

Streaming Power Iteration for Muon: 4. Principles

Su Jianlin

2026-04-13

After the three articles “Streaming Power Iteration for Muon: 1. First Glance”, “Streaming Power Iteration for Muon: 2. Acceleration”, and “Streaming Power Iteration for Muon: 3. Refinement”, readers should already have a good grasp of the idea, implementation, and acceleration details of Streaming Power Iteration. Overall, this can be regarded as a highly competitive method for implementing Muon, and thanks to its direct approximation of SVD computation, it also offers better scalability.

Due to space limitations, the mathematical principles behind the relevant operations were described relatively briefly at that time. Therefore, in this article, we will supplement some mathematical derivations regarding power iteration and QR decomposition, aiming to establish a more complete theoretical picture. However, these derivations still emphasize intuition over rigor, primarily to help clarify thinking (including the author’s), and we ask professional readers to be indulgent.

Coaxial Equivalence

Before beginning the derivation, we need to introduce the concept of “coaxial equivalence”. For matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}$, if there exists a signature matrix \mathbf{S} such that $\mathbf{A} = \mathbf{B}\mathbf{S}$, then \mathbf{A} and \mathbf{B} are said to be “coaxially equivalent”, and each is called a “coaxial matrix” of the other. Here, a “signature matrix” refers to a diagonal matrix whose diagonal entries are ± 1 , i.e., $\text{diag}(\pm 1, \pm 1, \dots, \pm 1)$.

It should be noted that the term “coaxial” is coined by the author to describe this equivalence relation, because from a coordinate system perspective, matrices \mathbf{A}, \mathbf{B} satisfying this condition actually describe the same coordinate system, differing only in the positive direction choices of certain axes. Some literature calls them “sign equivalent”, but the author finds “coaxial” more intuitive.

The reason for introducing the concept of coaxial equivalence is that many matrix decompositions are unique only up to coaxial equivalence (thus the author is puzzled why such a commonly used equivalence relation lacks a standard name). For example, in QR decomposition, suppose matrix $\mathbf{A} \in \mathbb{R}^{n \times m} (n \geq m)$ is full-rank, and $\mathbf{Q}_1\mathbf{R}_1$ and $\mathbf{Q}_2\mathbf{R}_2$ are two QR decompositions of \mathbf{A} ; then \mathbf{Q}_1 and \mathbf{Q}_2 are coaxially equivalent, and so are \mathbf{R}_1 and \mathbf{R}_2 .

This uniqueness is not hard to understand, because in the Gram-Schmidt orthogonalization process, we can arbitrarily reverse the direction of the orthogonalized vectors without changing orthogonality. Standard textbooks often ensure uniqueness of QR decomposition by requiring the diagonal entries of \mathbf{R} to be positive, which is another valid approach but can sometimes bring unnecessary complications. Additionally, the uniqueness of SVD also holds under coaxial equivalence.

Power Iteration

In “Streaming Power Iteration for Muon: 1. First Glance”, we introduced power iteration by sequentially solving for eigenvectors, then directly switched to the parallel version, asserting that it has the same convergence behavior. Here, we start directly from the parallel version and prove its convergence.

Let matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ ($n \geq m$) be full-rank (rank m), and consider the power iteration

$$\mathbf{V}_t = \text{QR}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}), \quad \mathbf{V}_0 = \mathbf{I} \quad (1)$$

Let the QR decomposition of $\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}$ be $\mathbf{Q}_t \mathbf{R}_t$, so that $\text{QR}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}) = \mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1} \mathbf{R}_t^{-1}$. Iterating step by step gives

$$\mathbf{V}_t = (\mathbf{A}^\top \mathbf{A})^t \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \dots \mathbf{R}_t^{-1} \quad (2)$$

Next, let the SVD of \mathbf{A} be $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$, so that $\mathbf{A}^\top \mathbf{A} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^\top$. Substituting into the above yields

$$\mathbf{V}_t = \mathbf{V} \mathbf{\Sigma}^{2t} \mathbf{V}^\top \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \dots \mathbf{R}_t^{-1} = \mathbf{V} \underbrace{(\mathbf{\Sigma}^{2t} \mathbf{V}^\top \mathbf{\Sigma}^{-2t})}_{\mathbf{B}} \underbrace{(\mathbf{\Sigma}^{2t} \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \dots \mathbf{R}_t^{-1})}_{\mathbf{C}} \quad (3)$$

The result is divided into three parts: $\mathbf{V}, \mathbf{B}, \mathbf{C}$. We aim to prove that $\mathbf{V}_t \rightarrow \mathbf{V}$ in the coaxial sense. It is known that both \mathbf{V}_t and \mathbf{V} are orthogonal matrices. By the closure property of upper triangular matrices, \mathbf{C} is upper triangular. If we can show that \mathbf{B} tends to an upper triangular matrix, then by the uniqueness of QR decomposition (up to coaxial equivalence), we obtain $\mathbf{V}_t \rightarrow \mathbf{V}$. For $\mathbf{B} = \mathbf{\Sigma}^{2t} \mathbf{V}^\top \mathbf{\Sigma}^{-2t}$, writing in components gives

$$B_{i,j} = V_{j,i} (\sigma_i / \sigma_j)^{2t} \quad (4)$$

Assuming $\sigma_1 > \sigma_2 > \dots > \sigma_m$, then for $i > j$, $(\sigma_i / \sigma_j)^{2t} \rightarrow 0$, meaning the sub-diagonal part of \mathbf{B} tends to zero; in other words, \mathbf{B} tends to an upper triangular matrix, satisfying the condition. This also implies that the convergence rate of power iteration depends on the ratio of adjacent singular values: the larger σ_i / σ_{i+1} , the faster the convergence. The iteration fails if there are equal singular values, but in practice, we can assume that the probability of two singular values being exactly equal is zero, thus avoiding this problem.

Essential Considerations

Let us carefully examine the entire proof process and understand what it essentially achieves.

First, the part $\mathbf{C} = \mathbf{\Sigma}^{2t} \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \dots \mathbf{R}_t^{-1}$ appears complex, but its only role is to be “an upper triangular matrix” that is, its precise form is irrelevant as long as it remains upper triangular. The core part driving the result is $\mathbf{B} = \mathbf{\Sigma}^{2t} \mathbf{V}^\top \mathbf{\Sigma}^{-2t}$, which tends to an upper triangular matrix, allowing the leading \mathbf{V} to be extracted via QR decomposition.

The step $(\mathbf{A}^\top \mathbf{A})^t$ is exactly what achieves this effect! In theory, $\mathbf{V}_t = \text{QR}((\mathbf{A}^\top \mathbf{A})^t)$, meaning the repeated QR operations beforehand are theoretically redundant—the orthogonalization could be done just once at the end. Of course, this equivalence is only of theoretical value; directly computing $(\mathbf{A}^\top \mathbf{A})^t$ would cause numerical overflow or underflow, making results unusable. Thus, performing QR periodically instead of only once at the end essentially maintains numerical stability.

We can also understand the multiple QR operations from the perspective of “preconditioning”. A preconditioner refers to pre-processing steps on inputs that theoretically do not change the output but are usually beneficial for numerical computation. For QR, right-multiplying any full-rank upper

triangular matrix \mathbf{R} does not alter the result, i.e., $\text{QR}(\mathbf{A}) = \text{QR}(\mathbf{A}\mathbf{R})$. Therefore, any such \mathbf{R} can be called a preconditioner for QR.

We know that QR itself can be written as right-multiplication by an upper triangular matrix. Thus, QR is a preconditioner for itself. Hence, we can insert QR operations into $(\mathbf{A}^\top \mathbf{A})^t$, transforming it into

$$(\mathbf{A}^\top \mathbf{A})^t \quad \rightarrow \quad \mathbf{A}^\top \mathbf{A} \text{QR}(\mathbf{A}^\top \mathbf{A} \text{QR}(\mathbf{A}^\top \mathbf{A} \cdots \text{QR}(\mathbf{A}^\top \mathbf{A}))) \quad (5)$$

Performing QR at both ends does not change the final result. However, since each QR produces an orthogonal matrix, right-multiplying by an orthogonal matrix carries almost no risk of numerical explosion, making QR a good self-preconditioner.

Related Variants

From the preconditioning viewpoint, we can understand many variants of power iteration. For example, the ColNorm version used at the beginning of the first article “Streaming Power Iteration for Muon: 1. First Glance”:

$$\mathbf{V}_t = \text{QR}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}) \quad \rightarrow \quad \mathbf{V}_t = \text{QR}(\mathbf{A}^\top \text{ColNorm}(\mathbf{A} \mathbf{V}_{t-1})) \quad (6)$$

The two \mathbf{V}_t are equal because column normalization ColNorm of a matrix is equivalent to right-multiplication by a diagonal matrix:

$$\text{ColNorm}([\mathbf{a}_1, \mathbf{a}_2, \cdots, \mathbf{a}_m]) = [\mathbf{a}_1, \mathbf{a}_2, \cdots, \mathbf{a}_m] \begin{bmatrix} \|\mathbf{a}_1\|^{-1} & 0 & \cdots & 0 \\ 0 & \|\mathbf{a}_2\|^{-1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \|\mathbf{a}_m\|^{-1} \end{bmatrix} \quad (7)$$

Diagonal matrices are a special case of upper triangular matrices, so ColNorm is also a preconditioner for QR and does not change the result. For the same reason, we can shift ColNorm one step outward, which also preserves the QR result:

$$\mathbf{V}_t = \text{QR}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}) \quad \rightarrow \quad \mathbf{V}_t = \text{QR}(\text{ColNorm}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1})) \quad (8)$$

Additionally, the protagonists of “Streaming Power Iteration for Muon: 2. Acceleration” and “Streaming Power Iteration for Muon: 3. Refinement” are the double QR power iteration:

$$\mathbf{V}_t = \text{QR}(\mathbf{A}^\top \mathbf{A} \mathbf{V}_{t-1}) \quad \rightarrow \quad \mathbf{V}_t = \text{QR}(\mathbf{A}^\top \text{QR}(\mathbf{A} \mathbf{V}_{t-1})) \quad (9)$$

From the discussions in those sections, it is easy to show that adding an extra QR to $\mathbf{A} \mathbf{V}_{t-1}$ does not theoretically change the result for \mathbf{V}_t .

Finite Precision Errors

We have repeatedly used phrases like “in theory”, because many results rely on infinite-precision, exact QR decompositions. In practice, QR is computed with finite precision, and for efficiency we often use the less accurate “Shifted Cholesky QR (SCQR)”, so analyzing the accumulation of errors becomes necessary.

By now, readers should be familiar with SCQR from the previous articles. It decomposes the QR of matrix \mathbf{A} into two steps: 1) perform Cholesky decomposition on $\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I}$ to obtain \mathbf{R} ,

2) compute the orthogonal matrix $\mathbf{Q} = \mathbf{A}\mathbf{R}^{-1}$. When $\lambda = 0$, SCQR is equivalent to exact QR. However, $\lambda = 0$ almost always fails in practice due to ill-conditioning. Setting an appropriate $\lambda > 0$ introduces error.

Fortunately, this error does not accumulate! Both $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I}$ yield upper triangular matrices upon Cholesky decomposition. Then \mathbf{A} is right-multiplied by the inverse of such a triangular matrix. As we know, right-multiplication by an upper triangular matrix does not alter the QR result. Thus, ignoring errors from matrix multiplication, we can consider SCQR to be “lossless” with respect to QR.

In other words, as long as an exact QR is used in the final step, no matter how many inaccurate SCQR steps are performed earlier, the final result remains exact. Conversely, the total error is equivalent to the error of a single SCQR step (the last one), and does not accumulate. This principle underlies also the “SCQR2” acceleration technique mentioned in “Streaming Power Iteration for Muon: 2. Acceleration”.

Conversely, if a transformation cannot be written as “right-multiplication by an upper triangular matrix”, it will lossily destroy information of the original matrix. Over long iterations, such errors accumulate until the process fails completely. The simplification scheme proposed by @Ji_Ha_Kim in the previous article “Streaming Power Iteration for Muon: 3. Refinement” is a classic counterexample.

Cholesky Decomposition

Finally, let us briefly review “Cholesky decomposition”. Its goal is to factor a given symmetric positive definite matrix \mathbf{B} into $\mathbf{L}\mathbf{L}^\top$, where \mathbf{L} is a lower triangular matrix. Denoting $\mathbf{R} = \mathbf{L}^\top$, we obtain the upper triangular form $\mathbf{R}^\top \mathbf{R}$.

Cholesky decomposition is fast because it admits a direct recursive analytic solution. Specifically, equating the entries from

$$\begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,m} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,m} \end{bmatrix} = \begin{bmatrix} l_{1,1} & & & \\ l_{2,1} & l_{2,2} & & \\ \vdots & \vdots & \ddots & \\ l_{m,1} & l_{m,2} & \cdots & l_{m,m} \end{bmatrix} \begin{bmatrix} l_{1,1} & l_{2,1} & \cdots & l_{m,1} \\ & l_{2,2} & \cdots & l_{m,2} \\ & & \ddots & \vdots \\ & & & l_{m,m} \end{bmatrix} \quad (10)$$

and solving entry by entry gives the recursive formulas

$$l_{j,j} = \pm \sqrt{a_{j,j} - \sum_{k=1}^{j-1} l_{j,k}^2}, \quad l_{i,j} = \frac{1}{l_{j,j}} \left(a_{i,j} - \sum_{k=1}^{j-1} l_{i,k} l_{j,k} \right) \quad \text{for } i > j \quad (11)$$

Cholesky QR is a classic application of Cholesky decomposition. Like Gram-Schmidt orthogonalization, it allows us to skip the orthogonalization procedure and directly obtain the triangular matrix \mathbf{R} . However, both face the same numerical difficulties: Gram-Schmidt proceeds vector by vector, and fails if linearly dependent or zero vectors are encountered, which manifests in singular values as low effective rank or large condition number; Cholesky decomposition similarly tends to fail under such conditions.

Another classic application of Cholesky decomposition is inversion of symmetric positive definite matrices. For example, solving $\mathbf{B}\mathbf{X} = \mathbf{C}$, where \mathbf{B} is symmetric positive definite, reduces to solving two triangular systems once we decompose \mathbf{B} as $\mathbf{L}\mathbf{L}^\top$. This process is efficient at each step. In particular, the iterative method proposed by @Ji_Ha_Kim for computing msign is based on inverting matrices via this approach.

Summary

This article mainly supplements some mathematical derivations regarding streaming power iteration, including convergence of power iteration, QR preconditioning, and a brief introduction to Cholesky decomposition, aiming to help readers better understand this method from a theoretical standpoint.

Please include this URL when reposting: <https://kexue.fm/archives/11710>

More details on republishing: [The Science Space FAQ](#)