

If $m > n$ and A has rank k , then

$$A = U_{m \times k} S_{k \times k} \underbrace{V^T}_{k \times n}$$

The pseudo-inverse of A is defined to be:

$$\underbrace{A^\dagger}_{A \text{ dagger}} = V S^{-1} U^T \quad (\text{also called Moore-Penrose inverse})$$

And even though A is not invertible,

$$A^\dagger A = (V S^{-1} U^T)(U S V^T)$$

$$= V S^{-1} U^T U S V^T$$

$$\underbrace{\quad \quad \quad}_{\substack{I \\ I}}$$

$$= V V^T = I \quad \text{when } V \text{ is square}$$

= projection onto row-space when
rank $A < n$.

Lastly the SVD affords the best 2-norm ^{rank k} matrix approximation:

Let A be $m \times n$ with rank r . Let $k < r$.

Then $\arg \min_{\substack{A', \\ \text{rank } A' = k}} \|A - A'\|_2 = \underbrace{U S_k V^T}_{\text{truncated SVD}} = (m \times n) \cdot \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \\ & & & \ddots \\ & & & & 0 \end{pmatrix} (n \times n)$

$n \times n$

$$\|A - U S_k V^T\|_2 = \|U S V^T - U S_k V^T\|_2$$

$$= \|U (S - S_k) V^T\|_2$$

$$= \|S - S_k\|_2 = \sigma_{k+1}$$

All methods are some form
of dimension reduction, etc.

Principal Component Analysis

Problem set up (in the continuous case)

Two main tenets:

- variation in data provides information
- correlation between data reduces amount of information.

(i.e. in the normal case, zero corr \Rightarrow indep.)

Idea: Transform (linearly) m variates into $k < m$ variates that contain almost as much variation as the original, but are approximately mutually independent.

Eg.
$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} a z_1 + b z_2 \\ c z_1 + d z_2 \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \quad z_1, z_2 \sim N(0,1)$$

If we knew $\text{Cov}(X_i, X_j)$ how could we decouple X_1, X_2 into z_1, z_2 ? What if z_1, z_2 were also correlated?

To re-phrase: If X_1, \dots, X_m can be linearly combined into Y_1, \dots, Y_m when Y_i, Y_j are independent and "contain" all of the variation, then Y_i 's are called the principal components.

(many other names exist too...)

We say X_1, \dots, X_m have been de-correlated.

Consider $\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$, Y_i a random variable (stocks, radiances, EEG signals, etc.)

then the covariance matrix is $C_{ij} = \text{cov}(Y_i, Y_j)$

$$\Rightarrow C = \mathbb{E} \left(\left(\vec{Y} - \mathbb{E}(\vec{Y}) \right) \left(\vec{Y}^T - \mathbb{E}(\vec{Y}^T) \right) \right)$$

WLOG, assume that $\mathbb{E}(Y_i) = 0$ and $\text{Var}(Y_i) = 1$.

Recall: C is a ^{SPSD} SPD matrix, and therefore can be diagonalized:

$$C = W D W^T \quad \text{where } W = (\vec{w}_1, \dots, \vec{w}_m) \quad \text{are}$$

the eigenvectors of C (with $\|\vec{w}_i\| = 1$) and

$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{pmatrix}$ are the eigenvalues.

The spectral representation of C is then:

$$C = \sum \lambda_j \vec{w}_j \vec{w}_j^T \quad (\text{assume } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0)$$

Define: $Y_{(k)} = \vec{w}_k^T \vec{Y}$ \leftarrow projection of \vec{Y} in the direction of maximum variance.

$$\begin{aligned} \Rightarrow \text{Var}(Y_{(k)}) &= \text{Var}(\vec{w}_k^T \vec{Y}) \\ &= \text{Var}\left(\sum_{i=1}^m w_{ki} Y_i\right) \\ &= \mathbb{E}\left(\sum_{i,j} w_{ki} Y_i \sum_j w_{kj} Y_j\right) \\ &= \vec{w}_k^T C \vec{w}_k = \vec{w}_k^T \lambda_k \vec{w}_k = \boxed{\lambda_k} \end{aligned}$$

The "total variation" of Y_1, \dots, Y_m can be defined as

$$\text{trace } C = \sum \lambda_i$$

And therefore $Y_{(k)}$ contains $\frac{\lambda_k}{\text{trace } C} = \frac{\lambda_k}{\sum \lambda_k}$ percent of the total variation.

To "approximate" the variats, pick p s.t. $\frac{\sum_1^p \lambda_k}{\sum_1^m \lambda_k}$ is

large enough, and then form

$$\vec{Y}_p = \begin{pmatrix} Y_{(1)} \\ \vdots \\ Y_{(p)} \end{pmatrix} = \begin{pmatrix} w_1^T \\ \vdots \\ w_p^T \end{pmatrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$$

← dimension reduction from m to p .

$\Rightarrow C_p = E(\vec{Y}_p \vec{Y}_p^T)$ is the "approximate" covariance matrix.

Thm C_p is the matrix closest to C in the Frobenius norm.

This was for random variats Y_1, \dots, Y_m — how about for actual data?

Eg. X_i models stocks, but you observe X_{ij} , $j=1, \dots, n$ prices/returns, how do you construct an "index" which captures "most" of the variance of the market?

PCA from data

Let X_{i1}, \dots, X_{in} , $i=1, \dots, m$ be n realizations of each of the m random variables.

- Assume $\bar{X}_i = \frac{1}{n} \sum_{j=1}^n X_{ij} = 0$

- The sample covariance matrix is S , with

$$S_{ij} = \frac{1}{n-1} \sum_{k=1}^n X_{ik} X_{jk}$$

- Assume that $S_{ii} = 1$.

- Let $X =$ "data matrix", $m \times n$, then $S = \frac{1}{n-1} X X^T$.

The matrix S is SPD and therefore also diagonalizable,

$$S = W D W^T \quad (W, D \text{ are not } \overset{\text{identically}}{\text{the same}} \text{ as before})$$
$$= \sum_{i=1}^m \lambda_i \vec{w}_i \vec{w}_i^T$$

And likewise the principal components of the data X

are

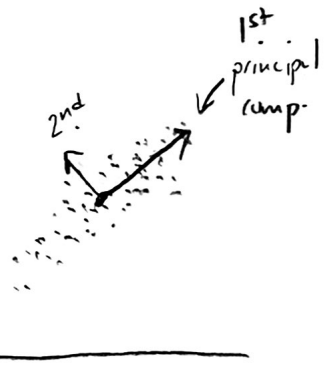
$$\tilde{X}_{(i)} = \vec{w}_i^T X$$

↑ ↑ ↑
1 × n 1 × m m × n

↑
a row vector
of "observations"
of the principal component.

w_{ij} = the contribution of

Graphically:



Practical Implementation of PCA

Recall, if the sample covariance matrix is

$$S = (X - \bar{X})(X - \bar{X})^T \frac{1}{n-1}$$
$$= W D W^T$$

However, we can work directly with the matrix $\frac{X - \bar{X}}{\sqrt{n-1}}$ by taking its SVD:

$$\frac{X - \bar{X}}{\sqrt{n-1}} = U S V^T$$

and then obviously:

$$\left(\frac{X - \bar{X}}{\sqrt{n-1}}\right) \left(\frac{X - \bar{X}}{\sqrt{n-1}}\right)^T = U S V^T V S^T U^T$$
$$= U S^2 U^T$$

↑ ↑ ↑
W D W^T

This is often much cheaper than forming S ($m n^2$) and then computing its eigenvectors (m^3), particularly if only one or two principal components are needed.

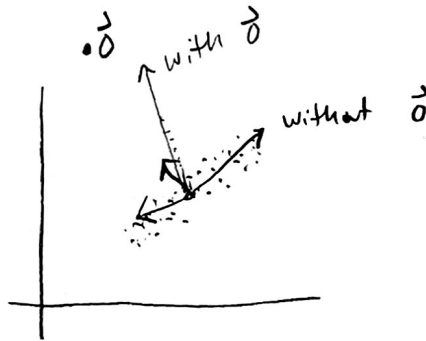
↑
(or k)

(Also, forming S squares the cond. number of $X - \bar{X}$, making it more accurate to work with the SVD instead.)

Robustness of PCA

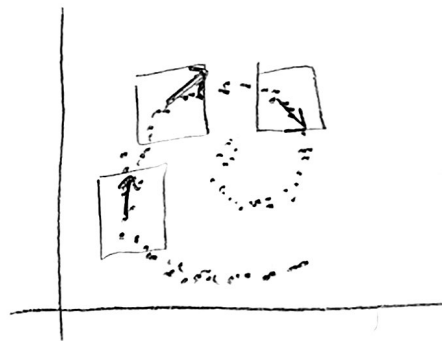
What types of things affect PCA?

Ex: outliers



Solution? Somehow ignore outliers:

- ① remove them, and do global PCA
- ② "locally cluster", and then do "local PCA"



What do PCs look like?

Clustering first avoids global effects.

Note: Much of this is "art"...

Related methods - Factor analysis (different generative model)

- Independent component analysis (different goal, independence)

(\leftrightarrow) principal