

Numerical Stability Analysis of Matrix Decomposition in Solving Linear Systems

Yicheng Zhao

Northeast Yucai Foreign Language School, Shenyang, China
zyc20071119@gmail.com

Abstract. This paper investigates the numerical stability of QR decomposition, Singular Value Decomposition (SVD), and Cholesky decomposition in least squares problems. Through theoretical analysis and numerical experiments, the computational errors and efficiency of the three decomposition methods are compared for matrices with different condition numbers. The experimental results show that SVD decomposition exhibits the best robustness for ill-conditioned matrices, while Cholesky decomposition is the most efficient for well-conditioned matrices. Additionally, this paper compares the performance of direct solving (without decomposition) with decomposition methods, demonstrating that decomposition methods significantly outperform direct solving in terms of numerical stability and computational efficiency. To further validate the findings, we conduct experiments on both synthetic and real-world datasets, covering a range of matrix sizes and condition numbers. The results highlight the trade-offs between accuracy and computational cost, providing practical insights for selecting the appropriate decomposition method based on specific problem requirements. This study not only reinforces the theoretical understanding of matrix decompositions but also offers actionable guidelines for their application in scientific computing and data analysis

Keywords: Matrix decomposition, Numerical stability, Condition number, QR decomposition

1. Introduction

The least squares method is a core approach for solving overdetermined linear systems widely used in data fitting, signal processing, and machine learning. However, directly solving the normal equations may lead to numerical instability, especially when the matrix condition number is large [1]

Matrix decomposition transforms the original matrix into a more manageable form, significantly improving computational efficiency and stability [2]. This paper focuses on QR decomposition, SVD decomposition, and Cholesky decomposition, analyzing their numerical properties in least squares problems and validating their applicability through experiments. Furthermore, this paper compares the performance of direct solving (without decomposition) with decomposition methods, highlighting the advantages of decomposition methods [3].

2. Theoretical background

2.1. Least squares problem (example)

The least squares solution x^* satisfies:

$$x^* = \operatorname{argmin}_x \|Ax - b\| \quad (1)$$

The analytical solution is

$$x = (A^T A)^{-1} A^T b \quad (2)$$

but direct computation may amplify errors due to the squared condition number of $A^T A$

2.2. Matrix decomposition methods

(1) QR Decomposition

Decompose A into an orthogonal matrix Q and an upper triangular matrix R :

$$A = QR \Rightarrow x^* = R^{-1} Q^T b \quad (3)$$

Features: Avoids computing $A^T A$, preserving the condition number $\kappa(A)$.

(2) SVD Decomposition

Decompose A into $U\Sigma V^T$, with the least squares solution:

$$x^* = V\Sigma^{-1}U^T b \quad (4)$$

Features: Optimal stability, capable of handling rank-deficient matrices.

(3) Cholesky Decomposition

When is A full-rank, decompose $A^T A = LL^T$, with the solution:

$$x^* = L^{-T} L^{-1} A^T b \quad (5)$$

Feature: High computational efficiency, but sensitive to ill-conditioned matrices, and can only be used in positive definite matrices

2.3. Condition number and numerical stability

(1) Condition Number

The matrix condition number

$$\kappa(A) = \|A\| \|A^{-1}\| \quad (6)$$

measures the sensitivity of the problem to the tiny perturbation of the initial values of the multiplication by matrices. For the least square problems specifically, the relatively error will satisfy the inequalities of

$$\frac{\|\delta x\|}{\|x\|} \leq \kappa(A) \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right) + o \quad (7)$$

As a result, the error of the solutions can be magnified exponentially if the condition number $\kappa(A)$ is significantly large [4].

For the Cholesky decomposition, direct calculation of $A^T A$ is required, so that the condition number is magnified squarely by

$$\kappa(A^T A) = \kappa(A)^2 \quad (8)$$

leading to poor stability for Cholesky decomposition.

For the QR decomposition, the original matrix is decomposed into the multiplication of the orthogonal matrix Q and an upper triangular matrix R . Given that the orthogonal matrix Q has the norm equal to 1 for itself and its inverse Q^{-1} , the condition number remains the same that [5].

$$\kappa(A) = \kappa(R) \quad (9)$$

For the SVD, the orthogonal matrices act as the same functions, while the algorithm for SVD also contains the truncation of the small singular values, it can reduce the condition number implicitly and avoid the ill-conditioning problems.

(2) Floating-point Arithmetic

Due to the methods of the floating-point arithmetic, the computation of matrices will propagate errors inevitably [6].

For the Cholesky decomposition, the prerequisite is that the matrix should be positive definite, so the floating-point mechanism may damage the positive-definiteness and the occurrence of having square root of the negative values. The square root and division of calculation will also lead to the round-off error by calculating in the process of Cholesky decomposition.

$$L_{ii} = \sqrt{A_{ii} - \sum_{k=1}^{i-1} L_{ik}^2} \quad (10)$$

For the QR decomposition, the Gram-Schmidt orthogonalization may lose orthogonality due to the round-off error by calculating

$$q_i = a_i - \sum_{j=1}^{i-1} (a_i^T q_j) q_j \quad (11)$$

Although improved methods such as Householder transformation and Givens rotation can maintain orthogonality, round-off errors still exist.

For SVD, the iteration algorithm will also produce round-off errors when calculating the singular values and singular vectors. Consequently, error propagation may lead to the loss of the matrix's orthogonality.

3. Experimental design and implementation

3.1. Experimental setup

In the numerical experiment, the matrix of a specific dimension will be generated randomly and be calculated by various decomposition approaches and the direct solving method. The final results will be represented by the errors between the accurate solution and the numerical solutions, and the computation time will also be considered as an index to exhibit the efficiency of the algorithm. Two

categories of the matrices will be generated with well-conditioning being random and full-rank, while the ill-conditioned ones are constructed by adding small perturbation to approximate singularity and high condition number $\kappa(A)$.

Given that the Cholesky decomposition can only be utilized in the scenario that the matrix is positive definite, two groups of matrices will be tested while the first group is the $A \in \mathbb{R}^{100 \times 50}$ and the other group contains the matrices with $A \in \mathbb{R}^{50 \times 50}$ and is positive definite.

Tools: Python + NumPy/SciPy.

3.2. Experimental results

Table 1. Results for random matrix(100*50)

Method	Well-conditioned error	Well-conditioned times	Ill-conditioned error	Ill-conditioned times
QR	1.03e-12	0.00027s	4.89e-06	0.0037s
SVD	8.39e-13	0.00084s	2.88e-08	0.0076s
Direct	2.57e-12	0.00030s	4.34e-01	0.0094s

Table 2. Results for positive-definite matrix (50*50)

Method	Well-conditioned error	Well-conditioned times	Ill-conditioned error	Ill-conditioned times
QR	3.12e-14	0.00097s	6.45e-07	0.0029s
SVD	2.98e-14	0.00182s	3.82e-09	0.0086s
Cholesky	1.05e-14	0.00062s	9.87e-03	0.0012s
Direct	1.33e-14	0.00046s	8.46e-01	0.0040s

3.3. Discussion

As shown in Table 1 and Table 2, from the above results, we can conclude that Cholesky decomposition has the smallest error and the highest efficiency when dealing with well - conditioned problems. However, it can only be applied to positive definite matrices under restricted conditions. Meanwhile, direct solving has a similar error but poorer stability.

For the ill-conditioned Matrix, SVD decomposition has the smallest error and the highest numerical stability, although with the longest calculating time, followed by QR decomposition which has moderate numerical stability. Cholesky decomposition and direct solving fail due to the instability in the circumstances of ill-conditioning problems.

4. Practical application case study

4.1. Structural engineering analysis: ill-conditioned stiffness matrix in bridge problems

(1) Background

In bridge structural mechanics analysis, the stiffness matrix K , which is usually positive definite, may become ill-conditioned ($\kappa(K) \gg 10^6$) due to material parameter errors or redundant degrees of freedom such as the over-constraint problems. Careful selection of numerical methods is required when solving the equation $Ku = F$ to reach the most precision results without considering the limit of the efficiency of the algorithm, where K represents the characteristics of the stiffness of

a specific engineering structure, u represents the displacement vector of nodes, and F is the load vector representing the force exerted on the nodes [7].

(2) Experimental Design

Generate Ill-conditioned Stiffness Matrix: Construct a matrix $K \in \mathbb{R}^{50 \times 50}$ with a condition number $\kappa(K) > 10^{12}$ by adding perturbation, Generate a load vector F with added Gaussian noise.

Solution Methods: Use Cholesky decomposition, QR decomposition, and SVD decomposition to solve for u

Evaluation Metrics: Relative error $\|Ku - F\| / \|F\|$.

Computation time.

(3) Results

Table 3. Results for engineering applications

Method	Relative Error	Computation Time
Cholesky	9.73e-01	0.00011s
QR	2.15e-06	0.00235s
SVD	1.04e-08	0.00731s
Direct	1.23e-00	0.00043s

As shown in Table 3, in engineering problems, efficiency can be negligible since instantaneity is not required, while precision is a key factor to consider. Ill - conditioned problems prevent the application of Cholesky Decomposition due to its numerical instability. As a result, SVD can be the most suitable matrix decomposition for this problem.

4.2. Finance modeling: well-conditioned covariance matrix in investment combination

(1) Background

In the optimization of investment combination, the covariance matrix Σ is positive definite and well-conditioned ($\kappa(K) < 10^3$) due to the feature of the data from the financial markets. The aim is to calculate the optimal combination of the investment with least risk or highest profits by minimize or maximize $\omega^T \Sigma \omega$, where Σ represents for the covariance of the assets and ω represents for the weight of the assets

(2) Experimental Design

Step1: Generate Well-conditioning Covariance Matrix:

Construct a matrix $\Sigma \in \mathbb{R}^{100 \times 100}$ by using the data from the history yield rate in 1000 days.

Step2:Solution Methods:

Use Cholesky decomposition, QR decomposition to get the answers.

Step3: Evaluation Metrics:

The difference of the objective function and the optimal solution theoretically.

Computation time.

(3) Results

Table 4. Results for financial applications

Method	Relative Error	Computation Time
Cholesky	1.22e-21	0.00012s
QR	4.50e-22	0.00161s

As shown in Table 4, the financial problems have a great demand in the efficiency of the algorithm in order to capture the opportunities to maximize the profits and avoid loss and need a flexible analysis of the trends of the market. Considering that the covariance matrices are usually well-conditioning and positive definite, it is an optimal choice to utilize the Cholesky Decomposition as an approach.

4.3. Signal processing: moderate condition number designing matrix in filter

(1) Background

In signal denoising, the design matrix A may form moderate ill-condition problems due to the redundant configuration of the sensors. The target is to compute the filter coefficient x requires both the numerical stability and the instantaneity. The equation containing $Ax = b$ where A is the design matrix and b is the observed signal containing the noise [8,9].

(2) Experimental Design

Step1: Generate Moderate-conditioning Matrix:

Construct a matrix $A \in \mathbb{R}^{100 \times 20}$ by using the data to simulate the signal received by the sensors

Step2: Solution Methods:

Use QR decomposition and SVD decomposition to get the answers.

Step3: Evaluation Metrics:

The SNR improvement by the calculation.

Computation time.

(3) Results

Table 5. Results for signal processing applications

Method	Relative Error	Computation Times
SVD	1.03e-01	0.00124s
QR	9.99e-02	0.00876s

As shown in Table 5, for signal processing problems, there are small differences in error, improvement, and computation times between the two approaches, and no numerical instability issues occur for them. As a result, both of the algorithm can be applied in solving the equations of signal processing.

5. Conclusion

In terms of numerical stability and computational efficiency, SVD decomposition performs best for ill - conditioned problems and has the highest precision in solving least - square problems. However, it has higher computational costs and is unable to obtain the answer instantaneously for high - dimensional matrices because of its high computational complexity. QR decomposition offers a balance between stability and efficiency. It can provide moderate numerical stability when facing ill

- conditioned problems and also has relatively higher computational efficiency. Cholesky decomposition is the fastest for well-conditioned matrices with low computational complexity but is sensitive to condition numbers and can only be restricted in solving the positive definite matrices.

As for the limitations of direct solving, Direct solving produces large errors for ill-conditioned matrices and is not recommended due to the lack of orthogonality that can reduce the condition number and easily to propagate significant errors in the multiplication and solving inverse.

This paper puts forward application recommendations, which can meet the requirements of high-precision : Prioritize SVD decomposition; For real-time computation: Use Cholesky decomposition for well-conditioned problems and QR decomposition for ill-conditioned problems.

References

- [1] Slock D T M .Backward consistency concept and round-off error propagation dynamics in recursive least-squares algorithms [J].Optical Engineering, 1992, 31(6): 1153-1169.
- [2] Björck, ke. Numerical Methods for Least Squares Problems [J].SIAM, 1996.
- [3] Demmel J .Applied Numerical Linear Algebra [M].Tsinghua University Press : Siam, 2011.
- [4] Watkins D S .Fundamentals of matrix computations [J].Wiley, 2019.
- [5] Gardiner R B J D .Fundamentals of Matrix Computations (David S. Watkins) [J].Siam Review, 2006, 35(3): 520-521.
- [6] Chen S C .Statistical Modeling, Matrix Computation and Stochastic Process. EXPERIENCES [J]. 2014. [1] Higham, N. J. (2002). -Accuracy and Stability of Numerical Algorithms (2nd Edition). SIAM.
- [7] Klaus-Jürgen Bathe, Ramm E , Wilson E L .Finite element formulations for large deformation dynamic analysis [J].International Journal for Numerical Methods in Engineering, 2010, 9(2): 353-386.
- [8] Noble J .Statistical Digital Signal Processing And Modeling [J]. 2016.
- [9] Narsimha B , Reddy P K A .Digital Signal Processing Lab Using MATLAB [J]. 2019.