

SI231b: Matrix Computations

Lecture 20: Low-rank Approximation and Regularized Least Square

Yue Qiu

qiuyue@shanghaitech.edu.cn

School of Information Science and Technology
ShanghaiTech University

Nov. 16, 2022

Motivation Example: Image Compression

Original Image

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

Compressed Image

- ▶ using truncated SVD of \mathbf{A} : store $\{\mathbf{u}_i, \sigma_i \mathbf{v}_i\}_{i=1}^k$ instead of the full \mathbf{A} .
- ▶ the compressed image is represented by $\mathbf{B} = \sum_i^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
- ▶ 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 470×641



Figure 1: original image

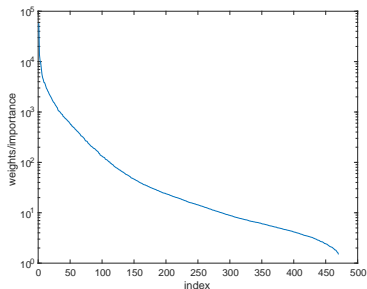


Figure 2: singular values

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

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

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

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

- ▶ then perform the aforementioned approximation by choosing $\mathbf{B} = \mathbf{A}_k$

Truncated SVD provides the best approximation in the LS sense:

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

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

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $k \in \{1, \dots, p\}$ with $p = \min\{m, n\}$ are given. The truncated SVD \mathbf{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:

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

The truncated SVD \mathbf{A}_k is an optimal solution to the above problem and the minimum is σ_{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,

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

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

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

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

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

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

$$\mathbf{A} = [\mathbf{u}_1, \dots, \mathbf{u}_k], \quad \mathbf{B} = [\sigma_1 \mathbf{v}_1, \dots, \sigma_k \mathbf{v}_k]^T$$

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:

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

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

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

Also, note the corollaries

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

- ▶ and many more...

Applying Weyl's inequality

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

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

- ▶ the equality above is attained if we choose $\mathbf{B} = \mathbf{A}_k$

Advantages of Using Low-rank Factorized Form

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

Computational Complexity Reduction

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

Memory Consumption Reduction

- ▶ $\mathcal{O}(mn)$ for \mathbf{A}
- ▶ $\mathcal{O}(r_k(m+n))$ for \mathbf{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 \mathbf{U} and \mathbf{V} rather than the explicit $\mathbf{B} = \mathbf{UV}^T$
- ▶ 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 $\mathbf{B} = \mathbf{U}_1 \mathbf{V}_1^T$, $\mathbf{C} = \mathbf{U}_2 \mathbf{V}_2^T$ and to compute $\mathbf{B} + \mathbf{C}$, we have

$$\mathbf{D} = \mathbf{B} + \mathbf{C} = \mathbf{U}_b \mathbf{V}_b^T + \mathbf{U}_c \mathbf{V}_c^T = \underbrace{\begin{bmatrix} \mathbf{U}_b & \mathbf{U}_c \end{bmatrix}}_{\mathbf{U}_d} \underbrace{\begin{bmatrix} \mathbf{V}_b^T \\ \mathbf{V}_c^T \end{bmatrix}}_{\mathbf{V}_d^T}.$$

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

1. compute a reduced QR factorization of \mathbf{U} , i.e., $\mathbf{U} = \mathbf{QR}$ with $\mathbf{Q} \in \mathbb{R}^{m \times r}$ and $\mathbf{R} \in \mathbb{R}^{r \times r}$ ($\mathcal{O}(r^2 m)$ cost)
2. form $\mathbf{C} = \mathbf{RV}^T$ with $\mathbf{C} \in \mathbb{R}^{r \times n}$ ($\mathcal{O}(r^2 n)$ cost)
3. compute the SVD of \mathbf{C} , i.e., $\mathbf{C} = \begin{bmatrix} \mathbf{U}_c^{(1)} & \mathbf{U}_c^{(2)} \end{bmatrix} \begin{bmatrix} \Sigma_c^{(1)} & \\ & \Sigma_c^{(2)} \end{bmatrix} \begin{bmatrix} (\mathbf{V}_c^{(1)})^T \\ (\mathbf{V}_c^{(2)})^T \end{bmatrix}$
with $\mathbf{U}_c^{(1)}$ having r' columns ($\mathcal{O}(r^2 n)$ cost)
4. $\tilde{\mathbf{A}} = \mathbf{QU}_c^{(1)}\Sigma_c^{(1)}(\mathbf{V}_c^{(1)})^T$ returns the best rank r' approximation of \mathbf{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 $\mathbf{A} \approx \mathbf{CUR}$ where \mathbf{C} is from columns of \mathbf{A} , \mathbf{R} contains rows of \mathbf{A} ;
- ▶ skelton/cross approximation;
- ▶ nonnegative matrix factorization (NMF) (widely used in NLP)

For high dimensional data, tensor computations are used.

For a **nonsingular** matrix \mathbf{A} , we are concerned with the solution of the linear system $\mathbf{Ax} = \mathbf{b}$.

Question: if there is a small perturbation in \mathbf{A} , what is the distance between the **perturbed solution** and **exact solution** \mathbf{x} ?

$$(\mathbf{A} + \Delta\mathbf{A})(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{b}$$

From **Lecture 6**, we know that

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|},$$

where $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ is defined as the condition number of the matrix \mathbf{A} and is denoted by $\kappa(\mathbf{A})$.

Note: $\kappa(\mathbf{A}) \geq 1$ (**how to prove?**)

When the matrix 2-norm is used,

$$\kappa_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A})}.$$

$\sigma_{\min}(\mathbf{A})$ measures the distance of \mathbf{A} to singularity. For orthogonal matrix \mathbf{A} , $\kappa_2(\mathbf{A}) = 1$.

When $\sigma_{\min}(\mathbf{A})$ is close to zero,

- ▶ $\kappa_2(\mathbf{A})$ gets large \rightsquigarrow small perturbation may lead to large solution error

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{b} = \sum_i \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i$$

- ▶ inverting \mathbf{A} gets more difficult and unstable

Note: $\kappa_2(\mathbf{A}^T\mathbf{A}) = \kappa_2(\mathbf{A}\mathbf{A}^T) = \kappa(\mathbf{A})^2$. (Can you prove this?)

This explains why forming problems with $\mathbf{A}^T\mathbf{A}$ or $\mathbf{A}\mathbf{A}^T$ is (almost) a bad idea.

Equivalence of Condition Number

The matrix \mathbf{A} is said to be **ill-conditioned** if $\kappa(\mathbf{A})$ is large. This statement is a norm dependent property.

Any two condition numbers $\kappa_\alpha(\cdot)$ and $\kappa_\beta(\cdot)$ are **equivalent** on $\mathbb{R}^{m \times n}$, which means that constants c_1 and c_2 can be found so that

$$c_1 \kappa_\alpha(\mathbf{A}) \leq \kappa_\beta(\mathbf{A}) \leq c_2 \kappa_\alpha(\mathbf{A}), \quad \forall \mathbf{A} \in \mathbb{R}^{m \times n}.$$

For example, for $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\begin{aligned} \frac{1}{n} \kappa_2(\mathbf{A}) &\leq \kappa_1(\mathbf{A}) \leq n \kappa_2(\mathbf{A}) \\ \frac{1}{n} \kappa_\infty(\mathbf{A}) &\leq \kappa_2(\mathbf{A}) \leq n \kappa_\infty(\mathbf{A}) \\ \frac{1}{n^2} \kappa_1(\mathbf{A}) &\leq \kappa_\infty(\mathbf{A}) \leq n^2 \kappa_1(\mathbf{A}) \end{aligned}$$

Therefore, if a matrix is ill-conditioned in the α -norm, it is also ill-conditioned in the β -norm.

Note: all vectors norms are equivalent and all matrix norms are also equivalent. (cf. Chapter 2.2 and 2.3 of [Golub & van Loan13'] for details.)

Recall from [Lecture 10](#), for $\mathbf{A} \in \mathbb{R}^{m \times n}$, the pseudoinverse of \mathbf{A} denoted by $\mathbf{A}^\dagger \in \mathbb{R}^{n \times m}$ satisfying the [Moore–Penrose conditions](#).

1. $\mathbf{A}\mathbf{A}^\dagger\mathbf{A} = \mathbf{A}$

2. $\mathbf{A}^\dagger\mathbf{A}\mathbf{A}^\dagger = \mathbf{A}^\dagger$

3. $(\mathbf{A}\mathbf{A}^\dagger)^T = \mathbf{A}\mathbf{A}^\dagger$

4. $(\mathbf{A}^\dagger\mathbf{A})^T = \mathbf{A}^\dagger\mathbf{A}$

- R. Penrose. A Generalized Inverse for Matrices. *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 51, pp. 406-413, 1955.

For a rank r matrix \mathbf{A} , its SVD is given by

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_1 & \mathbf{U}_2 \end{bmatrix} \begin{bmatrix} \tilde{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \end{bmatrix},$$

where $\tilde{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r)$, $\mathbf{U}_1 \in \mathbb{R}^{m \times r}$, $\mathbf{V}_1 \in \mathbb{R}^{n \times r}$. Then we get

$$\mathbf{A}^\dagger = \mathbf{V}_1 \tilde{\Sigma}^{-1} \mathbf{U}_1^T$$

Note: it is [not necessary](#) that $\mathbf{A}^\dagger\mathbf{A} = \mathbf{I}$ or $\mathbf{A}\mathbf{A}^\dagger = \mathbf{I}$

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

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

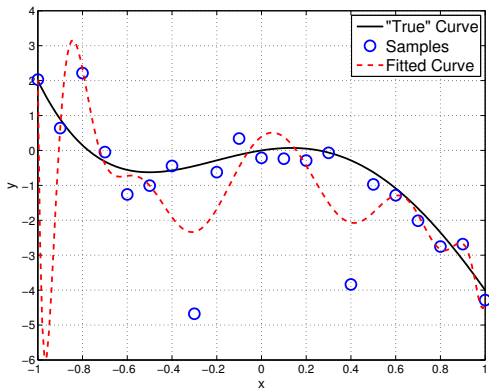
where $\bar{\mathbf{x}}$ is the true result; $\mathbf{A} \in \mathbb{R}^{m \times n}$ has full column rank; $\boldsymbol{\nu}$ is noise, modeled as a random vector, for example with mean zero and covariance $\gamma^2 \mathbf{I}$ (white noise).

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

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

Observation: the MSE becomes very large if some $\sigma_i(\mathbf{A})$'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 $\mathbf{x}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$ by

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

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

How may we make sense out of such a modification?

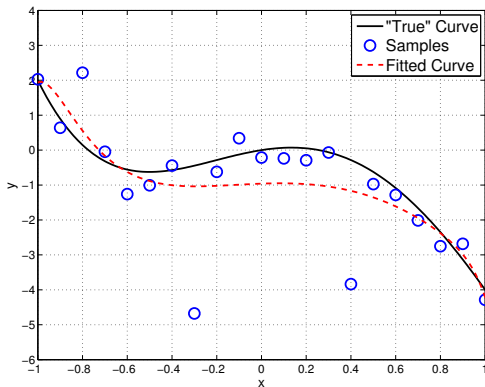
ℓ_2 -regularized LS: find an \mathbf{x} that solves

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

for some predetermined $\lambda > 0$.

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

Example: Curve Fitting Using ℓ_2 -Regularization



The fitted curve is done by ℓ_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.