

Principal Component Analysis (PCA)

Method for dimensionality reduction

We focus now on input variables $x_1, \dots, x_n \in \mathbb{R}^D$.

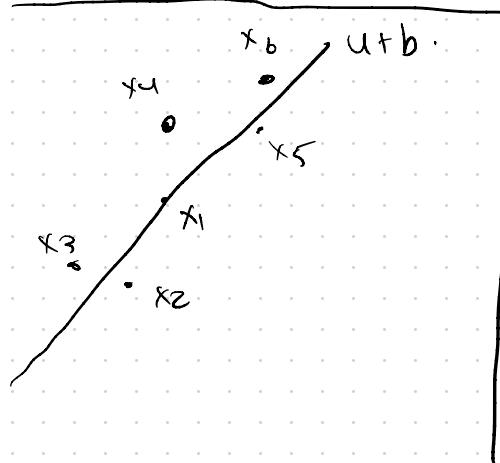
Why PCA?

- data compression
- Learn shape of data (when x_1, \dots, x_n are geometric)

Goal of PCA: find linear space $U \subseteq \mathbb{R}^D$

with $\dim(U) = d \leq D$ and a vector $b \in \mathbb{R}^D$ so that

$U + b$ is "close to" x_1, \dots, x_n .



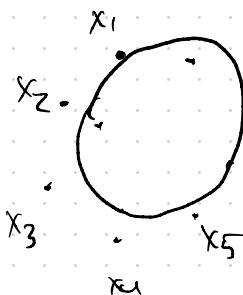
What about lower dimensional data that is not linear?

Compose with a feature map
 $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M$

$$z_i = \phi(x_i) \quad k \leq i \leq n$$

ex $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \left(1, x_1, x_2, x_1x_2, x_1^2, x_2^2 \right)$$



perhaps quadratic, so
map to all monomials degree ≤ 2 .

Image under ϕ very
close to $z_5 + z_6 = 1$

The x_i live close to
a zero set of a polynomial
Called an "algebraic variety"

The z_i live close to
linear space.

Suppose $x_1, \dots, x_n \stackrel{iid}{\sim} x \in \mathbb{R}^D$
 ↑ random variable $\phi = \text{feature map}$
 $\Rightarrow z_1, \dots, z_n \stackrel{iid}{\sim} z \in \mathbb{R}^M$ $\mathbb{R}^D \rightarrow \mathbb{R}^M$

Let $z = (z^{(1)} \dots z^{(n)})$, $\mu = \mathbb{E} z \in \mathbb{R}^M$
 Expected value

Covariance matrix

$$\Sigma := \left[\text{Cov}(z^{(i)}, z^{(j)}) \right]_{i,j=1}^M \in \mathbb{R}^{M \times M}$$

$$\text{Cov}(X, Y) = \mathbb{E} XY - \mathbb{E} X \mathbb{E} Y$$

Fact

Σ is positive semidefinite

\Leftrightarrow diagonalizable with orthogonal matrix P .

i.e. $P \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_M \end{pmatrix} P^T = \Sigma$

with eigenvalues $\lambda_1 \geq \dots \geq \lambda_M \geq 0$

We don't know distribution of Z , so we want to approximate Σ .

Definition

$$z_1, \dots, z_n \in \mathbb{R}^m$$

1. empirical average is $\bar{z} = \frac{1}{n} (z_1 + \dots + z_n)$

2. empirical covariance matrix is $S = (s_{ij})_{i,j=1}^M$

with

$$s_{ij} = \frac{1}{n} \sum_{k=1}^n [(z_k)_i - \bar{z}_i] \cdot [(z_k)_j - \bar{z}_j]$$

equivalently for Σ the feature matrix

$$\Sigma = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} \in \mathbb{R}^{n \times M}$$

Then

$$\bar{z} = \frac{1}{n} \Sigma^T e$$

$$e = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n$$

$$S = \frac{1}{n} (\Sigma \cdot e \bar{z}^T)^T (\Sigma - \bar{z} e)$$

Notion of distance for $\mathbb{I} \leq d < M$,

A space of maximal variance is a linear space $U \subseteq \mathbb{R}^M$ of $\dim(U) = d$ so

$$U \in \operatorname{argmax}_{U: \dim(U)=d} \mathbb{E} \|P_U(z-U)\|^2$$

Where P_U is orthogonal projection onto U .

Thm 3.39

$\lambda_1 \geq \dots \geq \lambda_d \geq 0$ eigenvalues of Σ ,
 u_i eigenvectors associated to λ_i so that

$$\langle u_i, u_j \rangle = \delta_{ij} \text{ - Then}$$

$$U = \operatorname{span}\{u_1, \dots, u_d\}$$

is a space of maximal variance of dimension d .

Moreover when $\lambda_d > \lambda_{d+1}$, U is uniquely determined.

Problem Σ hard to know, so we need to use $\hat{\Sigma}$.

Note on choosing d

Can choose $\lambda_{d+1} \approx 0$

or maximize $\frac{\lambda_d}{\lambda_{d+1}}$ or $\lambda_d - \lambda_{d+1}$

art of
Data
science.

Definition

u_1, \dots, u_d are called principal components of S .

From Thm 3.39 we compress data by mapping

$$z_i \mapsto P_u(z_i - \bar{z}) + \bar{z}$$

To make computationally viable, we minimize ^{distance to} data points:

$$\text{minimize } \sum_{i=1}^n \|(z_i - \bar{z}) - P_u(z_i - \bar{z})\|$$



Theorem

Let $\lambda_1 \geq \dots \geq \lambda_n \geq 0$ be eigenvectors of S .

Let u_i orthogonal basis of eigenvectors λ_i

so $\langle u_i, u_j \rangle = \delta_{ij}$. Then $U = \text{Span}\{u_1, \dots, u_d\}$

minimizes \star If $\lambda_d > \lambda_{d+1}$, U is uniquely determined.

Pf

Let u_i be any orthogonal basis of \mathbb{R}^M .

Center data set $w_i^* = z_i - \bar{z}$

let $A = [u_1 \dots u_d] \in \mathbb{R}^{M \times d}$, thus $\begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$

$$P_u = A A^T = \left[\begin{array}{c|c} u_1 & \dots & u_d \\ \hline \end{array} \right] \sum_{i=1}^d u_i u_i^T$$

$$\text{Set } W = (\Sigma - e\bar{z}^T)^T \in \mathbb{R}^{n \times n}$$

$$\text{So } nS = WW^T$$

$$\text{Since } I_M = \sum_{i=1}^M u_i u_i^T,$$

$$W - P_{\mathcal{U}} W = (I_M - A A^T) W = \left(\sum_{i=d+1}^M u_i u_i^T \right) W$$

$$\text{So } \star \text{ minimizes } \sum_{i=1}^n \|w_i - P_{\mathcal{U}}(w_i)\|^2$$

$$= \text{Trace} (W - P_{\mathcal{U}} W)^T (W - P_{\mathcal{U}} W)$$

$$= \text{Trace} \left(W^T \left(\sum_{i=d+1}^M u_i u_i^T \right)^T \left(\sum_{i=d+1}^M u_i u_i^T \right) W \right)$$

$$= \text{Trace} \left(\sum_{i=d+1}^M \sum_{j=d+1}^M u_i^T u_j u_i u_j^T \right) W W^T$$

(Since trace is cyclic invariant)

$$= \text{Trace} \left(\sum_{i=d+1}^M u_i u_i^T W W^T \right)$$

(since $\langle u_i, u_j \rangle = 1 \Leftrightarrow i=j$)

$$(\text{linearity of trace}) = \sum_{i=d+1}^M \text{Trace} (u_i u_i^T W W^T)$$

$$\begin{aligned} & \text{trace}(u^T u) \\ \text{where } u = & \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \\ u^T u \in \mathbb{R}^{n \times n} & \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \\ \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} u^T u = & \begin{bmatrix} u_1^T u_1 & \cdots & u_1^T u_n \\ \vdots & \ddots & \vdots \\ u_n^T u_1 & \cdots & u_n^T u_n \end{bmatrix} \end{aligned}$$

$$u_i^T u_i = \|u_i\|^2$$

$$\begin{aligned}
 & (nS = u^T) \\
 &= n \sum_{i=d+1}^M \text{Trace}(u_i^T S u_i) \\
 &\quad S u_i = \lambda_i u_i \quad \text{If } u_i \text{ eigenvector of } S \\
 &= n \sum_{i=d+1}^M \lambda_i \text{Trace}(u_i u_i^T) \\
 &= \boxed{n \sum_{i=d+1}^M \lambda_i} \quad \text{This is the squared distance! Minimized by taking smallest eigenvalues!}
 \end{aligned}$$

Proof of minimization using Lagrange multipliers

$$\text{minimize } n \sum_{i=d+1}^M u_i^T S u_i$$

subject to constraints $\langle u_i, u_j \rangle = \delta_{ij}$

Add parameter l_{ij} for $j = d+1, \dots, M$ such that $d+1 \leq i \leq j$

$$L(u_{d+1}, \dots, u_M, l_{ij}) = \sum_{i=d+1}^M u_i^T S u_i - \sum_{d+1 \leq i \leq j \leq M} (u_i^T u_j - \delta_{ij}) l_{ij}$$

Optimized when $\frac{\partial L}{\partial l_{ij}} = 0$ and $\frac{\partial L}{\partial u_i} = 0$

$$\frac{\partial \mathcal{L}}{\partial u_j} = 2S_{uj} - 2l_{ij} u_j - \sum_{d+1 \leq i < j} u_i l_{ij} = 0$$

$$d+1 \leq j \leq M$$

$$\frac{\partial \mathcal{L}}{\partial l_{ij}} = u_i^T u_j - \delta_{ij} = 0$$

When $j = d+1$

$$2S_{ud+1} = 2u_{d+1} l_{d+1,d+1}$$

$$\Rightarrow l_{d+1,d+1} = \lambda_{d+1}$$

when $j = d+2$

$$2S_{ud+2} - 2l_{d+2,d+2} u_{d+2} - \left(u_{d+1} l_{d+1,d+2} \right) = 0$$

Multiply by u_{d+1}^T on left

$$2u_{d+1}^T S_{ud+2} - \boxed{2l_{d+2,d+2} u_{d+1}^T u_{d+2}} - l_{d+1,d+2} = 0$$

$$2u_{d+1}^T S_{ud+2} = l_{d+1,d+2}$$

$$\frac{l_{d+1,d+2}}{2} = (U_{d+1}^T S_{d+2})^T = U_{d+2}^T S^T U_{d+1}$$

$$\begin{aligned} \left(\text{Since } S \text{ symmetric} \right) &= U_{d+2}^T S_{d+1} \\ &= \lambda_{d+1, d+1} U_{d+2}^T U_{d+1} = 0 \end{aligned}$$

$\Rightarrow l_{d+1, d+2} = 0$ and thus

$$S_{d+2} = l_{d+2, d+2} U_{d+2} \Rightarrow l_{d+2, d+2} = \lambda_{d+2}$$

\rightsquigarrow And so on to conclude $l_{ij} = \begin{cases} \lambda_j & i=j \\ 0 & i \neq j \end{cases}$

So the eigenvalues give optimal solutions
(i.e. derivative zero), the smallest eigenvalue
is the minimum.

\rightsquigarrow Note $\lambda_d > \lambda_{d+1}$ makes U unique
since there are no other choices for U .

Using SVD to get eigenvalues +
Notebook 6.

Next time