

Singular Value Decomposition. (SVD)

- SVD is a matrix factorization method that is useful for many applications, e.g. recommendation systems, least squares problems, etc.

- Using SVD, we can check:

- rank of a matrix
- a given matrix near a simpler one (e.g. a matrix of smaller rank)

- There are many algorithms to compute SVD but most of them are expensive. ("slow to compute").

⇒ How to compute a good approximation to the SVD of a big matrix fast is an active research topic in numerical linear algebra!

Consider $A \in \mathbb{R}^{m \times n}$.

"The image of the unit sphere under any $m \times n$ matrix is a hyperellipsoid."

For simplicity, assume $m \geq n$ and $\text{rank}(A) = n$.
(Full rank, tall matrix).

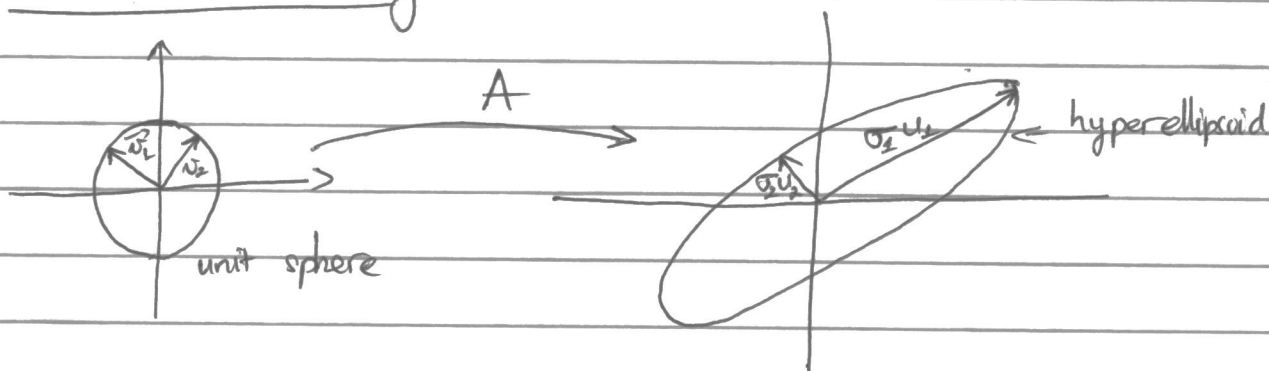
• Def: The singular values of A are the lengths of the n principal semi-axes of the hyperellipsoid $(A(\text{sphere}))$. σ_i

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0.$$

Def: The n left singular vectors of A are the (orthonormal) unit vectors in \mathbb{R}^n along the principal semi-axes of A (sphere). We denote them by $\{u_1, \dots, u_n\}$.
 $\Rightarrow \sigma_i u_i$ is the i th largest principal semi-axes of A (sphere).

Def: The n right singular vectors of A are the (orthonormal) unit vectors $\{v_1, \dots, v_n\}$ such that
 $Av_i = \sigma_i u_i, \quad i = 1, \dots, n.$

Geometric meaning:



• Reduced SVD:

$$Av_i = \sigma_i u_i \quad i = 1, \dots, n.$$

$$[Av_1 \quad Av_2 \quad \dots \quad Av_n] = [u_1 \quad u_2 \quad \dots \quad u_n]$$

$$A [v_1 \quad v_2 \quad \dots \quad v_n] = [u_1 \quad u_2 \quad \dots \quad u_n] \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix}$$

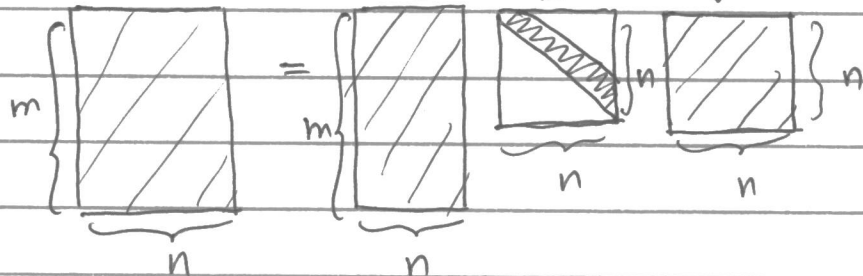
$$\begin{array}{ccc} A & V & = & U & \Sigma \\ \downarrow & \downarrow & & \downarrow & \downarrow \\ m \times n & n \times n & & m \times n & n \times n \end{array}$$

columns of U
are orthonormal.

$$\Rightarrow AV = U\Sigma$$

Since V is an orthogonal matrix,

$$A = U\Sigma V^T \quad \text{The reduced SVD of } A$$



* Pseudoinverse via SVD:

$$A^+ = V\Sigma^+ U^T$$

where $\Sigma^+ = \begin{bmatrix} 1/\sigma_1 & & \\ & \ddots & \\ & & 1/\sigma_n \end{bmatrix}$

[Exer: Show that $AA^+ = UU^T$
and $A^+A = VV^T$.]

The Moore - Penrose conditions: For a given matrix $A \in \mathbb{R}^{m \times n}$, if $B \in \mathbb{R}^{n \times m}$ satisfies the following:

i) $ABA = A$

ii) $BAB = B$

iii) $(AB)^T = AB$

iv) $(BA)^T = BA$

Then B is called the pseudoinverse (or the Moore - Penrose inverse) of A and written as A^+ .

Thm: (SVD) Any $m \geq n$ matrix A , with $m \geq n$, can be factorized

(Full version) $A = U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^T$

where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal

$\Sigma \in \mathbb{R}^{n \times n}$ is diagonal

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0.$$

Remark: There are accurate algorithms for computing the SVD.

Proposition: The 2-norm of a matrix is given by $\|A\|_2 = \sigma_1 \leftarrow$ the largest singular value of A .

Pf: (Sketch)

1) The norm is invariant under orthogonal transformations

$$\Rightarrow \|A\|_2 = \|\Sigma\|_2.$$

2) 2-norm of a diagonal matrix is equal to the absolute value of the largest diagonal element.

* Matrix properties via SVD.

Let $A \in \mathbb{R}^{m \times n}$.

$$p = \min\{m, n\}.$$

$r = \#$ of nonzero singular values \neq

$$(\Rightarrow r \leq p)$$

• Thm: $\text{rank}(A) = r$

$$\text{range}(A) = \text{span}\{u_1, \dots, u_r\}.$$

$$\text{Null}(A) = \text{span}\{u_{r+1}, \dots, u_n\}.$$

$$\|A\|_2 = \sigma_1 \quad \text{and} \quad \|A\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$$

• Thm: If $A^T = A$, $\sigma_i(A) = |\lambda_i(A)|$.
↑
ith eigenvalue
of A .

• Thm: For $A \in \mathbb{R}^{m \times m}$, (square matrix)
 $|\det(A)| = \prod_{i=1}^m \sigma_i$

Pf: Recall.

1) $\det(AB) = \det(A) \det(B)$

2) $\det(A^T) = \det(A)$

3) $\det(\text{diag}(a_1, \dots, a_n)) = \prod_{i=1}^n a_i$.

4) For any orthogonal Q , $|\det(Q)| = 1$.

$$\begin{aligned} \Rightarrow \det(A) &= \det(U \Sigma V^T) = \det(U) \det(\Sigma) \det(V^T) \\ &= \det(\Sigma) = \prod_{i=1}^m \sigma_i \end{aligned}$$

* Outer product:

Let $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$.

Then the outer product between u and v is:

$$uv^T = \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} [v_1 \quad \dots \quad v_n] = \begin{bmatrix} u_1 v_1 & u_1 v_2 & \dots & u_1 v_n \\ u_2 v_1 & u_2 v_2 & \dots & u_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ u_m v_1 & u_m v_2 & \dots & u_m v_n \end{bmatrix} \in \mathbb{R}^{m \times n}$$

This matrix has rank 1 because

$$\vec{u}\vec{v}^T = [v_1\vec{u} \quad v_2\vec{u} \quad \dots \quad v_n\vec{u}]$$

each column is just a scalar multiple of the same vector \vec{u} .

\Rightarrow we can view a matrix A as:

$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T, \quad r = \text{rank}(A)$$

sum of rank-1 matrices.

(Exer: show this)

* Best Rank-k Approximations:

Let $A \in \mathbb{R}^{n \times d}$. (That is, the rows of A are n points in d dimensional space).

Suppose that $\text{rank}(A) = r$. Then the SVD of A is

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T.$$

For $k \in \{1, 2, \dots, r\}$, let

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T.$$

A_k ~~is~~ called the sum truncated after k terms.

Fact: A_k has rank k (why?).

Lemma: The rows of A_k are the projections of the rows of A onto the subspace V_k spanned by the first k singular vectors of A , (i.e., $V_k = \text{span}\{v_1, \dots, v_k\}$).

pf: Let \vec{x} be an arbitrary vector in \mathbb{R}^d .

Since \vec{v}_i are orthonormal, the orthogonal projection of \vec{x} onto V_k is given by

$$\sum_{i=1}^k \langle \vec{v}_i, \vec{x} \rangle \vec{v}_i$$

$$= \left(\sum_{i=1}^k v_i v_i^T \right) \vec{x}$$

$\Rightarrow P = \sum_{i=1}^k v_i v_i^T$ is the projection matrix onto V_k .

Let $\vec{a}_1, \dots, \vec{a}_n$ be the columns of A^T .

\Rightarrow their transposes are the rows of A .

Then

$$PA^T = P \begin{bmatrix} \vec{a}_1 & \vec{a}_2 & \dots & \vec{a}_n \end{bmatrix} \\ = \begin{bmatrix} P\vec{a}_1 & P\vec{a}_2 & \dots & P\vec{a}_n \end{bmatrix}.$$

Transpose both sides:

$$P^T = P \leftarrow AP \stackrel{!}{=} \begin{bmatrix} (P\vec{a}_1)^T \\ \vdots \\ (P\vec{a}_n)^T \end{bmatrix}$$

$\Rightarrow AP$ is the matrix whose rows are the projections of the rows of A onto V_k .

Since $P = \sum_{i=1}^k v_i v_i^T$

$$AP = \sum_{i=1}^k A v_i v_i^T = \sum_{i=1}^k \sigma_i u_i v_i^T = A_k.$$

since $A v_i = \sigma_i u_i$.

Thm: For any matrix B of rank at most k .

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

That is, A_k is the best rank k approximation to A , where error is measured in the Frobenius norm.

(See textbook: Foundation of Data science
Theorem 3.6, p. 48).

pf. Let B minimize $\|A - B\|_F^2$ among all rank k or less matrices.

Let V be the space spanned by the rows of B .
 $\Rightarrow \dim(V) \leq k$.

Let a_i^T be the i th row of A .

b_i^T be the i th row of B .

$$\Rightarrow \|A - B\|_F^2 = \sum_{i=1}^n \|a_i^T - b_i^T\|_2^2$$

Since $b_i \in V$, $\|a_i^T - b_i^T\|_2^2$ is minimized if b_i^T is the projection of a_i^T onto V .

\Rightarrow each row of B is the projection of the corresponding row of A onto V .

Then $\|A - B\|_F^2$ is the sum of squared distance of rows of A to V . Since A_k minimizes the sum of squared distance of rows of A to any k -dim subspace, it by the previous lemma, $B = A_k$.

* Applications:

• Image Compression:

As each image is represented by a matrix, it is possible to compress an image by approximating it by a lower rank matrix.

* Power Method for SVD

Computing the SVD is an important research topic in ~~numb~~ numerical linear algebra.

power method is a basic method to establish and the approximate SVD of a matrix A in polynomial time.

Let A be matrix whose SVD is $\sum_{i=1}^r \sigma_i u_i v_i^T$.

Let $B = A^T A$. Then

$$B = A^T A \\ = \left(\sum_i \sigma_i \cancel{u_i} \cancel{u_i^T} \right) \left(\sum_j \sigma_j u_j v_j^T \right)$$

$$= \sum_{i,j} \sigma_i \sigma_j v_i u_i^T u_j v_j^T$$

$$\Rightarrow B = \sum_i \sigma_i^2 v_i v_i^T \quad \text{since } u_i \text{ are orthonormal.}$$

Facts: 1) The matrix B is square and symmetric, and has the same left and right-singular vectors.

2) v_j is an eigenvector of B with eigenvalue σ_j^2 .

$$B v_j = \left(\sum_i \sigma_i^2 v_i v_i^T \right) v_j$$

$$v_i^T v_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} = \sum_i \sigma_i^2 v_i v_i^T v_j$$

$$= \sigma_j^2 v_j v_j^T v_j$$

$$\Rightarrow B v_j = \sigma_j^2 v_j$$

Now compute B^2 ,

$$B^2 = \left(\sum_i \sigma_i^2 u_i u_i^T \right) \left(\sum_j \sigma_j^2 u_j u_j^T \right)$$

$$= \sum_{i,j} \sigma_i^2 \sigma_j^2 u_i u_i^T u_j u_j^T$$

$$= \sum_i \sigma_i^4 u_i u_i^T$$

Fact: $B^k = \sum_i \sigma_i^{2k} u_i u_i^T$

If $\sigma_1 > \sigma_2$, when k is large ($k \rightarrow \infty$), the first term $\sigma_1^{2k} u_1 u_1^T$ in the summation dominates.

$$\Rightarrow B^k \rightarrow \sigma_1^{2k} u_1 u_1^T \quad \text{as } k \rightarrow \infty$$

→ To find

⇒ To estimate u_1 : one can compute B^k , then normalize the first column of B^k .

But in practice, A may be a very large, sparse matrix, E.g. Netflix Rating matrix dramas.

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

in practice, A may be a $10^8 \times 10^8$ matrix
with 10^9 nonzero entries.

Though A is sparse, B need not be and in the
worse case may have all 10^{16} entries nonzero.

\Rightarrow computing B^k is very costly.

Faster way:

(randomly) choose a vector x and compute $B^k x$.
This is a little bit faster than computing B^k
because you do matrix-vector multiplication:

compute Bx = this is a vector-matrix-vector.

Then compute $B(Bx)$

and so on.

How to get v_i ?

Note that $x = \sum_{i=1}^d c_i v_i$ (suppose that B is
of full rank)

Since $\{v_1, \dots, v_d\}$ forms a basis for \mathbb{R}^d .

Then for k large

$$B^k x \approx (\sigma_i^{2k} v_i v_i^T) \left(\sum_{i=1}^d c_i v_i \right) = \sigma_i^{2k} c_i v_i$$

$$\Rightarrow B^k x \approx \sigma_i^{2k} c_i v_i$$

\Rightarrow normalizing $B^k x$, we can obtain an approximate
vector of v_i .

* SVD and condition number:

Recall the condition number for a square nonsingular matrix A is defined by

$$\kappa(A) := \|A\|_2 \|A^{-1}\|_2.$$

— $\kappa(A)$ small, then A is said well-conditioned.

— $\kappa(A)$ large, then A is said ill-conditioned



Finding the solution of $Ax = b$
by using computers may ~~be~~ give us a bad approximation

How is the condition number related to SVD?

Recall that

$$\|A\|_2 = \sigma_1.$$

and

$$\|A^{-1}\|_2 = \frac{1}{\sigma_m} \quad \text{where}$$

why? $A = U \Sigma V^T$

and $A^{-1} = V \Sigma^{-1} U^T$

$$= V \operatorname{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_m}\right) U^T.$$

since $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m > 0$,

$$\frac{1}{\sigma_m} \geq \frac{1}{\sigma_{m-1}} \geq \dots \geq \frac{1}{\sigma_1}.$$

↑
largest singular value of A^{-1}

$$\Rightarrow \|A^{-1}\|_2 = \frac{1}{\sigma_m}.$$

We can generalize the definition of the condition number for a rectangular matrix $A \in \mathbb{R}^{m \times n}$ using the pseudo-inverse A^+ and SVDs.

$$\begin{aligned} \kappa(A) &= \|A\|_2 \|A^+\|_2 \\ &= \frac{\sigma_1}{\sigma_r} \quad \text{where } r = \text{rank}(A). \end{aligned}$$