

## Lecture 9: Singular Value Decomposition.

Recall: If  $A \in \mathbb{R}^{n \times n}$  is symmetric, then

$$A = [u] [D] [v^T]$$

and if  $u_1, \dots, u_n$  are columns of  $A$ , then

- 1). diagonal entries of  $D$  are eigenvalues of  $A$ ,
- 2).  $u_i$  are eigenvectors of  $A$ , and are orthonormal.

$$[u_1 \dots u_n] [\lambda_1 \dots \lambda_n] [u_1^T \dots u_n^T]$$

$$= \underbrace{\lambda_1 u_1 u_1^T + \lambda_2 u_2 u_2^T + \dots + \lambda_n u_n u_n^T}_{n \times n \text{ matrix whose } i,j\text{-th entry is } u_i u_j}$$

In this lecture: the generalization of this to general matrices.

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

SVD for short

Thm: (Singular Value Decomposition)

Any  $A \in \mathbb{R}^{m \times n}$  can be written as

$$m \begin{bmatrix} A \\ \vdots \end{bmatrix} = \underbrace{n \begin{bmatrix} u \\ \vdots \end{bmatrix}}_{U} \underbrace{\begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & & \end{bmatrix}}_S \underbrace{\begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}}_V$$

where:

1.  $U \in \mathbb{R}^{m \times m}$  is orthogonal ( $u_1, \dots, u_n$  form orthonormal basis)
2.  $S$  is "diagonal",  $\sigma_i > 0 \ \forall i$ .
3.  $V \in \mathbb{R}^{n \times n}$  is also orthogonal

$\sigma_1 \geq \dots \geq \sigma_n > 0$  are called the singular values of  $A$ .

$u_1, \dots, u_m$  are left singular vectors of  $A$

$v_1, \dots, v_n$  are right singular vectors of  $A$ .

Equivalently:

$$A = \sum_{i=1}^{\min(m,n)} \sigma_i u_i v_i^\top$$

What is the action of  $A$ ?

$A v_i = \sigma_i u_i$ ,  $\forall i$ . so  $A$  maps  $\{v_1, \dots, v_n\}$  to  $\{u_1, \dots, u_n\}$ , and scales the  $i$ th coordinate by  $\sigma_i$ .

Similarly,  $u_i^\top A = \sigma_i v_i$  so  $A^\top$  does the "opposite" mapping, with the same scaling.

Connection to eigenvectors / eigenvalues, PCA

SVD is for general matrices, but spectral decomposition is only for symmetric.

However, for symmetric, SVD follows directly from spectral:

If  $A \in \mathbb{R}^{n \times n}$  is symmetric, then

$$A = U D U^\top$$

If all entries of  $D$  are  $> 0$ , then this is an SVD

$A$  is positive semi-definite (PSD)

otherwise, take

$$A = \begin{bmatrix} u_1 & \dots & u_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} \text{sign}(\lambda_1) u_1^\top \\ \vdots \\ \text{sign}(\lambda_n) u_n^\top \end{bmatrix}$$

and this is an SVD.

On the other hand, if  $A \in \mathbb{R}^{m \times n}$  is arbitrary, we can relate its

SVD to the Spectra of some related matrices:

$$[A] = [u_1 \dots u_m] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$

$$A^T A = [v_1 \dots v_n] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} [u_1 \dots u_m] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$

$\underbrace{\quad\quad\quad}_{= I \text{ since } \{u_i\} \text{ are orthonormal}}$

$$= [v_1 \dots v_n] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$

$$\underbrace{\quad\quad\quad}_{= \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}}$$

$$= [v_1 \dots v_n] \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$

This is spectral decomposition of  $A^T A$ !

so right singular vectors of  $A \Leftrightarrow$  eigenvectors of  $A^T A$   
 (nonzero) singular values of  $A \Leftrightarrow \sqrt{\text{eigenvalues of } A^T A}$ .

Similarly, left singular vectors of  $A \Leftrightarrow$  eigenvectors of  $A A^T$   
 " " " " if  $A A^T$ .

Relationship w/ PCA:

Note that in PCA, given data matrix  $X$ , the PC of  $X$  are top 2 eigenvectors of  $X^T X$ .

So PCA is just pretty much SVD!

Strictly speaking, PCA only cares about the right singular vectors of  $X$ , so SVD is strictly more general than PCA.  
 But often are used interchangeably.

Algos for SVD : Given  $A \in \mathbb{R}^{m \times n}$ , how to find SVD?

Idea: just use power method on  $A^T A$  and/or  $\bar{A} \bar{A}^T$ .

Once we've found right singular vectors  $\{v_1, \dots, v_n\}$ , left singular vectors are given by

$$u_1 = A v_1$$

$$\vdots$$
  
$$u_n = A v_n$$

$\begin{matrix} u_{n+1} \\ \vdots \\ u_m \end{matrix} \} \rightarrow$  any basis of space perpendicular to  $u_1, \dots, u_n$ .

Note: to apply power method to  $A^T A$ , we don't need to compute  $A^T A$ , which can be expensive. Instead:

$A^T A u = A^T (A u)$  so only need to do 2 matrix-vector multiplies, which is usually faster.

best algo in practice:

np.linalg.svd(A) (or equivalent pkg).

Application: Low-rank approximation.

Q: How can we fill in the following matrix?

$$A = \begin{bmatrix} 7 & ? & ? \\ ? & 8 & ? \\ ? & 12 & 6 \\ ? & ? & 2 \\ 21 & 6 & ? \end{bmatrix} \Rightarrow \begin{bmatrix} 7 & 2 & 1 \\ 56 & 8 & 4 \\ 42 & 12 & 6 \\ 28 & 4 & 2 \\ 21 & 6 & 3 \end{bmatrix}$$

Obviously impossible in worst case!

But what if  $A$  has structure? i.e. what if all rows are multiples of each other?

An example of a low-rank matrix.

Def:  $A \in \mathbb{R}^{m \times n}$  has rank 0 if  $A =$  all zeros.

Def:  $A \in \mathbb{R}^{m \times n}$  has rank 1 if  $A = u v^T$

$$\begin{bmatrix} u \end{bmatrix} \begin{bmatrix} v^T \end{bmatrix}$$

Def: A has rank k if A can be written as a sum of k rank 1 matrices, and cannot be written as a sum of k-1.

$$A = \sum_{i=1}^k u_i v_i^T. \quad \begin{matrix} \leftarrow \text{not necessarily unit vectors} \\ \text{or orthogonal.} \end{matrix}$$

$$A = \begin{bmatrix} u_1 & \dots & u_n \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_k^T \end{bmatrix}$$

"short, long"  
"tall, skinny"

Many equivalent definitions of rank:

1. The largest set of linearly independent columns of A has size k.
2. The largest set of linearly independent rows has size k.
3. A has k non-zero singular values.

Low-rank approximation.

Real-world data is unlikely to be exactly low-rank.

We can still ask for best rank-k approximation.

Q: Given A, find B which has rank-k that is "closest" to A.

A natural candidate: let

$$A = \sum_{i=1}^{\min(m,n)} \sigma_i u_i v_i^T \text{ be the SVD of } A$$

(recall  $\sigma_1 \geq \dots \geq \sigma_n \geq 0$ )

and take  $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$  as your rank-k approx.

It turns out that in many natural ways, this is optimal.

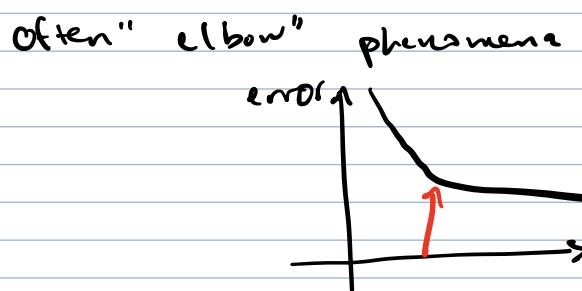
Thm: For any  $A \in \mathbb{R}^{m \times n}$ , any B rank k:

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

$$(\|M\|_F = \sqrt{\sum_{ij} M_{ij}^2} \text{ is } \underline{\text{Frobenius norm}} \text{ of } A).$$

Closely related to optimality of PCA!

How to choose  $k$ ?  $k$  trades off size of representation vs quality.



Another rule of thumb:

choose  $k$  s.t.

$$\sum_{i=1}^k \sigma_i \geq C \cdot \sum_{i=k+1}^n \sigma_i$$

3? 10?

Applications:

1. Compression
2. Denoising
3. Data completion.