

## Week 7-1: Power method

Lecturer: Bo Y.-C. Ning

May 10, 2022

**Disclaimer:** My notes may contain errors, distribution outside this class is allowed only with the permission of the Instructor.

## Last time

- PCA and SVD

## Today

- Power method

## 1 Review of singular value decomposition (SVD)

For a rectangular matrix  $A \in \mathbb{R}^{m \times n}$ , let  $p = \min\{m, n\}$ , then we have the SVD

$$A = U\Sigma V',$$

where  $U = (u_1, \dots, u_m)$  and  $V = (v_1, \dots, v_n)$  are orthogonal matrices and  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p)$  is a diagonal matrix such that  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ .  $\sigma_i$ s are called the *singular values*,  $u_i$ s are the left singular vectors and  $v_i$ s are the right singular vectors.

The matrix  $\Sigma$  is not a square matrix, one can define thin SVD, which factorizes  $A$  as

$$A = U_n \Sigma_n V' = \sum_{i=1}^n \sigma_i u_i v_i',$$

where  $U_n \in \mathbb{R}^{m \times n}$ ,  $U_n' U_n = I_n$ ,  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ . This is for  $m > n$ , if  $m < n$ , then we let  $V \in \mathbb{R}^{m \times n}$ ,

The following properties are useful: for  $\sigma(A) = (\sigma_1, \dots, \sigma_p)'$ , the rank of  $A$  is the number of nonzero singular values denoted as  $\|\sigma(A)\|_0$ . The Frobenius norm of  $A$ ,  $\|A\|_F = (\sum_{i=1}^p \sigma_i^2)^{1/2} = \|\sigma(A)\|_2$ , and the spectrum norm of  $A$ ,  $\|A\|_2 = \sigma_1 = \|\sigma(A)\|_\infty$ . Using the fact that  $U, V$  are both orthogonal matrices

$$\begin{aligned} A'A &= V\Sigma U'U\Sigma V' = V\Sigma^2 V', \\ AA' &= U\Sigma V'V\Sigma U' = U\Sigma^2 U' \end{aligned}$$

Last, the eigen-decomposition for a real symmetric matrix is  $B = W\Lambda W'$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ , which is the SVD of  $B$ .

## 2 Power method

To start, let's assume  $A \in \mathbb{R}^{n \times n}$  is a symmetric and p.s.d. matrix, the power method for obtaining the largest eigenvalue is given as:

- 1) Choose an initial guess of  $q^{(0)}$  (non-zero);
- 2) Repeat  $k = 1, \dots, K$ ,

$$\begin{aligned} z^{(k)} &= Aq^{(k-1)} \\ q^{(k)} &= \frac{z^{(k)}}{\|z^{(k)}\|_2}; \end{aligned}$$

- 3) Output:  $\lambda_1 \leftarrow q^{(K)'} A q^{(K)}$ .

## 3 Why the power method works?

Let's understand how the power method works. Before that, we need to recall a few facts:

- The eigenvalue  $v_i$  attached to  $i$ -th eigenvalue  $\lambda_i$  has the relation  $Av_i = \lambda_i v_i$
- Given  $A$ ,  $A = \sum_{i=1}^n \lambda_i v_i v_i'$ , where  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$  and  $\langle v_i, v_j \rangle = 0$  for  $i \neq j$
- $A^k = \sum_{i=1}^n \lambda_i^k v_i v_i'$ , why?

By inspecting the algorithm, we have

$$q^{(k)} = \frac{A^k q^{(0)}}{\|A^k q^{(0)}\|_2}.$$

Now given an initial guess of  $q^{(0)}$  of unit Euclidean norm, it is possible to express

$$q^{(0)} = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n,$$

for  $\alpha_1, \dots, \alpha_n$  are scalars. By the relation  $Av_i = \lambda_i v_i$ ,

$$A^k q^{(0)} = \alpha_1 \lambda_1^k \left( v_1 + \sum_{j=2}^n \frac{\alpha_j \lambda_j^k}{\alpha_1 \lambda_1^k} v_j \right)$$

For simplicity, let's denote  $y^{(k)} = \sum_{j=2}^n \frac{\alpha_j \lambda_j^k}{\alpha_1 \lambda_1^k} v_j$ , note that  $y^{(k)} \rightarrow 0$  as  $k \rightarrow \infty$  as long as  $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_n$  then

$$q^{(k)} = \frac{A^k q^{(0)}}{\|A^k q^{(0)}\|_2} = \frac{\alpha_1 \lambda_1^k (v_1 + y^{(k)})}{\|\alpha_1 \lambda_1^k (v_1 + y^{(k)})\|_2} \rightarrow v_1, \quad \text{as } k \rightarrow \infty$$

In practice,  $k$  will never goes to  $\infty$ , the algorithm will stop as some  $K$  when  $\min\{\|q^{(K)} - q^{(K-1)}\|_2, \|-q^{(K)} - q^{(K-1)}\|_2\} \leq \epsilon$  for some small  $\epsilon$ .

The output  $q^{(K)}$  is a close approximation of  $v_1$ , the leading eigenvector. How to obtain the leading eigenvalue  $\lambda_1$ ? (Hint: using  $Av_1 = \lambda_1 v_1$ ).

A few comments:

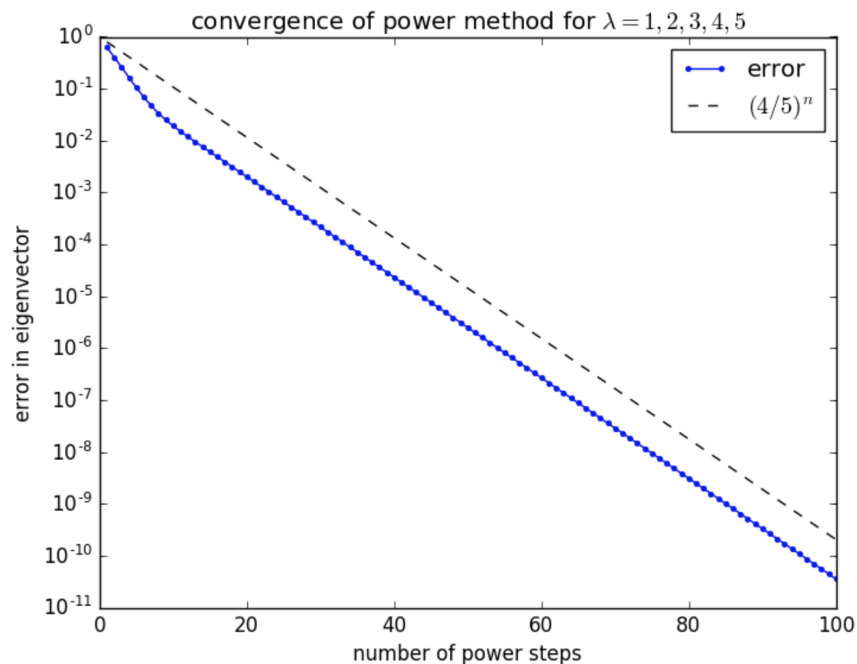


Figure 6.1: Convergence speed of the power method.

[Source: <https://web.mit.edu/18.06/www/Spring17/Power-Method.pdf>.]

- The power method works well if  $\lambda_1 > \lambda_2$ . It converges slowly if  $\lambda_1/\lambda_2 \approx 1$ .
- The convergence speed of the power method is proportional to  $(\lambda_2/\lambda_1)^k$ , the ratio between  $\lambda_2$  and  $\lambda_1$ .
- For a general matrix  $A^{n \times p}$ , we can apply the power method to  $A'A$  or  $AA'$  instead. Then the output is the absolute value of  $\lambda_1$ .
- How to get  $\lambda_2, \dots, \lambda_n$ ?
- Eigen-decomposition is implemented in LAPACK, see `eigen()` in R and `np.linalg.eig` in numpy.

The power method is the most basic algorithm for SVD. There are other methods such as

- Inverse power method for finding the eigenvalue of smallest absolute value (replace  $A$  with  $A^{-1}$  in the power method);
- QR algorithm for symmetric eigen-decomposition (takes  $4n^3/3$  for eigenvalues and  $8n^3/3$  for eigenvector)
- “Golub-Kahan-Reinsch” algorithm (Section 8.6 of Golub and Van Loan); used in `svd` function in R ( $4m^2n + 8mn^2 + 9n^3$  flops for an  $m > n$  matrix)
- Jacobi methods (Section 8.5 of Golub and Van Loan) (suitable for parallel computing).

Concluding remarks on numerical linear algebra:

- Numerical linear algebra forms the building blocks of most computation we do. Most lines of our code are numerical linear algebra.
- Be flop and memory aware! The form of a mathematical expression and the way the expression should be evaluated in actual practice may be quite different.
- Be alert to problem structure and make educated choice of software/algorithm — the structure should be exploited whenever solving a problem.
- Do not write your own matrix computation routines unless for good reason. Utilize BLAS and LAPACK as much as possible!
- In contrast, for optimization, often we need to devise problem specific optimization routines, or even “mix and match” them.

## 4 Ridge regression by SVD

In ridge regression, we minimize

$$\|y - X\beta\|_2^2 + \lambda\|\beta\|_2^2$$

If we obtain SVD of  $X$  such that  $X = U\Sigma V$ , then the equation is

$$(\Sigma^2 + \lambda I_p)V'\beta = \Sigma U'y.$$

We get

$$\hat{\beta}_\lambda = \sum_{j=1}^r \frac{\sigma_j u_j' y}{\sigma_j^2 + \lambda} v_j, \quad r = \text{rank}(X).$$

It is clear that  $\hat{\beta}_\lambda \rightarrow \beta_{OLS}$  as  $\lambda \rightarrow 0$  and  $\|\hat{\beta}_\lambda\|_2$  is monotone decreasing as  $\lambda \rightarrow \infty$ .