

TREX e-school with TurboRVB

Lecturer: Sandro Sorella

Targeting Real Chemical Accuracy at the Exascale project has received funding from the European Union Horizon 2020 research and innovation programme under Grant Agreement **No. 952165.**



Lecture III : QMC Wave function Optimization

Reminding Variational quantum Monte Carlo
Recent progress in stochastic optimization
Correlated sampling
Direct evaluation of energy derivatives
Stochastic reconfiguration or Natural gradients
Examples



- × 1D. The Jastrow- Fermi gas is the exact ground state for the Tomonaga-Luttinger model (B. Tayo and S.Sorella Phys. Rev.B, 2008)
- × D>1 Spin-Wave theory ground state is a Jastrow |MF=Neel>
- × Good description of Mott insulators (M. Capello et al. PRL 2005)
- × Spin liquid, the ground state of the Kitaev model can be written as a Gutzwiller projected pairing function with triplet correlations, again J(Gutzwiller-projection)-MeanField=Pfaffian function discussed before.
- Superconductivity in strongly correlated models (e.g. Hubbard) without the electronphonon mechanism.
- × Laughlin's WF FQHE is a (complex) Jastrow- MeanField (Jain PRB' 90)
- × Recently Valence bond Solids (Kekule') within a single determinant ansatz, taking into account several configurations.



Variational Monte Carlo

$$\pi(x) = \frac{|\Psi(x)|^2}{\int dx' |\Psi(x')|^2} \longrightarrow \{x_i\}_{\Psi} \quad x = \{\mathbf{r}_1 \sigma_1, \mathbf{r}_2 \sigma_2, \dots, \mathbf{r}_N \sigma_N\}$$

• Sampling the electron moves with Markov chains

$$\begin{array}{l} x \to x' \\ P(x \to x') = \min\left(1, \frac{|\langle x'|\Psi\rangle|^2}{|\langle x|\Psi\rangle|^2}\right) \\ \text{I principle} & \frac{\langle \Psi|H|\Psi\rangle}{\langle \Psi|\Psi\rangle} = \int dx \pi(x) e_L(x) \approx \frac{1}{M} \sum_{\text{M samples}} e_L(x_{\text{sample}}) \\ e_L(x) = \frac{\langle x|H|\Psi\rangle}{\langle x|\Psi\rangle} \quad \epsilon \text{ The local energy} \end{array}$$

• The variational principle



You believe in human learning and define a wave function $|\Psi_{\alpha}\rangle = J|\Psi_{\rm MF}\rangle$ where " α " indicates a set of variational parameters and J is a correlation term. For instance J = Gutzwiller projection, Jastrow correlation, three-body

The expectation value of the Hamiltonian, due to J, can be evaluated statistically:

$$E_{VMC}(lpha) = rac{\langle \Psi_lpha | H | \Psi_lpha
angle}{\langle \Psi_lpha | \Psi_lpha
angle}$$
 by QMC sampling

In the last decade ML and VMC were able to handle several parameters even when the target energy is known only statistically.



He₄ was studied with 2-parameters : $J = \prod_{i < j} \exp(-a_1/r_{ij}^{a_2})$ But also "recently" (M. Ogata et al. PRL 85 (2000) J= 1 :



The standard: parameters were optimized with several independent energy runs:

- .) at most two variational parameters
- 2) large computational resources for a single optimization



Many parameters \rightarrow Mean-field self consistency

e.g. Slater determinant in a lattice contains L (number of sites) x N (number of electrons) variational parameters. No problem with Hartree, Hartree-Fock, LDA, CGA, HSE...

In a many-body wf, self-consistent approaches are not available (what are the sc-fields?) and → Many parameters optimization problem → Moreover the energy evaluation is noisy



Correlated sampling, Umrigar 1988

We want to compute the energy $E(\alpha') = \frac{\langle \Psi_{\alpha'} | H | \Psi_{\alpha'} \rangle}{\langle \Psi_{\alpha'} | \Psi_{\alpha'} \rangle}$

Suppose to employ a Markov chain with the original weight

$$W_{\alpha}(x) = |\Psi_{\alpha}(x)|^2 \neq |\Psi_{\alpha'}(x)|^2, \text{ i.e. } \pi(x) = \frac{W_{\alpha}(x)}{\sum_{x'} W_{\alpha}(x')}$$

Then we can compute:

e:
$$R(x) = \frac{W_{\alpha'}(x)}{W_{\alpha}(x)} = \frac{|\Psi_{\alpha'}(x)|^2}{|\Psi_{\alpha}(x)|^2}$$

$$\sum e_L^{\alpha'}(x_{\text{sample}})R(x_{\text{sample}})$$

$$E_{\alpha'} = \frac{\sum_{x} \pi(x) R(x) e_L^{\alpha'}(x)}{\sum_{x} R(x) \pi(x)} = \frac{M \text{ samples}}{\sum_{x} R(x) \pi(x)} = \frac{M \text{ samples}}{M \text{ samples}}$$

 $R(x_{\text{sample}})$



Essentially variance R(x) ~ exp (# electrons)

Some improvement was to consider the minimization of the variance but the problem is only alleviated not solved.

Thus the approach is limited to very few electrons N<= 10



$$E(\alpha) = \frac{\int dx^{3N} \Psi_{\alpha}(x) H \Psi_{\alpha}(x)}{\int dx^{3N} \Psi_{\alpha}^{2}(x)}$$

$$\partial_{\alpha} E(\alpha) = 2 \operatorname{cov}(e_L(x), O(x))$$

$$\pi_{(x)} \propto \Psi_{\alpha}^2(x)$$

$$O(x) = \partial_{\alpha} \ln |\Psi_{\alpha}(x)|$$

$$e_L(x) = \frac{[H\Psi_{\alpha}](x)}{\Psi_{\alpha}(x)}$$



Next progress...

Next attempt ~ 2000: try to compute derivative and apply steepest descent:

$$\alpha'_k - \alpha_k = \delta \alpha_k = -\Delta \frac{\partial E}{\partial \alpha_k} = \Delta f_k$$

For small enough Δ the energy should go down. So one should optimize at each step by determining the optimal Δ :

$$\mathrm{Min}_{\Delta} E(\vec{\alpha} + \Delta \vec{f})$$

Go (explain SD)
$$\rightarrow$$



H2: simplest Jastrow+simplest geminal

Basis
$$:\phi_{1,a}(\vec{r}) = \exp(-Z|\vec{r} - \vec{R}_a|^2) \ a = 1, 2$$

 $g(\vec{r}, \vec{r}') = \sum_{ia,jb} \lambda_{ia,jb} \phi_{ia}(\vec{r}) \phi_{jb}(\vec{r}')$

 $\lambda_{11,11} = \lambda_{12,12} = 1$ $\lambda_{11,12} = \lambda_{12,11} = \lambda$

By symmetry only two parameters λ & Z in g



In practice one chooses Δ (tpar in TurboRVB) once for all



If you want to try steepest descent in TurboRVB: kl=-2 You can also use ADAMS (from ML) (yes_adams=.true.)



$$\operatorname{Cost}(\delta\alpha) = E(\vec{\alpha}) + \sum_{k} \left(\delta\alpha_{k} f_{k} + \frac{1}{2\Delta} \delta\alpha_{k}^{2} \right)$$

We implicitly optimize a Cost function at each step that penalizes big changes of variational parameters.

Is this the optimal thing to do?



There are directions that are extremely difficult to optimize, usually the most physical ones...

H₂ molecule (see tutorial) R=0.8 a.u. 4 parameters optimization







For each variational parameter α_k of a vector $\vec{lpha} = \{ lpha_1, lpha_2, \cdots lpha_p \}$

We can define an operator O_k :



But we are interested to normalized wf

$$|\psi_{\vec{\alpha}}\rangle = \frac{|\Psi_{\vec{\alpha}}\rangle}{||\Psi_{\vec{\alpha}}||}$$

Thus:

$$\begin{aligned} |\psi_{\vec{\alpha}+\delta\vec{\alpha}}\rangle &= \left(1 + \sum_{k=1}^{p} \delta\alpha_{k} (O_{k} - \bar{O}_{k}) |\psi_{\vec{\alpha}}\rangle\right)\\ \bar{O}_{k} &= \langle\psi_{\vec{\alpha}} | O_{k} |\psi_{\vec{\alpha}}\rangle = \langle\langle O_{k}(x)\rangle\rangle\end{aligned}$$



But now we can ask how much we actually change the wf. when we change $\alpha'_k \to \alpha_k + \delta \alpha_k$ $ds^2 = ||\frac{|\Psi_{\vec{\alpha}'}\rangle}{||\Psi_{\vec{\alpha}'}||} - \frac{|\Psi_{\vec{\alpha}}\rangle}{||\Psi_{\vec{\alpha}}||}|^2 = \sum_{k,k'} \delta\alpha_k \delta\alpha_{k'} S_{k,k'}$ $S_{k,k'} = \operatorname{cov}(O_k, O_{k'})$

This matrix is also known as the Fisher information metric F=4 S of the probability $p_{\vec{\alpha}}(x) \propto \Psi_{\vec{\alpha}}(x)^2$





Thus instead of using the Euclidean metric we can use the more appropriate Fisher information metric to define our cost function:

$$\operatorname{Cost}(\delta \vec{\alpha}) = E(\vec{\alpha}) - \sum_{k} f_k \delta \alpha_k + \frac{ds^2}{2\Delta}$$
$$ds^2 = \sum_{k,k'} \delta \alpha_k \delta \alpha_{k'} S_{k,k'}$$

$$\operatorname{Min}_{\delta\vec{\alpha}}\operatorname{Cost}(\delta\vec{\alpha}) \to \delta\alpha_k = \Delta \sum_{k'} S_{k,k'}^{-1} f_{k'}$$



Steepest descent:

$$\delta \alpha_k = \Delta f_k$$

Newton-Raphson:

SR or natural gradients:

$$\delta \alpha_{k} = \sum_{k'} H_{k,k'}^{-1} f_{k'}$$
$$\delta \alpha_{k} = \Delta \sum_{k'} S_{k,k'}^{-1} f_{k'}$$

See: S. S. PRL **80**, 4558 (1998), PRB **64**, 240512 (2001) Natural gradients: S.I. Amari *Neural Computation* **10**, **251–276** (1998)







Removing most of the slowness of steepest





A more difficult optimization 😇



 $\varepsilon = 0.0001$

Benzene dimer: S.S., M. Casula and D. Rocca JCP 127, 014105 (2007)



Why is that?

It is well known that the steepest is slow for Hessian matrices with large condition number r: Namely #Iterations >~ r



(a)



Here is the condition matrix of SR in H₂



entangled → SR ill conditioned soon That's why for R=0.8 we need 400000 iteration



- Optimization of several parameters is nowdays possible within QMC thanks to SR or natural gradients
- Hessian method is in principle faster and a further improvement is possible by using partial information of it→ the linear method
- More popular optimized known in ML (e.g. ADAMS) are much less accurate and useless for our accuracy target.



At variance of normal parameters α the atomic coordinates R appear also in the Hamiltonian H_R and we need to compute R, α derivatives of:

$$E_{R,\alpha} = \frac{\langle \Psi_{\alpha,R} | H_R | \Psi_{\alpha,R} \rangle}{\langle \Psi_{\alpha,R} | \Psi_{\alpha,R} \rangle} = \sum_x \pi(x) e_L^{\alpha,R}(x)$$

where now: $e_L^{\alpha,R}(x) = \frac{\langle x | H | \Psi_{\alpha,R} \rangle}{\langle x | \Psi_{\alpha,R} \rangle}$

is also dependent upon both α and R



Start from correlated sampling

$$E_{R+dR,\alpha} = \frac{\sum_{x \in \text{samples}} \left| \frac{\Psi_{R+dR,\alpha}}{\Psi_{R,\alpha}} \right|^2(x) e_L^{R+dR,\alpha}(x)}{\sum_{x \in \text{samples}} \left| \frac{\Psi_{R+dR,\alpha}}{\Psi_{R,\alpha}} \right|^2(x)}$$

Thus by differentiating with respect to dR we obtain:



$$F_R = -\frac{dE_{R,\alpha}}{dR} = F_{\text{Helmann-Feynman}} + F_{\text{Pulay}}$$

$$F_{HF} = -\langle \partial_R e_L^{R,\alpha} \rangle = -\frac{1}{M} \sum_{i=1}^M \partial_R e_L^{R,\alpha}(x_i)$$

$$F_P = 2\text{Cov}(O_R, e_L^{\alpha,R}) = -\frac{2}{M} \sum_{i=1}^M O_R(x_i) \times (e_L^{\alpha,R}(x_i) - E_{VMC})$$

$$O_R(x_i) = \partial_R \log |\Psi_{\alpha,R}(x_i)|$$

14/07/2021





The variational parameters $\alpha(R)$ depend on R because we should expect correction o(dR) if we reoptimize the wf at R+dR

Thus we should add:
$$F = F_{HF} + F_P - \sum_{\alpha} (\partial_R \alpha) (\partial_{\alpha} E_{\alpha,R})$$

But if we have done a good optimization and are at an energy minimum:

 $\partial_{\alpha} E_{\alpha,R} = 0 \quad \forall \alpha$ because they are just the Euler's conditions of minimum energy (not necessarily the absolute one)



Assume $n = \Psi(x)$ defines a variable indicating the distance from the so called nodal surface n=0.

Then for
$$n \rightarrow 0$$
 $-\partial_R e_L^{\alpha,R}$ or $O_R e_L^{\alpha,R} \propto \frac{1}{n^2}$ $\pi(n) \propto n^2$

The mean of the above random variable is finite but their variance:

$$\int dn(n^2)(1/n^4) \to \infty$$



$$\pi(x) = \frac{|\Psi^{\alpha,R}(x)|^2}{R(x)} \qquad R(x) \propto n^2 \text{ for } n \to 0$$

Then we have to compute average of finite random variables:

$$-\partial e_L^{\alpha,R} \times R(x)$$
 and $-O_R e_L^{R,\alpha} R(x) \to \text{FINITE for } n \to 0$



Since n vanishes as our determinant/Pfaffian of a matrix M. Each element of M⁻¹ blows up as 1/n and therefore the choice:

$$R(x) = \frac{1}{1 + \sum_{ij} |\epsilon M_{ij}^{-1}|^2} \propto n^2$$

Satisfies the requirements and make the variance FINITE







Sometimes you can definitively solve a boring problem in QMC



X,R and WF parameters

→ eloc,logpsi

TASK1: $x, R \rightarrow dist(I,i)$ I=1,2,...,#ion (nion) ,i=1,2,...,#el (nel) electron-ion distances.

TASK2: x,R,dist→ kel,vpseudolocal,prefactor,ivic,tmu (all electron coordinates in the pseudo mesh and LRDMC mesh ivic, with matrix element coefficients: prefactor for pseudo mesh, and tmu for LRDMC laplacian discretization)

TASK3: $R \rightarrow iond(I,J)$ (ion-ion distance)

- TASK4: x,R Z's (exponents basis) \rightarrow winv(1:nelorb,0:indt+4,1:nel),winvj(nelorbj,0:indt+4,1:nel) (basis array)
- TASK5: winv,winvj,jasmat,detmat,detmat_c,mu_c \rightarrow winvbar,winvjbar (Geminal Jastrow matrices contracted with basis)
- TASK6: iond \rightarrow vpot (Classical Coulomb ion potential)
- TASK7: vpseudolocal \rightarrow vpot (add Ewald if PBC and pseudo contribution if present)
- TASK8: winv, winvbar \rightarrow logpsidetln, ainv (Inverse Geminal matrix for fast updates)

TASK9: x,R,winvj,winvjbar, Z's \rightarrow jastrowall_ee(1:nel,1:nel,0:indt+4),jastrowall_ei(1:nion,1:nel),tabpip (for fast updates Jastrow)

→ logpsi

TASK10: x,R,ainv,winvbar \rightarrow winvup,winvdo

TASK11: tabpip,winvup,winvdo \rightarrow eloc

The subroutine name is compute_eloc_logpsi and is badly written as many input are explicitly passed (huge call) and several others are passed via f90 modules. A refactoring of this routine could help for future developments.





TASK(INPUT→OUTPUT)

$TASK_b(\overline{OUTPUT} \rightarrow \overline{INPUT})$

$\overline{\text{INPUT}}^{i} = \overline{\text{INPUT}}^{i} + \frac{\partial \sum_{k} \overline{\text{OUTPUT}}^{k} \text{OUTPUT}_{k}}{\partial \text{INPUT}_{i}}$

$$\bar{r}_{ASK11_b} \qquad \bar{e}_L = 1 \quad \rightarrow \quad \begin{cases} \bar{\vec{r}}_i &= \partial_{\vec{r}_i} e_L \\ \bar{\vec{R}}_i &= \partial_{\vec{R}_i} e_L \end{cases}$$

$$\bar{r}_{ASK1_b} \qquad \ln(\bar{\Psi}) = 1 \quad \rightarrow \quad \begin{cases} \bar{\vec{r}}_i &= \partial_{\vec{R}_i} \ln(\Psi) \\ \bar{\vec{R}}_i &= \partial_{\vec{R}_i} \ln(\Psi) \end{cases}$$

Thus we can compute all necessary for ion forces and/or optimization GPU implementation is trivial, each task BLAS3 offload and that's it.



Cpu time referenced to simple VMC (only energy) for computing all 3M force components in water.



It is a theorem it should be at most 4 times more expensive Use of pseudopotentials & Jastrow → 11 instances

For DMC the overhead to Compute all forces is ~0 but there is approximation



Follow us fip company/trex-eu @trex_eu



Targeting Real Chemical Accuracy at the Exascale project has received funding from the European Union Horizon 2020 research and innovation programme under Grant Agreement **No. 952165.**