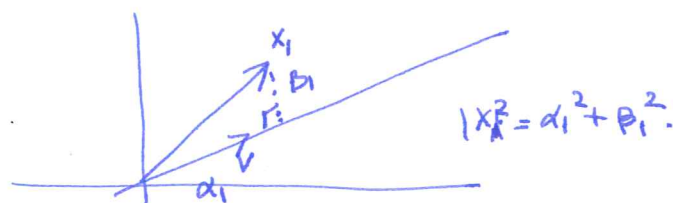


## Singular Value Decomposition

We're given a set of points  $X = \{x_1, x_2, \dots, x_n\}$  in  $\mathbb{R}^D$ . They're placed as rows of a matrix  $A = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ . We want to represent them using a small # of dimensions, but now not to preserve distances but preserve the "energy".

The 1-dimensional problem: find a 1-dimensional subspace such that the sum of projections (squared) is maximized - or the sum of distances squared is minimized.

These are equivalent (by Pythagoras).



So want  $v$  st.  $\sum_i (x_i \cdot v)^2$  is maximized st  $\|v\|=1$ . but  $\sum (x_i \cdot v)^2 = \|Av\|^2$ .

first (right) singular vector  $v_1 := \arg \max_{\|v\|=1} \|Av\|^2 = \arg \max_{\|v\|=1} \|Av\|$ .

first singular value  $\sigma_1(A) := \max_{\|v\|=1} \|Av\| = \|Av_1\|$ .

$\sigma_1^2 =$  sum of squared projections of points on  $v_1$ .

Suppose want 2-d subspace in some greedy fashion: - let's find

$v_2 = \arg \max_{\|v\|=1, v \perp v_1} \|Av\|$ . and  $\sigma_2(A) = \max_{\|v\|=1, v \perp v_1} \|Av\|$ .

Similarly:

$v_i = \arg \max_{\|v\|=1, v \perp v_1, v_2, \dots, v_{i-1}} \|Av\|$  and  $\sigma_i(A) = \max_{\|v\|=1, v \perp v_1, \dots, v_{i-1}} \|Av\|$ .

Observe: we've just defined things greedily. Find  $v_1$ , then  $v_2$ , then  $v_3, \dots, v_k$ .

What if the best 2-dim subspace did some global optimization nongreedily?

Thankfully, things are pretty good for us!

(2)

Claim: Suppose  $A \in \mathbb{R}^{n \times d}$  with singular vectors  $v_1, v_2, \dots, v_r$ . means  $\max_{v \perp v_1, \dots, v_r} \|Av\| = 0$ .

Let  $V_k = \text{span}(v_1, \dots, v_k)$ .  $\forall k \leq r$ .

then  $V_k$  is best-fit  $k$ -dimensional subspace for  $A$ .

(i.e. sps  $S$  is another  $k$ -dim subspace spanned by  $w_1, \dots, w_k$  (orthonormal basis))  
 $\Rightarrow \sum_{i=1}^k \|Aw_i\|^2 \leq \sum_{i=1}^k \|Av_i\|^2$ .

Pf: Base case  $k=1$  (trivially true)

$k=2$ :  $v_1, v_2$  are sig vectors. Say  $S$  is another subspace. Let  $w_1, w_2$  span  $S$ .

Choose  $w_2 \perp v_1$  in the subspace, and then  $w_1$  accordingly.

$\cdot \|Aw_2\|^2 \leq \|Av_2\|^2$  by choice of  $v_2$ . ✓

$\cdot \|Aw_1\|^2 \leq \|Av_1\|^2$  by choice of  $v_1$ .

$k$  general. By induction  $V_{k-1}$  is best fit  $(k-1)$ -dim subspace.

Choose  $w_1, \dots, w_k$  to span  $S$  such that  $w_k \perp v_1, \dots, v_{k-1}$ .

Again same argument. ▣

Fact:  $\sum \sigma_i^2 = \|A\|_F^2 = \sum_{ij} a_{ij}^2$ .

Pf:  $\sum_i \sigma_i^2 = \sum_i \|Av_i\|^2 = \sum_i \sum_j \langle a_j, v_i \rangle^2 = \sum_j \sum_i \langle a_j, v_i \rangle^2$

but  $v_i$  form an orthonormal basis for the row space of  $A$ .

$= \sum_j \|a_j\|^2 = \sum_{jk} a_{jk}^2 = \|A\|_F^2$ . ▣

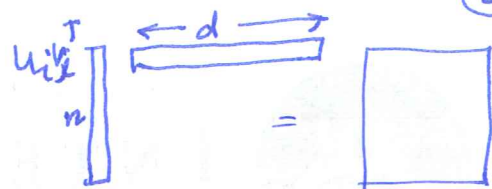
Now define:  $u_i = \frac{Av_i}{\|Av_i\|} = \frac{Av_i}{\sigma_i}$ . Clearly unit vectors. (also in fact orthonormal)

$u_i := i$ th left singular vector of  $A$ .

$= \arg \max_{\substack{u \perp u_1, \dots, u_{i-1} \\ \|u\|=1}} \|uA\|$  much like the right singular vectors.

Claim:  $A = \sum \sigma_i u_i v_i^T$

(BTW: if  $U = \begin{pmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_r \\ | & | & & | \end{pmatrix}$



$D = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \dots & \\ & & & \sigma_r \end{pmatrix}$   $V = \begin{pmatrix} | & | & & | \\ v_1 & v_2 & \dots & v_r \\ | & | & & | \end{pmatrix}$  then  $\sum \sigma_i u_i v_i^T = U D V^T$   
 $= \begin{matrix} \begin{matrix} | & | & & | \\ u_1 & u_2 & \dots & u_r \\ | & | & & | \end{matrix} \begin{matrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \dots & \\ & & & \sigma_r \end{matrix} \begin{matrix} \begin{matrix} | & | & & | \\ v_1 & v_2 & \dots & v_r \\ | & | & & | \end{matrix} \end{matrix}$

Pf: Subclaim: A and B are same iff  $Av = Bv \forall v \in \mathbb{R}^n$ .

Pf:  $\Rightarrow$  trivial  $\leftarrow$  ops not then choose  $v = e_i$ . where differ in  $i$ th column.

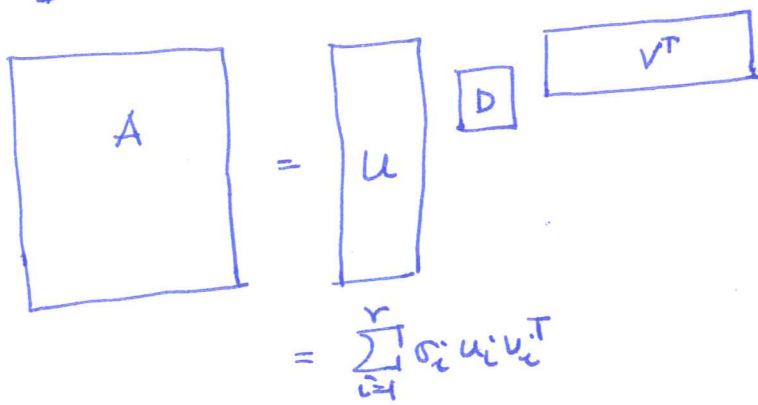
Subclaim: A and B are same iff  $Av = Bv \forall v$  in basis. Pf write each  $v$  as this basis.

Now: consider the basis of  $v_1, v_2, \dots, v_r, \underbrace{v_{r+1}, \dots, v_d}_{\text{extend to orthonormal basis}}$

for any  $j \in [r]$   ~~$Av_j = (\sum \sigma_i u_i v_i^T) \cdot v_j = \sigma_j u_j$~~   $\downarrow$   $Av_j$  (by defn of  $u_j$ )

$\Rightarrow$  A and  $\sum \sigma_i u_i v_i^T$  are identical.

for  $j \notin [r]$   $Av_j = 0 = (\sum \sigma_i u_i v_i^T) v_j$



the columns of  $U$  &  $V$  are orthonormal  
 (haven't proved it for columns of  $U$  yet  $\rightarrow$  do it in HW)

Claim: Suppose  $A = \tilde{U} \tilde{D} \tilde{V}^T$  for matrices  $\tilde{U}, \tilde{D}, \tilde{V}$  where  $\tilde{D}$  is diagonal,  $\tilde{U}, \tilde{V}$  have orthonormal columns. then  $\tilde{U} = U, \tilde{D} = D, \tilde{V} = V$  (apart from degeneracies due to subspaces). say all singular values are distinct then it is indeed unique.)

Thm: [Eckert Young]

Sps  $\sum_{i=1}^r \sigma_i u_i v_i^T = A_{n \times d}$ . For any  $k \leq r$ , define  $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ .

Then for any matrix  $B$  of rank at most  $k$ ,

$$\|A - B\|_F \geq \|A - A_k\|_F$$

Also: define the 2-norm or spectral norm of a matrix  $A$  to be  $\|A\|_2 = \max_{\|v\|=1} \|Av\|$   
 $= \max \text{ singular value of } A$ .

$$\Rightarrow \|A - B\|_2 \geq \|A - A_k\|_2$$

See HW exercises.

~~The top  $k$  singular values give the best  $k$ -dimensional subspace.~~

Lemma: the rows of  $A_k$  are the projections of rows of  $A$  onto  $V_k = \text{span}(v_1, \dots, v_k)$ .

Pf: if  $\bar{a}$  is a row of  $A$ ,  $\bar{a}$ 's projection is given by  $\sum_{i=1}^k (\bar{a} \cdot v_i) \cdot v_i^T$

$$\Rightarrow \text{All these projections are given by } \sum_{i=1}^k (A v_i) \cdot v_i^T = \sum_{i=1}^k \sigma_i u_i v_i^T = A_k$$

$\Rightarrow \|A - A_k\|_F$  is the remaining mass/energy

Fact: for a (square) symmetric matrix  $B$ , can find a basis of eigenvalues

$x_i$  st  $Bx_i = \lambda_i x_i \quad \forall i = 1 \dots n$ .  $x_i$  are orthonormal.

$$\Rightarrow B = X \Lambda X^T \quad X = \begin{pmatrix} | & & | \\ x_1 & & x_n \\ | & & | \end{pmatrix} \quad \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \end{pmatrix}$$

$$= \sum \lambda_i x_i x_i^T \quad (\text{again, unique upto rotation within subspaces})$$

Now: sps matrix is general  $A_{n \times d}$ . then

$$AA^T = (UDV^T)(UDV^T)^T = UD^2U^T \leftarrow \text{eigenvalues are } \sigma_i^2$$

eigenvectors are  $u_i$ .

$$A^T A = VD^2V^T \leftarrow \text{eigenvectors are now } v_i$$

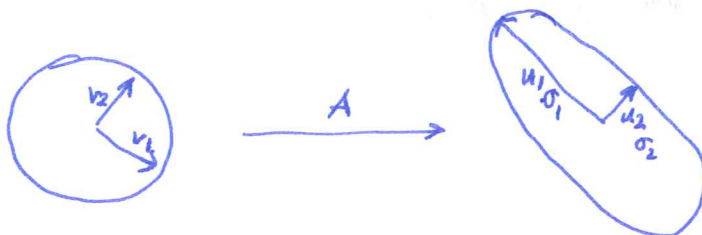
Moreover: having found (left) eigenvalues & eigenvectors (say)  $AA^T = X \Lambda X^T$  know  $U = X$ .

$$\text{and set } V^T = D^{-1} U^T A, \quad D^2 = \Lambda$$

One view,  $v_i$ 's are the directions along which  $A$  has most energy.

Another view:  $A = UDV^T \Rightarrow AV = UD$ .

cols of both  $U$  &  $V$  are orthonormal. So  $A$  takes these orthonormal cols of  $V$  into scaled orthonormal cols of  $U$ .



Example: handwriting recognition; face compression, etc.

~~Learning Structure of Gaussian~~

Topic Models (naive version)

$$A = \text{documents} \times \text{words} = UDV^T$$
$$= \begin{matrix} \text{docs} \\ \left( \begin{matrix} \end{matrix} \right) \\ \text{topics} \end{matrix} \begin{matrix} \uparrow \\ \text{weights of} \\ \text{topics} \end{matrix} \begin{matrix} \left( \begin{matrix} \end{matrix} \right) \\ \text{words} \\ \text{topics} \end{matrix}$$

Learning Gaussian Mixtures: [Venkatesh Wang]

Given  $k$  Gaussians, want to (a) cluster the points into  $k$  clusters, and with weights  $w_i$  (b) find the mean, variance, weight etc of the Gaussians. (this is fine given the clustering).

Known: if  $\frac{\text{inter-center}}{A}$  distances are at least  $\Omega(d^{1/4})$  then can cluster.

So: idea: find the SVD of space, with enough samples, the top  $k$ -dimensions

~~the~~ give us the subspace containing the  $k$  centers. Now: ~~this space~~ project down onto this space still gives us  $k$  Gaussians. (b/c of spherical symmetry) And hence a separation of  $\Omega(k^{1/4})$  suffices

## Subspace Embeddings:

Suppose  $W \subseteq \mathbb{R}^n$  is a linear subspace: an  $\varepsilon$ -subspace embedding <sup>into  $m$  dimensions</sup> ~~ensures that~~  
 $\Pi: \mathbb{R}^n \rightarrow \mathbb{R}^m$  ( $\Pi$  is a  $m \times n$  matrix) ensures that

$$1 - \varepsilon \leq \frac{\|\Pi x\|_2}{\|x\|_2} \leq 1 + \varepsilon \quad \forall x \in W.$$

So it maintains all ~~distances~~ vectors within subspace.

[Fact: JL applied to a fine enough net of points in  $W$  will give us such a property.]

[Suppose  $W$  is given by column space of matrix  $A \in \mathbb{R}^{n \times d}$ . Then can use SVDs to find an embedding of  $W$  into  $\mathbb{R}^d$ . How?

write  $A = UDV^T$  hence  $U^T A = DV^T$

$\uparrow$   $\uparrow$   $\uparrow$   
 $n \times r$   $r \times r$   $r \times d$   
 $(r \leq d)$

Any vector  $x \in W$  is a linear combination of columns of  $A$ . And  $U$  <sup>has</sup> orthonormal columns so  $U^T A$  rotates the columns of  $A$  into a  $r$ -dim subspace.  $\|U^T x\| = \|x\| \quad \forall x \in \text{span}(\text{col}(A))$

[N.b. this is a <sup>exact</sup>  $O$ -subspace embedding]

## Least Squares Regression:

Given  $A \in \mathbb{R}^{n \times d}$ ,  $b \in \mathbb{R}^n$  want an  $x \in \mathbb{R}^d$  st  $\min \|Ax - b\|$ .

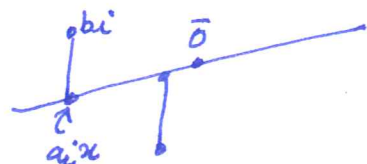
(Useful for solving overconstrained systems.)

$$= \min \sum_1^r (a_i^T x - b_i)^2$$

$\uparrow$   
 defines the subspace

Observe that  $Ax$  is in column span of  $A = UDV^T$ .

so  $Ax^*$  is the projection of  $b$  onto this column span.



$$\text{Proj}(b) = UU^T b$$

$\uparrow$   
 onto  $\text{colspan}(A)$

$\uparrow$  because  $U$  is <sup>orthonormal</sup> a basis for  $\text{colspan}(A)$ .

$$\text{(alternatively, proj onto colspan}(A) = A(A^T A)^{-1} A^T$$

$$= UDV^T (VDU^T UDV^T)^{-1} (UDV^T)^T = UU^T)$$

$$\text{hence } x^* = VD^{-1}U^T(UU^T b) = VD^{-1}U^T b.$$

$\Rightarrow$  least squares regression in  $SVD(n,d)$  time.

BTW: If  $A = UDV^T$

then  $VD^+U^T$  is called the pseudoinverse of  $A$  denoted by  $A^+$

( $D^+$  is really  $D^+$ , zeros stay 0)

and it ~~maps~~ ensures that for all vectors in the column space of  $A$ ,  $A^+$  acts like an inverse! ~~etc.~~

Formally: Restricted to  $\text{Colspace}(A^T)$  and  $\text{Colspace}(A)$  respectively, the operators  $A$  and  $A^+$  are inverses of each other.

Pf: sps  $x \in \text{Colspace}(A^T) \Rightarrow x = A^T y$  for some  $y$ .

$$\begin{aligned}
 \text{then } A^+Ax - x &= A^+A(A^T y) - A^T y \\
 &= VD^+U^T \cdot UDV^T(VDU^T y) - VDU^T y \\
 &= VD^+DDU^T y - VDU^T y \\
 &= V(D^+D - I)DU^T y = 0
 \end{aligned}$$

(Another motivation, of course, from the least squares) this multiplication gives zero.  $\square$

The pseudoinverses play an important role in Laplacians.

Recall:  $L(G) = \text{diag}(\text{degree vector}) - A$

but  $L(G) \cdot \mathbf{1} = 0$ . since the sum of all rows gives 0.

For connected graph, the <sup>values</sup> eigenvalues of  $L(G)$  are  $0 = \lambda_0 < \lambda_1 < \dots < \lambda_n$   
↑  
corresponding to eigenvector  $\mathbf{1}$ .

So can solve linear systems of the form  $Lx = b$  as long as  $b \perp \mathbf{1}$   
ie.  $\sum_v b_v = 0$ .  
by  $x = L^+b$ . (gives potentials corresponding to electrical flow demands)

Fact: Also, can do calculations and show:  $L^+L = LL^+ = I_n - \frac{1}{n}J$ .

$\Rightarrow$  for any vector  $b \perp \mathbf{1}$ ,  $LL^+b = b$  since  $b^T \mathbf{1} = 0$ .