|A - CR\|_F^2 = \sum_{j=k+1}^r \sigma_j^2$

Constructive proof:

Connection with PCA

$$Y_{N \times P}^T = \begin{pmatrix} x_1^{(1)} & \dots & x_p^{(1)} \\ \vdots & & \vdots \\ x_1^{(N)} & \dots & x_p^{(N)} \end{pmatrix} = \begin{pmatrix} \bar{y}_1^T \\ \vdots \\ \bar{y}_N^T \end{pmatrix}$$

$$\bar{y}_i \in \mathbb{R}^P \quad i = 1, \dots, N$$

$$\vec{1}_{N \times 1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad \text{vector of ones.}$$

Let: $\vec{1}_{1 \times N}^T Y_{N \times P}^T = \begin{pmatrix} \sum_{i=1}^N x_1^{(i)} & \dots & \sum_{i=1}^N x_p^{(i)} \end{pmatrix}_{1 \times P}$

$$\frac{1}{N} \vec{1}^T (\vec{1}^T Y^T) = \begin{pmatrix} \langle x_1 \rangle_i & \dots & \langle x_p \rangle_i \\ \vdots & & \vdots \\ \langle x_1 \rangle_i & \dots & \langle x_p \rangle_i \end{pmatrix}_{N \times P}$$

Then define:

$$\tilde{Y}_{N \times P}^T = Y_{N \times P}^T - \frac{1}{N} \vec{1} \vec{1}^T Y^T = \underbrace{\left(I - \frac{1}{N} \vec{1} \vec{1}^T \right)}_{\text{Centering matrix}} Y^T$$

It then follows that

$$(\tilde{Y} \tilde{Y}^T)_{P \times P} = C_X \quad \text{sample covariance matrix}$$

Using the SVD of \tilde{Y} which is the centered data matrix.

So we conclude that PCA just considers the SVD of the centered data matrix \tilde{Y} , which is equivalent to obtaining the eigenvectors of the covariance matrix $C_X = (\tilde{Y} \tilde{Y}^T)$

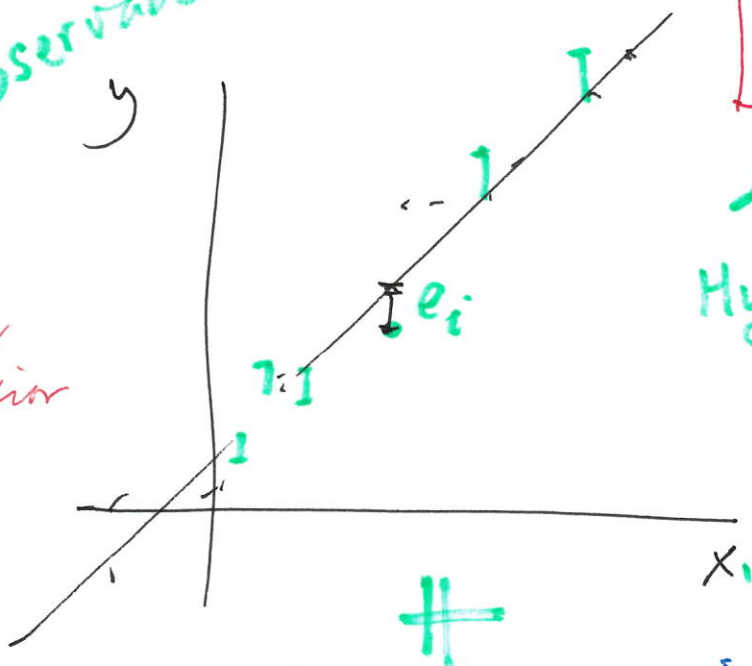
Observable:

Geometric picture of PCA/SVD

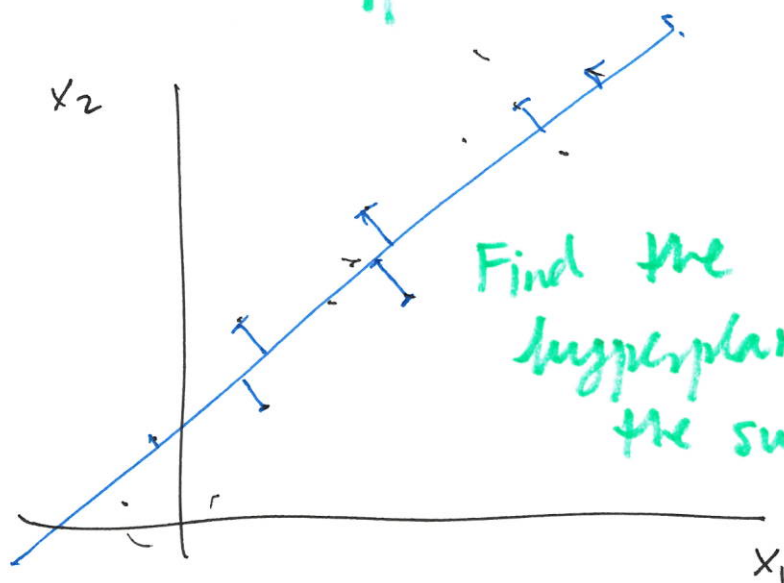
Linear regression:

Hyperplane that minimises the sum of squared error

(A)
Linear regression



(B)
PCA



PCA

Find the hyperplane such that the sum of ~~residual~~ distances that are normal to the plane are minimised.

The hyperplanes we obtain are not the same for linear regression and PCA.

PCA through SVD is close to total least squares or proper orthogonal decompositions.