

PCA

①

We are given a set of vectors (x_1, \dots, x_m) , $x_i \in \mathbb{R}^d$

Goal: "reduce the dimensionality" of these vectors: i.e., for some $n < d$, find a matrix $W \in \mathbb{R}^{n \times d}$ s.t. Wx is compressed representation. Let $U \in \mathbb{R}^{d \times n}$ be the reconstruction matrix s.t. UWx is reconstruction of x .

Formally, want to solve

$$\min_{\substack{U \in \mathbb{R}^{d \times n} \\ W \in \mathbb{R}^{n \times d}}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2 \quad (*)$$

Note: if a_i are columns of A , then $\sum \|a_i\|_2^2 = \|A\|_F^2$, so we are minimizing the Frobenius norm

$$\|X - UW X\|_2^2, \text{ where } X = [x_1, \dots, x_m] \in \mathbb{R}^{d \times m}$$

The following claim simplifies the problem

Claim: Let (U, W) be a solution to $(*)$ then the columns of U are orthonormal, & $U^T U = I_n$, & $W = U^T$.

pf

1. Fix any U, W .
2. Let $R := \{UWx : x \in \mathbb{R}^d\}$ be range of map. Note $R \subset \mathbb{R}^d$ but $\dim R = n$. Let $V \in \mathbb{R}^{d \times n}$ be matrix whose cols form orthonormal basis for R .
3. Let x be given & consider $\min_y \|x - Vy\|_2^2$. Using calculus, this is

$$V^T V x = \arg \min_{\tilde{x} \in R} \|x - \tilde{x}\|$$

To see this, note

$\|x - Vy\|_2^2 = \|x\|_2^2 + y^T \underbrace{V^T V}_{I_n} y - 2y^T V^T x = \|x\|_2^2 + \|y\|_2^2 - 2y^T V^T x$. This is strongly convex. Set gradient in y to 0 to get $y^* = V^T x$.

Replacing u, w w/ v, v^T doesn't increase the objective.

$$\sum_i \|x_i - uw^T x_i\|^2 \geq \sum_i \|x_i - vv^T x_i\|^2.$$

We had w, u fixed. But since this holds $\forall w, u$, we're done. □

So, an equivalent optimization problem is

$$\begin{aligned} \min_{u \in \mathbb{R}^{d \times n}} & \sum_i \|x_i - uu^T x_i\|^2 \quad (***) \\ & u^T u = I_n \end{aligned}$$

Trick observe

$$\|x - u^T u x\|^2 = \|x\|^2 - 2x^T u u^T x - \underbrace{x^T u u^T u u^T x}_{I_n}$$

$$= \|x\|^2 - x^T u u^T x$$

$$= \|x\|^2 - \text{tr}(u^T x x^T u)$$

\uparrow

Note: If $z \in \mathbb{R}^n$, then $\underbrace{z^T z}_1 = \text{tr}(z z^T)$

Let $A = \sum_{i=1}^m x_i x_i^T$. A is symmetric w/ decomposition $\sum z_i^2 = \sum z_i^2$

$$A = V D V^T$$

where $D = \text{diag}$ & $V^T V = V V^T = I_d$

Since A is PSD, diag of D is ≥ 0 .

So problem is equivalent to

$$\begin{aligned} \arg \max_{u \in \mathbb{R}^{d \times n}} & \text{tr}(u^T A u) \quad (***) \\ & u^T u = I_n \end{aligned}$$

Thm Let $(x_1, \dots, x_m) \subset \mathbb{R}^d$. Let $A = \sum_{i=1}^m x_i x_i^T$ & let u_1, \dots, u_n be n eigenvectors corresponding to n largest eigenvalues. Then solution of ~~(*)~~ is to set $U = (u_1, \dots, u_n)$ & $W = U^T$.

pf.

1. Suffices to look at ~~(*)~~

2. Will show that $\text{tr}(U^T A U) \leq \sum_{i=1}^n D_{j,j}$ \swarrow n largest eigenvalues of A .

But this value is attained by our solution. So if we show this, we're done.

3. Fix some ^{arbitrary} matrix $U \in \mathbb{R}^{d \times n}$ w/ orthonormal cols. Let $B = U^T U$ so $VB = U U^T U = U$. Have

$$U^T A U = \underbrace{B^T}_{W^T = U^T} \underbrace{V^T V D V^T}_{= A} V B = B^T D B$$

Have

$$DB = \begin{pmatrix} D_1 (\text{row 1 of } B) \\ \vdots \\ D_{\perp} (\text{row } \perp \text{ of } B) \end{pmatrix} \quad B \in \mathbb{R}^{d \times n}$$

$$\text{tr}(B D B) = \text{tr} \left(\underbrace{\begin{bmatrix} B_1^T & \dots & B_{\perp}^T \end{bmatrix}}_1 \begin{bmatrix} D_1 B_1 \\ \vdots \\ D_{\perp} B_{\perp} \end{bmatrix} \right) \quad B_i \in \mathbb{R}^{1 \times n}$$

Multiply out 1st block to

$$\text{get } D_1 \sum_{j=1}^n B_{1,j}^2$$

$$\text{tr}(\cdot) = D_1 \sum_{j=1}^n b_{1,j}^2 + \dots + D_{\perp} \sum_{j=1}^n b_{\perp,j}^2 = \sum_{i=1}^{\perp} D_{i,i} \sum_{j=1}^n b_{i,j}^2$$

4. This is annoying, but here $\sum_{i=1}^d \sum_{j=1}^n b_{i,j}^2 = n$. by orthogonality,
 to sum $(\text{col } i) \cdot (\text{col } j)$

$$\sum_{j=1}^n \underbrace{\sum_{i=1}^d b_{i,j}^2}_{(\text{col } j) \cdot (\text{col } j)} = n$$

If we augment B to \tilde{B} w/ $\tilde{B} = [B, \text{other}]$

$$\sum_{j=1}^d \tilde{B}_{i,j}^2 = 1 \quad \forall i \quad \Rightarrow \quad \sum_{j=1}^n B_{i,j}^2 \leq 1$$

w/ $\tilde{B}^T \tilde{B} = I$, then

Hence,

$$\text{tr}(U^T A U) \leq \max_{\substack{\beta \in [0,1]^d \\ \|\beta\| \leq n}} \sum_{i=1}^d D_{i,i} \beta_i$$

The max on RHS is attained by vector β place all mass on $n \log d$. Hence

$$\text{tr}(U^T A U) \leq \sum_{i=1}^n D_{i,i} \leq n \log d$$

least $n \log d$ of d

as desired

□

Perspective 2: An SVD based perspective

(5)

Want to solve

$$\min_{\Sigma} \left\| \Sigma - \underbrace{UW}_{\text{rank } n} \right\|_F^2$$

equivalently, solve

$$\min_{\substack{\text{rank } B \leq n \\ \dim B = \dim \Sigma}} \left\| \Sigma - B \right\|_F^2$$

Thm (Eckart-Young-Mirsky) Let $A = U\Sigma V^T$ be the SVD of A . The best rank n approx to A in Frobenius Norm is

$$A_n := \sum_{i=1}^n \sigma_i u_i v_i^T$$

pf 1. As mtr. of norms, s.t. $\sigma_1(B) \geq \sigma_2(B) \geq \dots$ for a matrix.

Groundwork claim:

$$\sigma_{i+j-1}(\Sigma + \Delta) \leq \sigma_i(\Sigma) + \sigma_j(\Delta)$$

To do: prove this. Weyl's thm for eigenvalues/singular values?

Would like to show that if $\text{rank } B_n = n$, $\|A - A_n\|_F \leq \|A - B_n\|_F$

observe

$$\|A - A_n\|_F^2 = \left\| \sum_{i=n+1}^r \sigma_i u_i v_i^T \right\|_F^2 = \sum_{i=n+1}^r \sigma_i^2(A)$$

rank of A

Let $u = u_i, v = v_i$,

$$\|uv^T\|_F = \sum_k \sum_l (u_k v_l)^2 = \underbrace{\sum_k u_k^2}_{=1} \underbrace{\sum_l v_l^2}_{=1}$$

Note that $\sigma_{n+1}(B_n) = 0$. In above claim, let $\gamma = B_n$ & $A = A - B_n$, & $j = n+1$

Get,

$$\sigma_{i+n}(A) \leq \sigma_i(A - B_n). \quad \text{Hence,}$$

$$\|A - B_n\|_F^2 \geq \sum_{i=1}^r \sigma_i^2(A - B) \geq \sum_{i=n+1}^r \sigma_i^2(A) = \|A - A_n\|_F^2$$

Exact same
reasoning as before

What's the connection between the SVD low rank approx perspective & the "compression" perspective?

First, note that PCA perspective gives us a low dimensional representation of \mathbf{X} , namely $\mathbf{U}^T \mathbf{X}$. The SVD perspective just gives us a low rank matrix - not helpful for simplicity, dimensionality, issues.

Let's clarify what we have

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

$d \times m \quad d \times d \quad d \times m \quad m \times m$

$$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_d]$$

$$\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_d]$$

To disambiguate, call matrix \mathbf{U} from PCA stuff \mathbf{W} , so

$\mathbf{W} \mathbf{W}^T \mathbf{X}$ is low-rank approx to \mathbf{X} , $\mathbf{W} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$.

$\mathbf{W}: d \times n$, $\mathbf{W}^T: n \times d$. ↑
also cols of \mathbf{U} in SVD

\mathbf{W} compresses compresses cols of \mathbf{X} as $\mathbf{W}^T \mathbf{x}_i$, $\mathbf{W} \mathbf{W}^T \mathbf{x}_i$ brings them back.

So, it appears we have that $\mathbf{W} \mathbf{W}^T \mathbf{X} \approx \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ are both optimal low rank approximations, what gives? Are these the same? ... Yes!

To make the main obvious, call $\mathbf{W} = [\mathbf{u}_1, \dots, \mathbf{u}_n] =: \mathbf{U}_n$.

Not

$$\sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \begin{matrix} \mathbf{1} \times n & \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} & \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} \\ & n \times n & n \times m \end{matrix}$$

But, from PCA,

$$\begin{aligned} \mathbf{W} \mathbf{W}^T \mathbf{X} &= \mathbf{U}_n \mathbf{U}_n^T \mathbf{X} \stackrel{\text{SVD}}{=} \mathbf{U}_n \underbrace{\mathbf{U}_n^T \mathbf{U}}_{\mathbf{I}_n} \mathbf{\Sigma} \mathbf{V}^T \\ &= \mathbf{U}_n [\mathbf{I}_n \ 0] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n & & \end{bmatrix} \mathbf{V}^T = \begin{bmatrix} \mathbf{I}_n & \mathbf{0}_{n \times d} \end{bmatrix} \\ &= \mathbf{U}_n \underbrace{\mathbf{\Sigma}_n}_{\substack{\text{w/ obvious modifications} \\ \text{w/ obvious modifications}}} \mathbf{V}_n^T \end{aligned}$$

$\mathbf{U}_n^T \mathbf{U} = \begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_n^T \end{bmatrix} [\mathbf{u}_1 \dots \mathbf{u}_d \dots \mathbf{u}_m]$

A couple of quick follow up thoughts:

- There is a nice 2d ^{2d} demo of PCA in my notes folder

- Also, a demo using eigenfaces (pulled the demo from the internet)

Thinking quickly now, I think eigenfaces are just $\sigma_i u_i v_i^T$ for $i = 1, \dots, k$
 k Singular Values/Vectors.