

Notes on numerical stabilization in determinant quantum Monte Carlo

Carsten Bauer¹

¹*Institute for Theoretical Physics, University of Cologne, 50937 Cologne, Germany*

(Dated: May 17, 2019)

In these notes we will empirically compare different matrix decompositions for use in numerical stabilization in many-body quantum Monte Carlo simulation, in particular when calculating equal time Green's functions. We will both benchmark their speed and accuracy. We will further show that special care has to be taken when calculating time-displaced Green's functions and will review a well-known but somewhat hard to find method for stabilizing the necessary inversion in this case. We will focus on and use the Julia programming language for all calculations. An implementation of all discussed methods is available in the open-source software library StableDQMC.jl.

I. THE ISSUE

A. Determinant quantum Monte Carlo in a nutshell

For the description of the DQMC algorithm¹, we consider a quantum field theory that can be split into a purely bosonic part S_B and a part S_F , which comprises fermion kinetics T and boson-fermion interactions V . An example is the famous Hubbard model after decoupling the on-site interaction $Un_{i,\uparrow}n_{i,\downarrow}$ by means of a Hubbard-Stratonovich or Hirsch transformation in either the spin or charge channel². As per usual, the central quantity of interest is the partition function

$$\mathcal{Z} = \int D(\psi, \psi^\dagger, \phi) e^{-S_B - S_F}. \quad (1)$$

The basic idea is to switch from the d dimensional quantum theory to a $D = d + 1$ dimensional classical theory as in the path integral framework. The first step is to discretize, $\beta = M\Delta\tau$, and Trotter decompose the imaginary time τ ,

$$\mathcal{Z} = \int D\phi e^{-S_B} \text{Tr} \left[\exp \left(-\Delta\tau \sum_{l=1}^M \psi^\dagger [T + V_\phi] \psi \right) \right]. \quad (2)$$

The exponential is then separated and the appearing commutator ignored which results in a systematic error of the order $\mathcal{O}(\Delta\tau^2)$,

$$\begin{aligned} e^{A+B} &\approx e^A e^B \\ e^{-\Delta\tau(T+V)} &\approx e^{-\frac{\Delta\tau}{2}T} e^{-\Delta\tau V} e^{-\frac{\Delta\tau}{2}T} + \mathcal{O}(\Delta\tau^3), \\ \mathcal{Z} &= \int D\phi e^{-S_B} \text{Tr} \left[\prod_{l=1}^m B_l \right] + \mathcal{O}(\Delta\tau^2). \end{aligned} \quad (3)$$

Note that the potential part $e^{-\Delta\tau\psi^\dagger V_\phi \psi}$ in the time-slice propagators $B_l = e^{-\frac{\Delta\tau}{2}\psi^\dagger T \psi} e^{-\Delta\tau\psi^\dagger V_\phi \psi} e^{-\frac{\Delta\tau}{2}\psi^\dagger T \psi}$ depends on the boson ϕ due to the present fermion-boson coupling. Rewriting the trace in (3) as a determinant, an identity which can be proven [CITE](#), yields the fundamen-

tal form

$$\mathcal{Z} = \int D\phi e^{-S_B} \det G_\phi^{-1} + \mathcal{O}(\Delta\tau^2), \quad (4)$$

where

$$G = (\mathbb{1} + B_M B_{M-1} \cdots B_1)^{-1} \quad (5)$$

is the equal-time Green's function.

Only under specific circumstances, such as the presence of a symmetry, can the integral kernel be safely be interpreted as a probability weight, as G_ϕ and its determinant are generally complex valued (*sign problem*).

B. Numerical instabilities - Model system

We consider the following simple non-interacting model system in one dimension,

$$H = -t \sum_{\langle i,j \rangle} c_i^\dagger c_j + \mu \sum_i n_i, \quad (6)$$

where we set the hopping amplitude to $t = -1$ and the chemical potential to $\mu = -0.1$. Generally, we are interested in the equal time Green's function G , Eq. (5), as an observable and its determinant, which enters into the probability weight in Eq. (4). Let's make the latter point a bit more precise. In a Metropolis scheme it is actually the ratio of determinants,

$$\frac{\det G'^{-1}}{\det G^{-1}} = \frac{\det G}{\det G'}, \quad (7)$$

which determines the acceptance or rejection of a proposed (primed G) Markov walker step.

We will showcase the issue of numerical instabilities arising in the computation of G by discussing the accuracy of the calculation of the building block of G , the slice matrix product chain

$$B_M B_{M-1} \cdots B_1 = \underbrace{B B \cdots B}_{M \text{ factors}}. \quad (8)$$

Here the second equality stems from the fact that our model, Eq. 6, is non-interacting and the slice matrices B_l

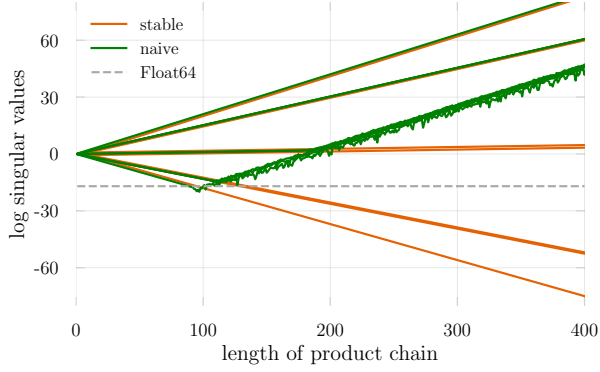


FIG. 1. **Numerical instabilities** (green) due to finite machine precision (**Float64**) in the calculation of the slice matrix product chain $B_M B_{M-1} \cdots B_1$ for model (6).

do neither depend on any Hubbard-Stratonovich boson field nor on imaginary time.

In Fig. 1 we see that a naive computation of Eq. 8 is doomed to fail for $M \geq M_c \approx 100$. Leaving a discussion of how to stabilize the product for the next section, let us understand the origin of this instability. The eigenvalues of the system are given by

$$\epsilon_k = -2t \cos(k) + \mu. \quad (9)$$

Neglecting the contribution by the chemical potential for simplicity, the energy values are bounded by $-2t \leq \epsilon_k \leq 2t$. The positive definite slice matrix $B = e^{-\Delta\tau T}$ has a condition number of about $\kappa \approx e^{4|t|\Delta\tau}$ and the product chain, Eq. 8, has $\kappa \approx e^{4|t|M\Delta\tau} = e^{4|t|\beta}$. We therefore see that the condition number diverges at low temperatures $T = 1/\beta = 1/(M\Delta\tau)$ where M is large. In this case roundoff errors due to finite machine precision will spoil the result. We can estimate the breakdown of the computation for the data type **Float64**, that is double floating-point precision³, by solving $\kappa(M) \sim 10^{-17}$ for M_c . We find $M_c \approx 100$ in good agreement with what we observe in Fig. 1. Switching to the non-IEEE data type **Double64**, we see in Fig. 2 that the onset of roundoff errors is shifted to lower temperatures.

II. STABILIZATION SCHEMES

One way to solve all problems is to switch to arbitrary precision numerics. In Julia the latter is provided through the **BigFloat** data type. However, this comes at the expense of (very) slow performance since, among other things, **BigFloat** lacks hardware support. Arbitrary precision numerics is still a valuable tool and we will use it to benchmark the accuracy of stabilization methods below.

How can we remedy the numerical issue outlined above in a floating point precision computation? The idea is to keep different scales separated throughout the computation as long as possible and only mix them in the final

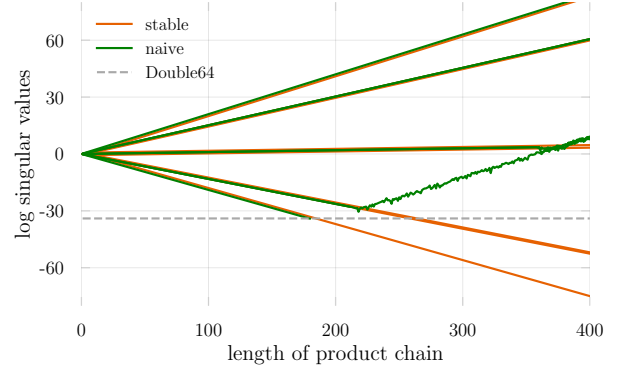


FIG. 2. **Numerical instabilities** due to finite machine precision (**Double64**) in the calculation of the slice matrix product chain $B_M B_{M-1} \cdots B_1$ for model (6).

step if necessary. A way to realize this is by deploying matrix decompositions, which we write generally as

$$B = UDX. \quad (10)$$

Here, U and X are matrices of scale ~ 1 and D is a real diagonal matrix containing all the different scales of B on the diagonal. We will refer to the values in D as singular values independent of the particular decomposition.

Instead of calculating the matrix product $B_2 B_1$, we compute

$$\begin{aligned} B_2 B_1 &= \underbrace{U_2 D_2 X_2}_{B_2} \underbrace{U_1 D_1 X_1}_{B_1} \\ &= U_2 \underbrace{(D_2 ((X_2 U_1) D_1))}_{U' D' X'} X_1 \\ &= U_r D_r X_r, \end{aligned} \quad (11)$$

with $U_r = U_2 U'$, $D_r = D'$, and $X_r = X' X_1$. Here, brackets indicate the order of operations. If we follow this scheme, largely different scales contained in the diagonal matrices won't be mixed throughout the computation. Repeating this procedure for the full slice matrix product chain, Eq. (8), we obtain an accurate UDX decomposition of the result.⁴

Using this result, we have to be careful to separate scales in the equal-time Green's function calculation, Eq. 5, as well. A straightforward procedure to add the unit matrix and perform the inversion in a stabilized manner is given by

$$\begin{aligned} [\mathbb{1} + UDX]^{-1} &= [U \underbrace{(U^\dagger X^{-1} + D)}_{udx} X]^{-1} \\ &= [(Uu)d(xX)]^{-1} \\ &= (xX)^{-1} d(Uu)^{-1} \\ &= U_r D_r X_r, \end{aligned} \quad (12)$$

with $U_r = (xX)^{-1}$, $D_r = d$, $X_r = (Uu)^{-1}$.

Another prescription^{5,6}, trying to take even more care of scale mixing, is

$$\begin{aligned} [1 + UDX]^{-1} &= [(D_p^{-1}X^{-1} + UD_m)D_pX]^{-1} \\ &= X^{-1} [D_p^{-1} \underbrace{(D_p^{-1}X^{-1} + UD_m)^{-1}}_{\substack{udx \\ udx}}] \quad (13) \\ &= U_r D_r X_r, \end{aligned}$$

with $U_r = X^{-1}u$, $D_r = d$, and $X_r = x$. Here, we initially separate the scales in D as $D_p = \max(D, 1)$ and $D_m = \min(D, 1)$ and perform two intermediate decompositions to stabilize the inversion and recombination. We will see below that it is sometimes necessary to rely on this second procedure to obtain accurate results for G .

So far we haven't specified the concrete decomposition $B = UDX$. In fact, there are a couple of choices. We will consider presumably most popular ones below.

A. SVD (UDV^\dagger)

A SVD is given by

$$B = UDV^\dagger, \quad (14)$$

where U is unitary, D is a real diagonal matrix, and V^\dagger is unitary. In this case we can use the unitarity of U and V^\dagger to calculate inverse terms like, for example, $(Uu)^{-1}$ in the last line of 12 as $(Uu)^{-1} = u^\dagger U^\dagger$, which is generally cheaper.

Julia offers a couple of purely-Julia SVD implementations, like `GenericSVD.jl`, which we will use for `BigFloat` computations. However, some of the most optimized algorithms are part of LAPACK⁷ and Julia defaults to those algorithms for regular floating point types. Concretely, there are three SVD functions⁸ implementing different algorithms:

- `gesdd` (default): Divide-and-conquer (D&C)
- `gesvd`: Simple
- `gesvj`: Jacobi algorithm (through `JacobiSVD.jl`)

We will compare all of them below.

B. QR (UDT)

A QR decomposition reads

$$B = QR = UDT, \quad (15)$$

where we have split R into a diagonal, D , and upper triangular piece T . Hence, $U = Q$ is unitary, $D = \text{diag}(R)$ is a real diagonal matrix, and T is upper triangular.

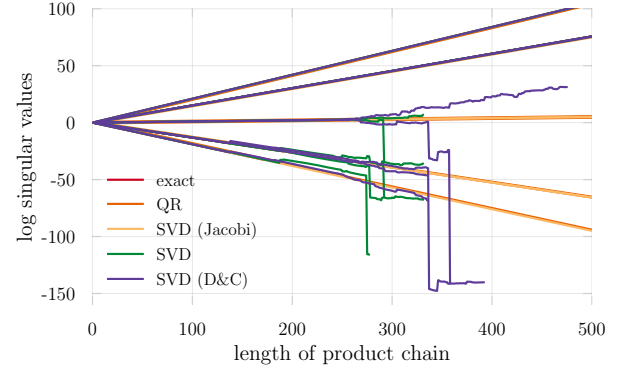


FIG. 3. **Comparison of matrix decompositions** to heal the numerical instabilities in the calculation of the slice matrix product chain $B_M B_{M-1} \cdots B_1$ for model (6). The QR and Jacobi SVD singular values seem to lie on top of the exact ones whereas regular SVD and divide-and-conquer SVD show large deviations at low temperatures $\beta \gtrsim 25$ ($\Delta\tau = 0.1$).

III. BENCHMARKS

In the following we want to assess how the mentioned matrix decompositions perform in stable computations of the slice matrix product chain $B_M \cdots B_1$, the Green's function G , and its determinant $\det G$, both with respect to accuracy and speed, for our model system in Eq. 6.

A. Accuracy

Before benchmarking the efficiency of an algorithm, it is crucial to check its correctness first. Fig. 3 shows the log singular values of the slice matrix product chain Eq. 8 stabilized with different matrix decompositions as a function of its length M . While QR and Jacobi SVD seem to lie on top of the numerically exact result, we observe large deviations for the simple and D&C SVD algorithms at low temperatures ($\beta \gtrsim 25$).

LAPACK SVD error bounds⁹ 'Thus large singular values (those near σ_1) are computed to high relative accuracy and small ones may not be.'

Turning to the equal-time Green's function, Eq. (5), we take the results for the slice matrix chains and perform the inversions according to the schemes presented above. We take the maximum absolute difference between the obtained Green's functions and the exact G as an accuracy measure. The findings for the simple inversion scheme in Eq. 12 are shown in Fig. 4. At high temperatures, all decompositions give the correct Green's function up to some limit close to floating point precision. However, at low temperatures only the QR decomposition reproduces G_{exact} reliably. It has the highest accuracy by a large margin, followed by the Jacobi SVD as the best of the SVD methods, which all fail to reproduce the exact result accurately. As displayed in Fig. 5, switching

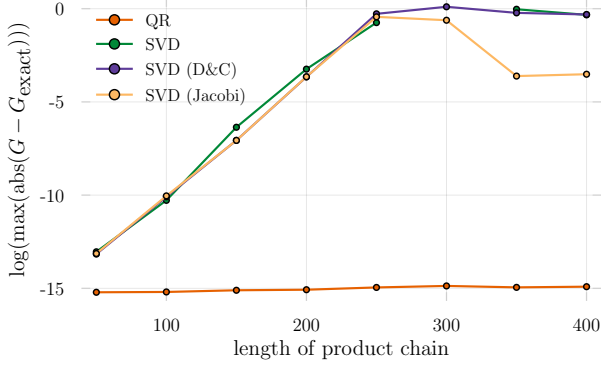


FIG. 4. **Accuracy of the Green's function** obtained from stabilized computations using the listed matrix decompositions and the inversion scheme from Eq. (12).

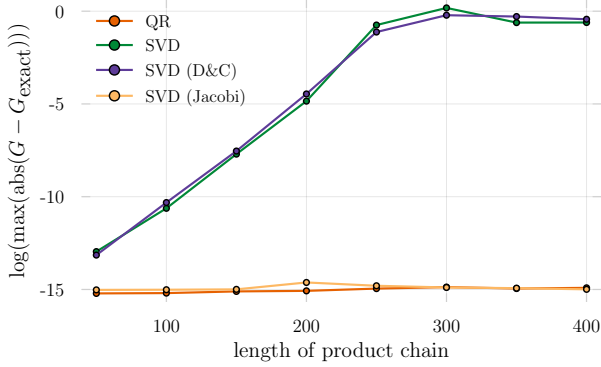


FIG. 5. **Accuracy of the Green's function** obtained from stabilized computations using the listed matrix decompositions and the careful inversion scheme, Eq. **TODO Scalettar**.

to the more careful procedure, Eq. 13, does improve the accuracy of the Jacobi SVD dramatically while the deviations seen for the other two SVD based schemes are still of order unity. This suggests that only the QR decomposition, independent of the inversion procedure, and the Jacobi SVD in combination with Eq. 13 are suited for computing the equal time Green's function reliably.

In Figs. 6, 7 we show the logarithm of the relative error of the Green's function determinant, relevant in the Metropolis acceptance, obtained for all combinations of matrix decompositions and inversion schemes. Similar to the results of the last paragraph, only the QR decomposition leads to accurate results for all accessed temperatures, irrespective the employed inversion scheme, whereas all SVD methods show large relative errors. Interestingly, using the Jacobi SVD and Eq. 13 for the inversion doesn't remedy the numerical issues for the Green's function determinant. We also note that the latter inversion scheme does generally lead to more imprecise results than the simpler one in Eq. 12.

B. Efficiency

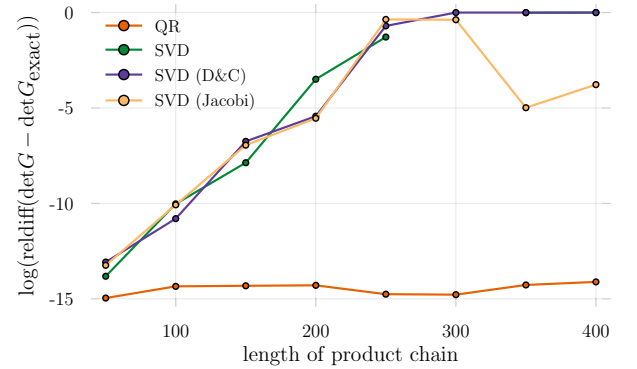


FIG. 6. **Accuracy of the determinant** of the equal-time Green's function obtained from stabilized computations using the listed matrix decompositions and the inversion scheme from Eq. (12).

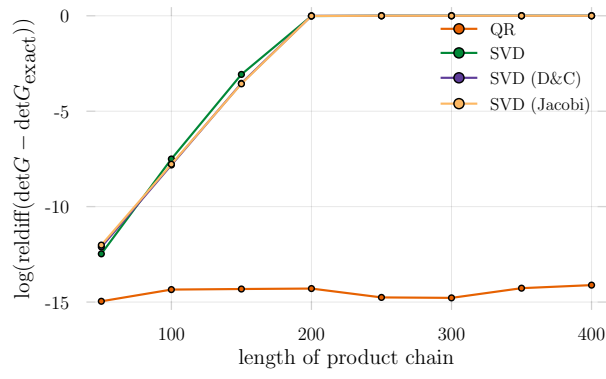


FIG. 7. **Accuracy of the determinant** of the equal-time Green's function obtained from stabilized computations using the listed matrix decompositions and the careful inversion scheme, Eq. [Scalettar](#).

IV. THE TIME-DISPLACED GREEN'S FUNCTION

Scalettar and other's method⁵ Hirsch method¹⁰

-
- ¹ R. Blankenbecler, D. J. Scalapino, and R. L. Sugar, “Monte Carlo calculations of coupled boson-fermion systems. I,” *Physical Review D* **24**, 2278–2286 (1981).
- ² J. E. Hirsch, “Discrete Hubbard-Stratonovich transformation for fermion lattice models,” *Physical Review B* **28**, 4059–4061 (1983).
- ³ David Goldberg, “What every computer scientist should know about floating-point arithmetic,” *ACM Comput. Surv.* **23**, 5–48 (1991).
- ⁴ Note that we intentionally do not discuss the option to calculate B^M as UD^MX since typically the system will be interacting and the B matrices in the product chain will differ.
- ⁵ E. Y. Loh, J. E. Gubernatis, R. T. Scalettar, S. R. White, D. J. Scalapino, and R. L. Sugar, “Numerical Stability and the Sign Problem in the Determinant Quantum Monte Carlo Method,” *International Journal of Modern Physics C* **16**, 1319–1327 (2005).
- ⁶ E. Y. Loh, J. E. Gubernatis, R. T. Scalettar, R. L. Sugar, and S. R. White, “Stable Matrix-Multiplication Algorithms for Low-Temperature Numerical Simulations of Fermions,” (1989) pp. 55–60.
- ⁷ E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen, *LAPACK Users’ Guide*, 3rd ed. (Society for Industrial and Applied Mathematics, Philadelphia, PA, 1999).
- ⁸ Note that Fortran LAPACK functions are named according to realness and symmetries of the matrix. In Julia multiple-dispatch takes care of routing different matrix types to different *methods*. The Julia function `gesdd` works for both real and complex matrices, i.e. there is no (need for) `cgesdd`.
- ⁹ Susan Blackford, “Error Bounds for the Singular Value Decomposition,” <http://www.netlib.org/lapack/lug/node96.html> (1999), [Online; accessed 16-May-2019].
- ¹⁰ J. E. Hirsch, “Stable Monte Carlo algorithm for fermion lattice systems at low temperatures,” **38**, 12023 (1988).