# Low-Rank Matrix Approximations, SVD and \_\_\_\_\_ Connections to PCA

Archer Yang

McGill University

October 4, 2024

# Reading assignment: read chapter 4.5-4.6 of Mathematics for Machine Learning book.

# LOW-RANK MATRIX

### **Applications**

Data compression

Feature learning

## **Data compression**







The size of the four compressed images is 45Kb, 345Kb, 1.1Mb and 3.4Mb<sup>1</sup>.

<sup>1</sup> Multivariate Statistics, Richard Wilkinson

#### Feature representation learning



Let's review the concept of matrix rank.

**Rank-**0 **matrix** There is only one rank-zero matrix of a given size, namely the all-zero matrix.

**Rank-1 matrix** all rows are multiples of each other. In the example in (4), all columns are also multiples of each other; this is not an accident. An equivalent definition of a rank-1  $m \times n$  matrix is as the outer product  $uv^{\top}$  of an *m*-vector *u* and an *n*-vector *v*.

$$\boldsymbol{A} = \boldsymbol{u}\boldsymbol{v}^{\top} = \begin{pmatrix} u_1\boldsymbol{v}^{\top} \\ u_2\boldsymbol{v}^{\top} \\ \vdots \\ u_m\boldsymbol{v}^{\top} \end{pmatrix} = \begin{pmatrix} v_1\boldsymbol{u} & v_2\boldsymbol{u} & \cdots & v_n\boldsymbol{u} \end{pmatrix}$$

Note that each row is a multiple of  $v^{\top}$ , and each column is multiple of u.

**Rank-2 matrix** A rank-two matrix is just a sum of two rank-1 matrices

$$A = u_{1}v_{1}^{\top} + u_{2}v_{2}^{\top} = \begin{pmatrix} u_{11}v_{1}^{\top} + u_{21}v_{2}^{\top} \\ u_{12}v_{1}^{\top} + u_{22}v_{2}^{\top} \\ \vdots \\ u_{1m}v_{1}^{\top} + u_{2m}v_{2}^{\top} \end{pmatrix}$$
$$= \begin{pmatrix} | & | \\ u_{1} & u_{2} \\ | & | \end{pmatrix} \begin{pmatrix} - & v_{1}^{\top} & - \\ - & v_{2}^{\top} & - \end{pmatrix} = \begin{pmatrix} | & | \\ u_{1} & u_{2} \\ | & | \end{pmatrix} \begin{pmatrix} | & | & | \\ v_{1} & v_{2} \\ | & | \end{pmatrix} \begin{pmatrix} | & | \\ v_{1} & v_{2} \\ | & | \end{pmatrix}^{\top}$$

**Note:** a rank-2 matrix is one that can be written as the sum of two rank-1 matrices and is not itself a rank-0 or rank-1 matrix.

**Rank-**r matrix A rank-r matrix can be written as the sum of r rank-1 matrices, and **cannot** be written as the sum of r - 1 or fewer rank-1 matrices.

$$\boldsymbol{A} = \boldsymbol{u}_{1}\boldsymbol{v}_{1}^{\top} + \dots + \boldsymbol{u}_{r}\boldsymbol{v}_{r}^{\top} = \begin{pmatrix} u_{11}\boldsymbol{v}_{1}^{\top} + \dots + u_{r1}\boldsymbol{v}_{r}^{\top} \\ u_{12}\boldsymbol{v}_{1}^{\top} + \dots + u_{r2}\boldsymbol{v}_{r}^{\top} \\ \vdots \\ u_{1m}\boldsymbol{v}_{1}^{\top} + \dots + u_{rm}\boldsymbol{v}_{r}^{\top} \end{pmatrix}$$
$$= \begin{pmatrix} | & | \\ u_{1} & \dots & u_{r} \\ | & | \end{pmatrix} \begin{pmatrix} - \boldsymbol{v}_{1}^{\top} & - \\ \vdots \\ - \boldsymbol{v}_{r}^{\top} & - \end{pmatrix} = \underbrace{\begin{pmatrix} | & | \\ u_{1} & \dots & u_{r} \\ | & | \end{pmatrix}}_{\boldsymbol{U}} \underbrace{\begin{pmatrix} | & | \\ v_{1} & \dots & v_{r} \\ | & | \end{pmatrix}}_{\boldsymbol{V}^{\top}}$$

Rank-r matrix A can be factored into the product of a "tall and skinny" matrix U and a "short and fat" matrix V<sup>T</sup>

$$\mathbf{M} \times \mathbf{n} \qquad \mathbf{M} \times \mathbf{r} \mathbf{r} \times \mathbf{n}$$
$$\mathbf{A} = \mathbf{U} \mathbf{V}^{\mathsf{T}}$$

Equivalent definitions:

- The largest linearly independent subset of columns of A has size r. i.e., all n columns of A arise as linear combinations of only r of them.
- The largest linearly independent subset of rows of *A* has size *r*.



## Low-Rank Matrix Approximations: Motivation

**Low-rank approximation:** To approximate a rank-r matrix A with a rank-k matrix A(k) ( $k \le r$ )

$$\boldsymbol{A} \approx \boldsymbol{A}(k)$$

or

$$\min_{\boldsymbol{A}(k)} \|\boldsymbol{A} - \boldsymbol{A}(k)\|$$

Applications:

**Compression**:  $m \times n$  numbers in A; k(m + n) numbers in A(k)

De-noising

$$A = A(k) + error$$

# SINGULAR VALUE DECOMPOSITION

# Singular Value Decomposition (SVD)

- To solve low-rank matrix approximation problem, we apply the singular value decomposition (SVD), a central matrix decomposition method in linear algebra.
- It applied to **all matrices**, not only to square matrices.

## Singular Value Decomposition (SVD)

#### Theorem 4.22 MML (SVD Theorem)

Let  $A \in \mathbb{R}^{m \times n}$  be a rectangular matrix. The SVD of A is a decomposition of the form

$$A = U \Sigma V^{\top}$$

- $U = (u_1, \dots, u_m)$  is an  $m \times m$  orthogonal matrix,  $U^{\top}U = I_{m \times m}$ . Columns of U are left singular vectors.
- $V = (v_1, \dots, v_n)$  is an  $n \times n$  orthogonal matrix,  $V^\top V = I_{n \times n}$ . Columns of V are right singular vectors.
- $\Sigma$  is an  $m \times n$  diagonal matrix with (ordered) nonnegative entries. Diagonal entries  $\sigma_1 \ge \cdots \ge \sigma_n \ge 0$  are singular values.

# SVD (rank-r)



SVD (full column rank rank(A) = n)



# SVD vs Eigen-decomposition

In contrast to the eigen-decomposition

$$A = V \Lambda V^{-1}$$

in SVD:

- Orthogonal matrices U and V are not the same, U and V need not even have the same dimension.
- A need not be square.

Reduced SVD (rank(A) =  $r \le n$ )



MXY

rxr

Reduced SVD (full column rank rank(A) = n)



21

#### **Reduced SVD for Rank-***r* **matrix**

If rank(A) = r, the SVD expresses A as a nonnegative linear combination of r rank-1 matrices,

$$\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^{\mathsf{T}} = \sum_{i=1}^{m} \sum_{j=1}^{n} \sigma_{ij} \mathbf{u}_{i} \mathbf{v}_{j}^{\mathsf{T}}$$

$$= \sum_{i=1}^{r} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{\mathsf{T}} = \sum_{i=1}^{r} \sigma_{i} \mathbf{A}_{i}$$
(1)

where  $\sigma_{ij}$  is the *i*-th row *j*-th column of  $\Sigma$ .  $u_i$  and  $v_i$  are the corresponding left and right singular vectors.

#### **Reduced SVD for Rank-***r* **matrix**

Eq. (1) can be further simplified

$$A = \sum_{i=j}^{r} \sigma_{ij} u_i v_j^{\top} + \sum_{i\neq j}^{r} \sigma_{ij} u_i v_j^{\top}$$
$$= \sum_{i=j\leq r}^{r} \sigma_{ij} u_i v_j^{\top} + \underbrace{\sum_{i=j>r}^{r} \sigma_{ij} u_i v_j^{\top}}_{=0} + \underbrace{\sum_{i\neq j}^{r} \sigma_{ij} u_i v_j^{\top}}_{=0}$$
$$= \sum_{i=1}^{r} \sigma_i u_i v_i^{\top}$$
$$\equiv \sum_{i=1}^{r} \sigma_i A_i$$

where  $\sigma_i \equiv \sigma_{ii}$  and  $A_i = u_i v_i^{\top}$  is a rank-1 matrix.

#### Feature representation learning



**Question:** why we observe vertical and horizontal lines in  $A_1, \dots, A_5$ ?

#### **Reduced SVD for rank-***r* **matrix**

Let  $A \in \mathbb{R}^{m \times n}$  be a rank-*r* matrix, where  $1 \le r \le \min(n, p)$ . The reduced SVD of A is

$$\boldsymbol{A} = \boldsymbol{U}_r \boldsymbol{\Sigma}_r \boldsymbol{V}_r^{\top} = \sum_{i=1}^{\prime} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^{\top}$$

$$\mathbf{I}$$
  $U_r = (u_1, \dots, u_r)$  is an  $m \times r$  orthogonal matrix,  $U_r^{\top} U_r = I_{r \times r}$ .

**V**<sub>r</sub> =  $(v_1, \dots, v_r)$  is an  $n \times r$  orthogonal matrix,  $V_r^\top V_r = I_{r \times r}$ .

■  $\Sigma_r$  is an  $r \times r$  diagonal matrix with singular values  $\sigma_1 \ge \cdots \ge \sigma_r > 0$ (all positive).

# LOW-RANK MATRIX APPROXIMATION

## Rank-k approximation

How do we approximate a rank-r matrix A by a rank-k matrix  $(k \le r)$ ?

Recall that in reduced SVD, a rank-r matrix A can be expressed as

$$\boldsymbol{A} = \boldsymbol{U}_{r}\boldsymbol{\Sigma}_{r}\boldsymbol{V}_{r}^{\mathsf{T}} = \sum_{i=1}^{r}\sigma_{i}\boldsymbol{u}_{i}\boldsymbol{v}_{i}^{\mathsf{T}}$$
(2)

where  $\sigma_1 \geq \cdots \geq \sigma_r > 0$ 

**Rank-***k* approximation of *A*: keep only the first *k* terms in (2)

$$\widehat{\boldsymbol{A}}(k) = \sum_{i=1}^{k} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^{\mathsf{T}}$$

we see that  $rank(\widehat{A}(k)) = k$  and

$$\boldsymbol{A} - \widehat{\boldsymbol{A}}(k) = \sum_{i=k+1}^{r} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^{\top}$$

#### Algorithm

- **1** Compute the SVD  $A = U\Sigma V^{\top}$ , where  $\Sigma$  is a nonnegative  $m \times n$  diagonal matrix with diagonal entries sorted from high to low.
- 2 Keep only the top k right singular vectors: Set V<sub>k</sub><sup>T</sup> equal to the first k rows of V<sup>T</sup> (a k × n matrix)
- 3 Keep only the top k left singular vectors: Set  $U_k$  equal to the first k rows of U (a  $m \times k$  matrix)
- 4 Keep only the top k singular values: Set Σ<sub>k</sub> equal to the first k rows and columns of Σ (a k × k matrix)
- 5 The rank-k approximation is then

$$\widehat{\boldsymbol{A}}(k) = \boldsymbol{U}_k \boldsymbol{\Sigma}_k \boldsymbol{V}_k^\top = \sum_{i=1}^k \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^\top$$



(a) Original image A.

(b) Rank-1 approximation  $\widehat{A}(1)$ .(c) Rank-2 approximation  $\widehat{A}(2)$ .



(d) Rank-3 approximation  $\widehat{A}(3)$ .(e) Rank-4 approximation  $\widehat{A}(4)$ .(f) Rank-5 approximation  $\widehat{A}(5)$ .



The size of the four compressed images is 45Kb, 345Kb, 1.1Mb and 3.4Mb $^2.$  More demos are here.  $^3$ 

<sup>&</sup>lt;sup>2</sup>*Multivariate Statistics*, Richard Wilkinson

<sup>&</sup>lt;sup>3</sup>https://timbaumann.info/svd-image-compression-demo/

# Quantifying approximation error

- How to quantify the approximation error  $||A \widehat{A}(k)||$ ?
- Need the notion of a matrix norm  $\|\cdot\|$ .
- e.g. the Frobenious norm or the spectral norm.

#### **Frobenius norm**

Let  $A \in \mathbb{R}^{m \times n}$ . The Frobenius norm of A is

$$\|\boldsymbol{A}\|_{F} = \sqrt{\sum_{i=1}^{m} \sum_{j=1}^{n} |a_{ij}|^{2}} = \sqrt{\operatorname{tr}(\boldsymbol{A}^{\top}\boldsymbol{A})}$$

(Homework) Show that for any matrix A and orthogonal matrices U and V:

$$\|\boldsymbol{U}^{\mathsf{T}}\boldsymbol{A}\|_{F}^{2} = \|\boldsymbol{A}\|_{F}^{2} \qquad \|\boldsymbol{A}\boldsymbol{V}\|_{F}^{2} = \|\boldsymbol{A}\|_{F}^{2} \qquad (3)$$

• (Homework) Combining the result in (3) and using the SVD  $A = U\Sigma V^{T}$  to show that

$$\|\boldsymbol{A}\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}$$

where  $r \equiv \operatorname{rank}(A)$  and  $\sigma_i$  are the singular values of A.

#### **Spectral norm**

Let  $A \in \mathbb{R}^{m \times n}$ . The spectral norm of A is defined as

$$\|A\|_{2} = \max_{x} \frac{\|Ax\|_{2}}{\|x\|_{2}} = \max_{x:\|x\|=1} \|Ax\|_{2}$$

#### **Spectral norm**

(Homework) Let 
$$A \in \mathbb{R}^{m \times n}$$
. Show that

 $\|\boldsymbol{A}\|_2 = \sigma_1$ 

where  $\sigma_1$  is the first (largest) singular value of A.

## **Eckart-Young-Mirsky Theorem**

For either the spectral norm  $\|\cdot\|_2$  or the Frobenious norm  $\|\cdot\|_F$ , the rank-*k* approximation

$$\widehat{\boldsymbol{A}}(k) = \boldsymbol{U}_k \boldsymbol{\Sigma}_k \boldsymbol{V}_k^{\top} = \sum_{i=1}^k \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^{\top}$$

minimizes the approximation error

$$\|A - \widehat{A}(k)\| \le \|A - C\|$$
 for all rank-k matrices C

In other words

 $\widehat{A}(k) = \underset{\operatorname{rank}(C)=k}{\operatorname{arg\,min}} ||A - C||$ 

Moreover,

$$\|\boldsymbol{A} - \widehat{\boldsymbol{A}}(k)\| = \begin{cases} \sigma_{k+1} & \text{for the } \|\cdot\|_2 \text{ norm} \\ \sqrt{\sum_{i=k+1}^r \sigma_i^2} & \text{for the } \|\cdot\|_F \text{ norm} \end{cases}$$

# **RELATIONSHIP WITH EIGEN-DECOMPOSITION**

#### Relation to eigenvalue decomposition

Given  $\boldsymbol{X}$  as a  $N \times p$  matrix. If  $\boldsymbol{X}$  has SVD

Then it follows

$$S = \frac{1}{N} X^{\mathsf{T}} X = \frac{1}{N} (U \Sigma V^{\mathsf{T}})^{\mathsf{T}} U \Sigma V^{\mathsf{T}}$$
$$= \frac{1}{N} V \Sigma^{\mathsf{T}} U^{\mathsf{T}} U \Sigma V^{\mathsf{T}} = V \frac{\Sigma^{\mathsf{T}} \Sigma}{N} V^{\mathsf{T}}$$

where  $\boldsymbol{\Sigma}^{\mathsf{T}}\boldsymbol{\Sigma} = \operatorname{diag}(\sigma_1^2, \cdots, \sigma_r^2, 0, \cdots, 0) \in \mathbb{R}^{p \times p}$  with  $r = \operatorname{rank}(\boldsymbol{X})$ .

Meanwhile, by the eigen-decomposition

$$S = V \Lambda V^{\top}$$

#### Relation to eigenvalue decomposition

Therefore

$$\boldsymbol{V}\left(\frac{\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma}}{N}\right)\boldsymbol{V}^{\top} = \boldsymbol{S} = \boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{\top}$$

- Right-singular vectors of X (the columns of V) are eigenvectors of  $S = \frac{1}{N} X^{\top} X$ .
- The non-zero eigenvalues of S are related to the non-zero singular values of X via

$$\mathbf{\Lambda} = \frac{\mathbf{\Sigma}^{\top} \mathbf{\Sigma}}{N} \quad \text{i.e.} \quad \lambda_d = \frac{\sigma_d^2}{N} \quad \text{for } d = 1, \cdots, p$$

### Relation to eigenvalue decomposition

We also have

$$\frac{1}{N} \boldsymbol{X} \boldsymbol{X}^{\top} = \frac{1}{N} \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\top} (\boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\top})^{\top}$$
$$= \frac{1}{N} \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\top} \boldsymbol{V} \boldsymbol{\Sigma}^{\top} \boldsymbol{U}^{\top} = \frac{1}{N} \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{\top} \boldsymbol{U}^{\top}$$

■ Therefore, left-singular vectors of X (the columns of U) are eigenvectors of  $\frac{1}{N}XX^{\top}$ .

# PCA AS LOW-RANK MATRIX APPROXIMATION

#### PCA as Low-rank Matrix Approximation

PCA solves

$$\min_{\boldsymbol{B}} \frac{1}{N} \| \boldsymbol{X}^{\top} - \boldsymbol{B} \boldsymbol{B}^{\top} \boldsymbol{X}^{\top} \|_{F}^{2} \equiv \frac{1}{N} \| \boldsymbol{X} - \boldsymbol{X} \boldsymbol{B} \boldsymbol{B}^{\top} \|_{F}^{2}$$
subject to  $\boldsymbol{B}^{\top} \boldsymbol{B} = \boldsymbol{I}_{k \times k}$ 

On the other hand, rank-k matrix approximation of X solves

$$\widehat{oldsymbol{X}}(k) = \operatorname*{arg\,min}_{\mathsf{rank}(oldsymbol{C})=k} \|oldsymbol{X}-oldsymbol{C}\|_F \quad \mathsf{where} \ \widehat{oldsymbol{X}}(k) = oldsymbol{U}_k oldsymbol{\Sigma}_k oldsymbol{V}_k^ op$$

• One can show that  $V_k = B$ ,  $U_k \Sigma_k = XB$  and

$$\boldsymbol{U}_k \boldsymbol{\Sigma}_k \boldsymbol{V}_k^{\top} = \boldsymbol{X} \boldsymbol{B} \boldsymbol{B}^{\top}$$

### PCA as Low-rank Matrix Approximation

This relation suggest that SVD can also be used to produce PCA for X (*k*-dim. subspace), as follows

- **1** Perform rank-k matrix approximation of  $X \approx \widehat{X}(k) = U_k \Sigma_k V_k^{\top}$ .
- 2 Let *B* denote the top *k* principal components of  $S = \frac{1}{N}X^{\top}X$ , set  $B = V_k = (v_1, \dots, v_k).$
- **3** (Homework) The **PC scores** can be obtained by  $XB = U_k \Sigma_k$ .
  - e.g. for the *n*-th observation, the projected coordinates on the subspace spanned by v<sub>1</sub>, ···, v<sub>k</sub> are

$$(z_{n1}, \dots, z_{nk}) = (\boldsymbol{v}_1^{\mathsf{T}} \boldsymbol{x}_n, \dots, \boldsymbol{v}_k^{\mathsf{T}} \boldsymbol{x}_n) = (\boldsymbol{v}_1^{\mathsf{T}} \boldsymbol{x}_n,$$

# **Computation complexity**

Computing PCA for  $X \in \mathbb{R}^{N \times p}$  using SVD has less computation complexity then using eigen-decomposition

#### Eigen-decomposition approach:

- **S** covariance matrix computation  $O(Np\min(N, p))$
- Eigen-decomposition  $O(p^3)$
- Total:  $O(Np\min(N,p)) + O(p^3)$

#### SVD approach:

SVD  $O(Np\min(N, p))$ 

# MATRIX COMPLETION VIA LOW RANK

**APPROXIMATION** 

## What Are The Missing Entries?

Consider the following matrix,

$$\boldsymbol{X} = \begin{bmatrix} 7 & ? & ? \\ ? & 8 & ? \\ ? & 12 & 6 \\ ? & ? & 2 \\ 21 & 6 & ? \end{bmatrix}$$

Matrix completion problem: What are the missing entries?

Consider an extreme assumption: that all rows are multiples of each other.

$$\boldsymbol{X} = \begin{bmatrix} 7 & ? & ? \\ ? & 8 & ? \\ ? & 12 & 6 \\ ? & ? & 2 \\ 21 & 6 & ? \end{bmatrix}$$

Under such assumption, here is the completed matrix

$$\widehat{\boldsymbol{X}} = \begin{bmatrix} 7 & 2 & 1 \\ 28 & 8 & 4 \\ 42 & 12 & 6 \\ 14 & 4 & 2 \\ 21 & 6 & 3 \end{bmatrix}$$
(4)

When you know something about the "structure" of a partially known matrix, then sometimes it's possible to recover all of the "lost" information.

# Matrix completion using matrix approximation

#### Applications:

#### Matrix completion:

 First impute X to obtain X'. Fill the missing entries of X using one of the following methods

**0**;

the average of the known entries in the same column;

- the average of the known entries in the same row;
- the average of the known entries of the matrix.
- 2 Compute the best rank-k approximation to X'

$$\widehat{X}(k) = \underset{\operatorname{rank}(C)=k}{\operatorname{arg\,min}} \|X' - C\|$$

• Why not we just do the first step to use X'?

# SVD FOR HIGH-DIMENSIONAL DATA

## **SVD (***m* < *n***, rank-***r***)**



SVD (m < n, rank(A) = m)



# Reduced SVD (m < n, rank-r)



53

Reduced SVD (m < n, rank(A) = m)

