

# SI231b: Matrix Computations

## Lecture 21: Low-rank Approximation and Regularized Least Square

Yue Qiu

[qiuyue@shanghaitech.edu.cn](mailto:qiuyue@shanghaitech.edu.cn)

School of Information Science and Technology  
ShanghaiTech University

Nov. 25, 2021

## Original Image

- ▶ Let  $A \in \mathbb{R}^{m \times n}$  be a matrix whose  $(i, j)$ th entry  $a_{ij}$  represents the  $(i, j)$ th pixel of an image.
- ▶ memory consumption for storing  $A$ :  $m * n$

## Compressed Image

- ▶ using truncated SVD of  $A$ : store  $\{u_i, \sigma_i v_i\}_{i=1}^k$  instead of the full  $A$ .
- ▶ the compressed image is represented by  $B = \sum_i^k \sigma_i u_i v_i^T$
- ▶ memory consumption for truncated SVD:  $(m + n) * k$ 
  - much less than  $m * n$  if  $k \ll \min\{m, n\}$

# Image Compression Illustration

original image, sizes  $470 \times 641$



Figure 1: original image

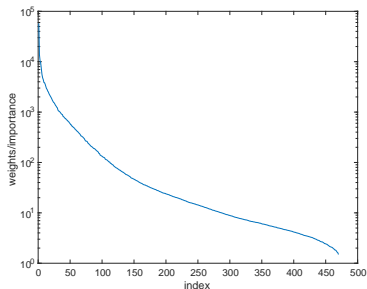


Figure 2: singular values

# Image Compression Illustration

compressed image with  $r = 10$



compressed image with  $r = 20$



compressed image with  $r = 30$



compressed image with  $r = 40$



**Aim:** given a matrix  $A \in \mathbb{R}^{m \times n}$  and an integer  $k$  with  $0 \leq k \leq \text{rank}(A)$ , find a matrix  $B \in \mathbb{R}^{m \times n}$  such that  $\text{rank}(B) \leq k$  and  $B$  best approximates  $A$

- ▶ it is somehow unclear about what a “best approximation” means, and we will specify one later
- ▶ applications: PCA, dimensionality reduction, . . . . . the same kind of applications in matrix factorization
- ▶ **truncated SVD:** denote

$$A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where the  $k$ th “partial sum” captures as much of the energy of  $A$  as possible, and the meaning of “energy” will be specified later

- ▶ then perform the aforementioned approximation by choosing  $B = A_k$

Truncated SVD provides the best approximation in the LS sense:

**Theorem**[Eckart-Young-Mirsky]. Consider the following problem

$$\min_{B \in \mathbb{R}^{m \times n}, \text{rank}(B) \leq k} \|A - B\|_F^2$$

where  $A \in \mathbb{R}^{m \times n}$  and  $k \in \{1, \dots, p\}$  with  $p = \min\{m, n\}$  are given. The truncated SVD  $A_k$  is an optimal solution to the above problem and the minimum is  $\sum_{i=k+1}^p \sigma_i^2$

- ▶ also note the matrix 2-norm version of the Eckart-Young-Mirsky theorem:

$$\min_{B \in \mathbb{R}^{m \times n}, \text{rank}(B) \leq k} \|A - B\|_2^2$$

The truncated SVD  $A_k$  is an optimal solution to the above problem and the minimum is  $\sigma_{k+1}^2$

(cf. Theorem 2.4.8 in [Golub & van Loan 13'])

- ▶ the energy mentioned before is defined by either the Frobenius norm or the 2-norm

# Low-rank Factorization Approximation

In practice, we are more interested in the **factorized form** of low-rank approximation,

$$\min_{A \in \mathbb{R}^{m \times k}, B \in \mathbb{R}^{k \times n}} \|Y - AB\|_F^2$$

where  $k \leq \min\{m, n\}$ ; A denotes a basis matrix; B is the coefficient matrix.

- the matrix factorization problem may be reformulated as (verify)

$$\min_{Z \in \mathbb{R}^{m \times n}, \text{rank}(Z) \leq k} \|Y - Z\|_F^2,$$

and the truncated SVD  $Y_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ , where  $Y = U \Sigma V^T$  denotes the SVD of Y, is an optimal solution by the **Eckart-Young-Mirsky** theorem.

- thus, an optimal solution to the matrix factorization problem is given by

$$A = [u_1, \dots, u_k], \quad B = [\sigma_1 v_1, \dots, \sigma_k v_k]^T$$

Similar to variational characterization of eigenvalues of real symmetric matrices, we can derive various variational characterization results for singular values, e.g.,

- ▶ Courant-Fischer characterization:

$$\sigma_k(A) = \min_{\dim S_{n-k+1} \subseteq \mathbb{R}^n} \max_{x \in S_{n-k+1}, \|x\|_2=1} \|Ax\|_2$$

- ▶ Weyl's inequality: for any  $A, B \in \mathbb{R}^{m \times n}$ ,

$$\sigma_{k+l-1}(A+B) \leq \sigma_k(A) + \sigma_l(B), \quad k, l \in \{1, \dots, p\}, \quad k+l-1 \leq p.$$

Also, note the corollaries

- $\sigma_k(A+B) \leq \sigma_k(A) + \sigma_1(B)$ ,  $k = 1, \dots, p$
- $|\sigma_k(A+B) - \sigma_k(A)| \leq \sigma_1(B)$ ,  $k = 1, \dots, p$  (important results of perturbation theory)

- ▶ and many more...



## Applying Weyl's inequality

- ▶ for any  $B$  with  $\text{rank}(B) \leq k$ , we have
  - $\sigma_l(B) = 0$  for  $l > k$
  - (Weyl)  $\sigma_{i+k}(A) \leq \sigma_i(A - B) + \sigma_{k+1}(B) = \sigma_i(A - B)$  for  $i = 1, \dots, p - k$
  - and consequently

$$\|A - B\|_F^2 = \sum_{i=1}^p \sigma_i(A - B)^2 \geq \sum_{i=1}^{p-k} \sigma_i(A - B)^2 \geq \sum_{i=k+1}^p \sigma_i(A)^2$$

- ▶ the equality above is attained if we choose  $B = A_k$

# Advantages of Using Low-rank Factorized Form

Let  $A \in \mathbb{R}^{m \times n}$  being approximated by  $B = UV^T$  with  $U \in \mathbb{R}^{m \times r_k}$ ,  $V \in \mathbb{R}^{n \times r_k}$  and  $r_k \ll \{m, n\}$ , i.e.,  $B \approx A$ .

## Computational Complexity Reduction

- ▶ matrix-vector product with  $z \in \mathbb{R}^n$ 
  - $\mathcal{O}(mn)$  for  $Az$
  - $\mathcal{O}(r_k(m+n))$  for  $Bz$
- ▶ matrix-matrix product with  $Z \in \mathbb{R}^{n \times n}$ 
  - $\mathcal{O}(mn^2)$  for  $AZ$
  - $\mathcal{O}(r_k(m+n)n)$  for  $BZ$

## Memory Consumption Reduction

- ▶  $\mathcal{O}(mn)$  for  $A$
- ▶  $\mathcal{O}(r_k(m+n))$  for  $B$

# Key Ingredients for Using Low-rank Approximation

The key of low-rank approximation lies in the fact that

- ▶ all computations should be performed using low-rank factors  $U$  and  $V$  rather than the explicit  $B = UV^T$
- ▶ the rank  $r_k \ll \{m, n\}$

## Rank Growth

In computations, to keep the results in factorized form, the rank will increase. For example, for  $m \times n$  matrices  $B = U_1 V_1^T$ ,  $C = U_2 V_2^T$  and to compute  $B + C$ , we have

$$D = B + C = U_b V_b^T + U_c V_c^T = \underbrace{\begin{bmatrix} U_b & U_c \end{bmatrix}}_{U_d} \underbrace{\begin{bmatrix} V_b^T \\ V_c^T \end{bmatrix}}_{V_d^T}.$$

The rank of  $D$  turns to be  $r_b + r_c$  in the general case and continues growing when more computations are performed.

We need to **reduce the rank** for less computational complexity.

**Keeping the rank bounded is the key** in applying low-rank approximation for computations.

For an  $m \times n$  matrix  $A = UV^T$  with low-rank factors  $U \in \mathbb{R}^{m \times r}$  and  $V \in \mathbb{R}^{n \times r}$ , the following procedure returns a best rank  $r'$  of  $A$  with  $r' < r$

1. compute a reduced QR factorization of  $U$ , i.e.,  $U = QR$  with  $Q \in \mathbb{R}^{m \times r}$  and  $R \in \mathbb{R}^{r \times r}$  ( $\mathcal{O}(r^2 m)$  cost)
2. form  $C = RV^T$  with  $C \in \mathbb{R}^{r \times n}$  ( $\mathcal{O}(r^2 n)$  cost)
3. compute the SVD of  $C$ , i.e.,  $C = \begin{bmatrix} U_c^{(1)} & U_c^{(2)} \end{bmatrix} \begin{bmatrix} \Sigma_c^{(1)} & \\ & \Sigma_c^{(2)} \end{bmatrix} \begin{bmatrix} (V_c^{(1)})^T \\ (V_c^{(2)})^T \end{bmatrix}$   
with  $U_c^{(1)}$  having  $r'$  columns ( $\mathcal{O}(r^2 n)$  cost)
4.  $\tilde{A} = QU_c^{(1)}\Sigma_c^{(1)}(V_c^{(1)})^T$  returns the best rank  $r'$  approximation of  $A$

Can you prove the optimality?

# Summary of Low-rank Approximation

We have seen from the previous analysis that the key to keep the computational complexity low using low-rank approximation is

- ▶ using low-rank factorized form
- ▶ reducing the increased rank while performing computations

To perform computations using low-rank approximations, we need to **start with low-rank factorized form**,

- ▶ may be already given
- ▶ using SVD to compute (**one time cost**)
- ▶ using **randomized algorithm** to find one if SVD is too expensive, cf. the following reference by Caltech
  - N. Halko, P. G. Martinsson, and J. A. Tropp. Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions. *SIAM Review*, vol. 53, pp. 217–288, 2011.

We have introduced the low-rank approximation using SVD in this lecture, which in turn gives **optimal** results. Other related low-rank approximation methods which are **less accurate but computationally cheaper** include

- ▶ CUR factorization  $A \approx CUR$  where  $C$  is from columns of  $A$ ,  $R$  contains rows of  $A$ ;
- ▶ skelton/cross approximation;
- ▶ nonnegative matrix factorization (NMF) (widely used in NLP)

For high dimensional data, tensor computations are used.

**Question:** how sensitive is the LS solution when there is noise?

$$y = A\bar{x} + \nu,$$

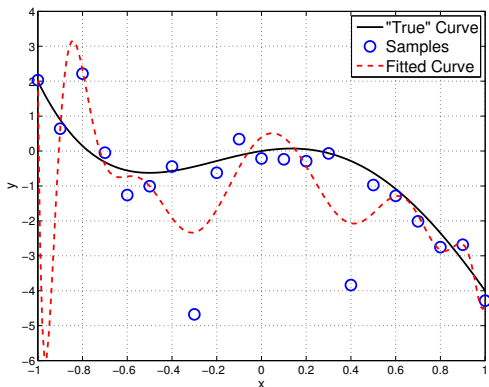
where  $\bar{x}$  is the true result;  $A \in \mathbb{R}^{m \times n}$  has full column rank;  $\nu$  is noise, modeled as a random vector, for example with mean zero and covariance  $\gamma^2 I$  (white noise).

**Mean square error (MSE) analysis:** from  $x_{LS} = A^\dagger y = \bar{x} + A^\dagger \nu$  we get

$$\begin{aligned} E[\|x_{LS} - \bar{x}\|_2^2] &= E[\|A^\dagger \nu\|_2^2] = E[\text{tr}(A^\dagger \nu \nu^T (A^\dagger)^T)] = \text{tr}(A^\dagger E[\nu \nu^T] (A^\dagger)^T) \\ &= \gamma^2 \text{tr}(A^\dagger (A^\dagger)^T) \\ &= \gamma^2 \sum_{i=1}^n \frac{1}{\sigma_i^2(A)} \end{aligned}$$

**Observation:** the MSE becomes very large if some  $\sigma_i(A)$ 's are close to zero.

## Example: Curve Fitting



The same curve fitting example in [Lecture 7](#). The “true” curve is the true  $f(x)$  with polynomial order  $n = 4$ . In practice, the model order may not be known and we may have to guess. The fitted curve above is done by LS with a guessed model order  $n = 16$ .



**Intuition:** replace  $x_{LS} = (A^T A)^{-1} A^T y$  by

$$x_{RLS} = (A^T A + \lambda I)^{-1} A^T y,$$

for some  $\lambda > 0$ , where the term  $\lambda I$  is added to improve the conditioning of the system, i.e., **move the singular values of  $A^T A$  away from zero**, thereby attempting to reduce noise sensitivity.

How may we make sense out of such a modification?

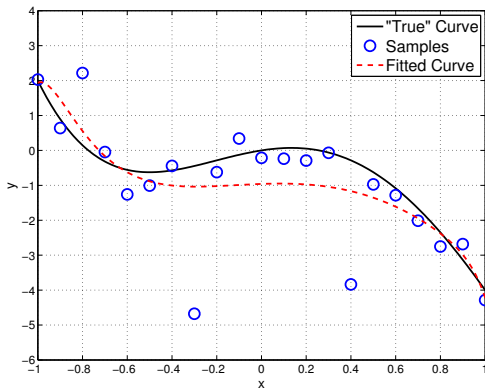
**$\ell_2$ -regularized LS:** find an  $x$  that solves

$$\min_{x \in \mathbb{R}^n} \|Ax - y\|_2^2 + \lambda \|x\|_2^2$$

for some predetermined  $\lambda > 0$ .

- ▶ the solution is uniquely given by  $x_{RLS} = (A^T A + \lambda I)^{-1} A^T y$
- ▶ the formulation says that we try to minimize both  $\|y - Ax\|_2^2$  and  $\|x\|_2^2$ , and  $\lambda$  controls which one should be more emphasized in the minimization

## Example: Curve Fitting Using $\ell_2$ -Regularization



The fitted curve is done by  $\ell_2$ -regularized LS with a guessed model order  $n = 18$  and with  $\lambda = 0.1$ .

If you are interested in the modified least squares problems and the their solution via SVD, you are [suggested](#) to read

- ▶ Gene H. Golub and Charles F. Van Loan. *Matrix Computations*, *Johns Hopkins University Press*, 2013.

Chapter 6.1 – 6.4.