

# More Accurate QMC Calculations by Using Machine Learning

by

Even Marius Nordhagen

THESIS

for the degree of

MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences  
University of Oslo

March 3, 2019

# Abstract

Abstract should be written here

# Acknowledgements

First of all I would like to thank Christine for cooking delicious food and feeding me when I need it the most. She and her food have maintained my motivation in the best way, and they are the main reason why I've survived those two years of masters.

Secondly, my annoying cohabitants have forced me to spend more time at university. I'm so glad you behave the way you do, if you were more cozy I would certainly spend more time at home and this thesis would not be the same.

Acknowledgements should be written here

-One of my mottos used to be that everything has a reason, which I used to point out whenever people were talking about things that apparently could not be described immediately. It was first when I learned about quantum mechanics that I understood I was wrong, and that is one of the reasons why quantum mechanics caught me so hardly.



# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	Many-Body problem . . . . .	10
1.2	Machine learning . . . . .	10
1.3	Goals and milestones . . . . .	10
<b>2</b>	<b>Quantum Many-Body Physics</b>	<b>11</b>
2.1	Introductory Quantum Physics . . . . .	12
2.1.1	The Schrödinger Equation . . . . .	12
2.1.2	The Variational Principle . . . . .	13
2.1.3	Postulates of Quantum Mechanics . . . . .	13
2.2	The Trial Wave Function . . . . .	14
2.2.1	Bosons and fermions . . . . .	14
2.2.2	Slater determinant . . . . .	15
2.2.3	Electron structure calculations . . . . .	15
2.2.4	Cusp condition . . . . .	16
2.2.5	Basis set . . . . .	16
2.2.6	Jastrow factor . . . . .	16
<b>3</b>	<b>Systems</b>	<b>17</b>
3.1	Quantum dots . . . . .	17
3.2	Atomic systems . . . . .	18
3.3	Molecular systems . . . . .	18
3.4	Bose-Einstein Condensation . . . . .	18
3.5	Electron gas . . . . .	18
3.6	Helium gas . . . . .	18
<b>4</b>	<b>Quantum Monte-Carlo Methods</b>	<b>19</b>
4.0.1	Isotropic processes . . . . .	19
4.0.2	Anisotropic processes . . . . .	19
4.0.3	Variational Monte Carlo . . . . .	19
<b>5</b>	<b>Ab Initio Quantum Methods</b>	<b>23</b>
5.1	Hartree-Fock . . . . .	23
5.2	Configuration Interaction . . . . .	23
5.3	Coupled Cluster . . . . .	23
<b>6</b>	<b>Machine Learning</b>	<b>25</b>
6.1	Supervised Learning . . . . .	25
6.1.1	Linear regression . . . . .	25
6.1.2	Logistic regression . . . . .	26

6.1.3	Neural network . . . . .	29
6.2	Unsupervised Learning . . . . .	32
6.2.1	Statistical foundation . . . . .	32
6.2.2	Boltzmann Machines . . . . .	33
6.2.3	Restricted Boltzmann Machines . . . . .	33
<b>7</b>	<b>Optimization</b>	<b>35</b>
7.1	Optimization of Trial Wave Function . . . . .	35
7.1.1	Energy Minimization . . . . .	35
7.1.2	Variance Minimization . . . . .	35
7.2	Minimization Algorithms . . . . .	35
7.2.1	Gradient Descent . . . . .	35
7.2.2	Conjugate Gradient . . . . .	35
<b>8</b>	<b>Scientific Programming</b>	<b>37</b>
8.1	Object Orientated Programming . . . . .	37
8.1.1	Inheritance . . . . .	38
8.1.2	Pointers . . . . .	38
8.1.3	Virtual Functions . . . . .	38
<b>9</b>	<b>Implementation</b>	<b>39</b>
9.1	Foundation . . . . .	39
9.2	Structure . . . . .	39
9.3	Optimization algorithms . . . . .	40
9.3.1	Stochastic Gradient Descent . . . . .	40
9.3.2	ADAM . . . . .	40
<b>10</b>	<b>Results</b>	<b>41</b>
10.1	Comparison of Wave Functions . . . . .	41
10.2	Ground state energies . . . . .	41
<b>11</b>	<b>Conclusion and future work</b>	<b>43</b>
<b>A</b>	<b>Dirac notation</b>	<b>45</b>
<b>B</b>	<b>Scaling</b>	<b>47</b>
B.0.1	Harmonic Oscillator - Natural units . . . . .	47
B.0.2	Atomic systems - Atomic units . . . . .	47
<b>C</b>	<b>Machine Learning Wave Function</b>	<b>49</b>
<b>D</b>	<b>Wave Function Elements</b>	<b>51</b>
D.1	Kinetic Energy Calculations . . . . .	51
D.2	Parameter Update . . . . .	52
D.3	Derivatives . . . . .	52
D.3.1	Simple Gaussian . . . . .	52
D.3.2	Padé-Jastrow Factor . . . . .	53
D.3.3	Slater Determinant . . . . .	53
D.3.4	NQS-Gaussian . . . . .	53
D.3.5	NQS-Jastrow Factor . . . . .	53
D.3.6	Hydrogen-Like Orbitals . . . . .	53

## List of abbreviations

Letters	Meaning
RBM	- Restricted Boltzmann Machine
MC	- Monte Carlo
VMC	- Variational Monte Carlo
DMC	- Diffusion Monte Carlo
ML	- Machine Learning
WF	- Wave Function
SPF	- Single Particle Function

Table 1: List of symbols used with explanation.

## Source Code

The source code is given in <https://github.com/evenmn>



# Chapter 1

## Introduction

Properties and behavior of quantum many-body systems are determined by the laws of quantum physics which have been known since the 1930s. The time-dependent Schrödinger equation describes the bounding energy of atoms and molecules, as well as the interaction between particles in a gas. In addition, it has been used to determine the energy of artificial structures like quantum dots, nanowires and ultracold condensates. [Electronic Structure Quantum Monte Carlo Michal Bajdich, Lubos Mitas]

Even though we know the laws of quantum mechanics, many challenges are encountered when calculating real-world problems. First, interesting systems often involve large number of particles, which causes expensive calculations. Second, we do not have a good model for the three-body interaction, which is vital when it comes to strong correlations. Paul Dirac recognized those problems already in 1929,

*”The general theory of quantum mechanics is now almost complete... ...The underlying physical laws necessary for the mathematical theory of a large part of physics and the whole of chemistry are thus completely known, and the difficulty is only that the exact application of these laws leads to equations much too complicated to be soluble.”*

-Paul Dirac, Quantum Mechanics of Many-Electron Systems, 1929

”At the same time, advent of computer technology has offered us a new window of opportunity for studies of quantum (and many other) problems. It spawned a “third way” of doing science which is based on simulations, in contrast to analytical approaches and experiments. In a broad sense, by simulations we mean computational models of reality based on fundamental physical laws. Such models have value when they enable to make predictions or to provide new information which is otherwise impossible or too costly to obtain otherwise. In this respect, QMC methods represent an illustration and an example of what is the potential of such methodologies.” [Electronic Structure Quantum Monte Carlo Michal Bajdich, Lubos Mitas]

History of QMC before and after the invention of electronic computers. Enrico Fermi 1930s similarities between imaginary time Schrödinger equation and stochastic processes in statistical mechanics. Metropolis VMC early 1950s. Kalos Greens’s function Monte Carlo late 1950s. Ceperly and Alder 1980 homogeneous electron gas.

Multi scale calculations are... A field of interest is how a systems behave when the interaction gets weaker. One way to model this, is to have a harmonic oscillator with a decreasing system frequency.

With machine learning, a function can be fitted to everything as long as we can define a cost function to minimize. The purpose of this thesis is to construct optimal wave functions for different systems by using machine learning.

- Low frequency (weakly interacting electrons) field of interest - Multi scale calculations
- Cartesian - Introduce the wavefunction - Mention the uncertainty principle and also quantum entanglement to catch the readers interest

## **1.1 Many-Body problem**

Not possible to solve analytically

## **1.2 Machine learning**

Branch of artificial intelligence

## **1.3 Goals and milestones**

- Investigate a new method to solve the Many-Body problem

# Chapter 2

## Quantum Many-Body Physics

If you are not completely  
confused by quantum  
mechanics, you do not  
understand it.

---

John Wheeler

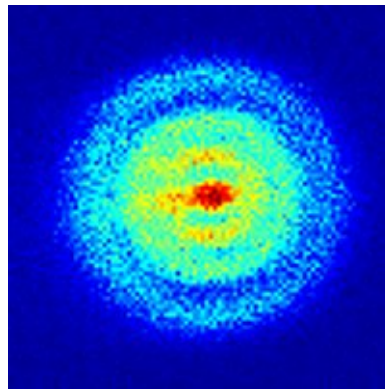


Figure 2.1: The first photograph of a Hydrogen atom, which shows the true wave function. Published in Phys. rev. lett. 110, 213001 (2017), *Hydrogen atoms under magnification*.

Around 1900, some physicists thought that there were nothing new to be discovered in physics and all that remained was more precise measurements, as Lord Kelvin famously pointed out. He could not be more wrong. In the following years, things were observed that could only be described by a quantized theory, led by Albert Einstein's explanation of the photoelectric effect in 1905.

Immense efforts were placed on completing the theory, and contributions from an array of scientists over a period of 20 years were necessary to get it finished. In 1929, Paul Dirac stated something similar to what Lord Kelvin said 30 years earlier, but hopefully with greater accuracy.

## 2.1 Introductory Quantum Physics

In this section we will present the fundamentals of the quantum theory, that will make up the framework of this project. The theory is based on David Griffith's incredible textbook, *Introduction to Quantum Mechanics*, where the reader is relegated for in-depth information.

Before we get started, we make a few assumptions in order to simplify our problem. Their most important ones are specified below with an explanation why they are valid.

- **Point-like particles:** First, all particles involved will be assumed to be point-like, such that they lack spatial extension. This includes the nucleus in atomic systems, but it still makes sense since the distance from the nucleus to the electrons is known to be much larger than the nucleus extent.
- **Non-relativistic spacetime:** Second, we operate in the non-relativistic spacetime, which is valid as long as we do not approach the speed of light and the forces are not strong. Applying classical physics, we can find that the speed of the electron in a hydrogen atom is about 1% of the speed of light, and even though the electrons get higher speed in heavier atoms, we do not need to worry about it. The forces acting are the weak Coulomb forces.
- For specific systems we might make new assumptions and approximate. For instance, for atomic systems we will assume that the nucleus is at rest. Those approximations will be discussed consecutively.

### 2.1.1 The Schrödinger Equation

The Schrödinger equation is a natural starting point, which gives the energy eigenstates of a system defined by a Hamiltonian  $\hat{\mathcal{H}}$  and its eigenfunctions,  $\Psi$ , which are the wave functions. The time-independent Schrödinger equation reads

$$\hat{\mathcal{H}}\Psi_n(\mathbf{r}) = \epsilon_n\Psi_n(\mathbf{r}) \quad (2.1)$$

where the Hamiltonian is the total energy operator. By analogy with the classical mechanics, this is given by

$$\hat{\mathcal{H}} = \hat{\mathcal{T}} + \hat{\mathcal{V}} \quad (2.2)$$

with  $\hat{\mathcal{T}}$  and  $\hat{\mathcal{V}}$  as the kinetic and potential energy operators respectively.

The kinetic energy yields  $T = p^2/2m$ , such that the kinetic energy operator can be represented as

$$\hat{\mathcal{T}} = \frac{\hat{\mathbf{p}}^2}{2m} \quad (2.3)$$

according to Ehrenfest's theorem. Further, the momentum operator is  $\hat{\mathbf{p}} = -i\hbar\nabla$ .

The potential energy can be split into an external part and an interaction part, where the latter is given by the Coulomb interaction. For two identical particles of charge  $q$ , the repulsive interaction gives the energy

$$V_I = k \frac{q^2}{r_{12}} \quad (2.4)$$

where  $r_{12}$  is the distance between the particles. The total Hamiltonian of a system of  $N$  identical particles takes the form

$$\hat{\mathcal{H}} = - \sum_i^N \frac{\hbar^2}{2m} \nabla_i^2 + \sum_i^N u_i + \sum_i^N \sum_{j>i}^M k \frac{q^2}{r_{ij}} \quad (2.5)$$

which is the farthest we can go without specifying the external potential  $u_i$ .  $r_{ij}$  is the relative distance between particle  $i$  and  $j$ , defined by  $r_{ij} \equiv |\mathbf{r}_i - \mathbf{r}_j|$ .

Setting up equation (2.1) with respect to the energies, we obtain an integral,

$$\epsilon_n = \frac{\int d\mathbf{r} \Psi_n^*(\mathbf{r}) \hat{\mathcal{H}} \Psi_n(\mathbf{r})}{\int d\mathbf{r} \Psi_n^*(\mathbf{r}) \Psi_n(\mathbf{r})}, \quad (2.6)$$

which not necessarily is trivial to solve. If we take the wave function squared we get the probability of finding a particle at a certain position,

$$P(\mathbf{r}) = \Psi_n^*(\mathbf{r}) \Psi_n(\mathbf{r}) = |\Psi_n(\mathbf{r})|^2 \quad (2.7)$$

so the nominator is simply the integral over all probabilities. If the wave function is normalized correctly, this should always give 1. Assuming that is the case, the expectation value can be expressed more elegantly by using Dirac notation,

$$E[\Psi] = \langle \Psi | \hat{\mathcal{H}} | \Psi \rangle, \quad (2.8)$$

where the first part,  $\langle \Psi |$  is called a bra and the last part,  $|\Psi \rangle$  is called a ket. At first this might look artificial and less informative, but it simplifies the notation significantly. More information about the notation is found in Appendix A.

In many cases we do not know the exact wave function, and need to rely on a trial wave function guess. Henceforth, we will use  $\Psi$  as the exact total wave function,  $\psi$  as the exact single particle function (SPF),  $\Psi_T$  as the total trial wave function and  $\phi$  as the trial SPF. [Gri05]

### 2.1.2 The Variational Principle

In the equations above, the presented wave functions are assumed to be the exact eigen functions of the Hamiltonian. But often we do not know the exact wave functions, and we need to guess what the wave functions might be. In those cases we make use of the variational principle, which states that only the exact ground state wave function is able to give the ground state energy. All other wave functions that fulfill the required properties (see section 2.2) give higher energies, and mathematically we can express the statement

$$\epsilon_0 \leq \langle \Psi_T | \hat{\mathcal{H}} | \Psi_T \rangle. \quad (2.9)$$

Variational Monte-Carlo is a method based on (and named after) the variational principle, where we vary the trial wave function in order to obtain the lowest energy. It will be detailed in chapter (4).

### 2.1.3 Postulates of Quantum Mechanics

The quantum theory is built on a few fundamental postulates, which will always be true. Before we go further, we will take a quick look at some of them.

1. Associated with any particle moving in a conservative field of force is a wave function which determines everything that can be known about the system.
2. With every physical observable  $q$  there is associated an operator  $Q$ , which when operating upon the wavefunction associated with a definite value of that observable will yield that value times the wavefunction.
3. Any operator  $Q$  associated with a physically measurable property  $q$  will be Hermitian.
4. The set of eigenfunctions of operator  $Q$  will form a complete set of linearly independent functions.
5. For a system described by a given wavefunction, the expectation value of any property  $q$  can be found by performing the expectation value integral with respect to that wavefunction.
6. The time evolution of the wavefunction is given by the time dependent Schrodinger equation.

## 2.2 The Trial Wave Function

By the first postulate of quantum mechanics, the wave function contains all the information specifying the state of the system. This means that all observable in classical mechanics can also be measured from the wave function, resulting finding the wave function as our main goal.

The true wave function is in most cases unknown, and we need an initial wave function based on educated guesses, called a trial wave function. This function has to meet some requirements in order to be used in the variational principle, including the CUSP condition, normalizable and ... In addition, the wave function should be symmetric or anti-symmetric under exchange of two coordinates.

### 2.2.1 Bosons and fermions

An introductory sentence is needed to increase the FLOW. Assume we have a permutation operator  $\hat{P}$  which exchanges two coordinates in the wave function,

$$\hat{P}(i \rightarrow j)\Psi_n(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_M) = p\Psi_n(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_M), \quad (2.10)$$

where  $p$  is just a factor which comes from the transformation. If we again apply the  $\hat{P}$  operator, we should switch the same coordinates back, and we expect to end up with the initial wave function. For that reason,  $p = \pm 1$ .<sup>1</sup>

The particles that have an antisymmetric (AS) wavefunction under exchange of two coordinates are called fermions, named after Enrico Fermi, and have half integer spin. On the other hand, the particles that have a symmetric (S) wavefunction under exchange of two coordinates are called bosons, named after Satyendra Nath Bose, and have integer spin.

---

<sup>1</sup>This was true until 1976, when J.M. Leinaas and J. Myrheim discovered the anyon, <https://www.uio.no/studier/emner/matnat/fys/FYS4130/v14/documents/kompendium.pdf>.

It turns out that because of their antisymmetric wavefunction, two identical fermions cannot be found at the same position at the same time, known as the Pauli principle. This causes some difficulties when dealing with multiple fermions, because we always need to ensure that the total wavefunction becomes zero if two identical particles happen to be at the same position. To do this, we introduce a Slater determinant as described in the next chapter. In this particular project, we are going to focus on electrons and therefore fermions. Anyway, much of the theory applies for bosons as well.

Read <https://manybodyphysics.github.io/FYS4480/doc/pub/secondquant/html/secondquant-bs.html>

### 2.2.2 Slater determinant

For a system of more particles we can define a total wavefunction, which is a composition of all the single particle wavefunctions (SPF) and contains all the information about the system. For fermions we need to compile the SPFs such that the Pauli principle is fulfilled at all times. One way to do this is by setting up the SPFs in a determinant, known as a Slater determinant.

Consider a system of two identical fermions with SPFs  $\phi_1$  and  $\phi_2$  at positions  $\mathbf{r}_1$  and  $\mathbf{r}_2$  respectively. The way we define the wavefunction of the system is then

$$\Psi_T(\mathbf{r}_1, \mathbf{r}_2) = \begin{vmatrix} \phi_1(\mathbf{r}_1) & \phi_2(\mathbf{r}_1) \\ \phi_1(\mathbf{r}_2) & \phi_2(\mathbf{r}_2) \end{vmatrix} = \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) - \phi_2(\mathbf{r}_1)\phi_1(\mathbf{r}_2), \quad (2.11)$$

which is set to zero if the particles are at the same position. This is called a Slater determinant, and yields the same no matter how big the system is.

Notice that we denote the wave function with the ' $T$ ', which indicates that it is a trial wave function. We do this because the spin part is avoided with  $\psi$  as the radial parts only, thus this wavefunction is not the true wavefunction. We will look closer at how we can factorize out the spin part later. The spin part is assumed to not affect the energies.

A general Slater determinant for a system of  $N$  particles takes the form

$$\Psi(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) = \begin{vmatrix} \psi_1(\mathbf{r}_1) & \psi_2(\mathbf{r}_1) & \dots & \psi_N(\mathbf{r}_1) \\ \psi_1(\mathbf{r}_2) & \psi_2(\mathbf{r}_2) & \dots & \psi_N(\mathbf{r}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \psi_1(\mathbf{r}_N) & \psi_2(\mathbf{r}_N) & \dots & \psi_N(\mathbf{r}_N) \end{vmatrix} \quad (2.12)$$

where the  $\psi$ 's are the true single particle wave functions, which are the tensor products

$$\psi = \phi \otimes \xi \quad (2.13)$$

with  $\xi$  as the spin part.

There is a similar "Slater permanent" which applies for bosons. It is similar to a determinant, but all negative signs are replaced by positive signs.

### 2.2.3 Electron structure calculations

For our purpose we will study fermions with spin  $\sigma = \pm 1/2$  only, which can be seen as electrons. In this particular case, the SPFs can be arranged in spin-up and spin-down

parts, such that the Slater determinant can be simplified to

$$\Psi(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) = \begin{vmatrix} \phi_1(\mathbf{r}_1)\xi_{\uparrow} & \phi_1(\mathbf{r}_1)\xi_{\downarrow} & \dots & \phi_{N/2}(\mathbf{r}_1)\xi_{\downarrow} \\ \phi_1(\mathbf{r}_2)\xi_{\uparrow} & \phi_1(\mathbf{r}_2)\xi_{\downarrow} & \dots & \phi_{N/2}(\mathbf{r}_2)\xi_{\downarrow} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(\mathbf{r}_N)\xi_{\uparrow} & \phi_1(\mathbf{r}_N)\xi_{\downarrow} & \dots & \phi_{N/2}(\mathbf{r}_N)\xi_{\downarrow} \end{vmatrix}. \quad (2.14)$$

This is called the wavefunction ansatz, because assumptions are raised, like two particles with opposite spins are found to be at the same position all the time, i.e, an equal number of fermions have spin up as spin down, are applied. Further the...

SHOULD END UP WITH SPLITTED DETERMINANTS HERE

For a detailed walkthrough, see appendix I in REF(The Stochastic Gradient Approximation: an application to Li nanoclusters, Daniel Nissenbaum).

### 2.2.4 Cusp condition

The Coulomb potential gives a cusp when two charged particles approach each other. This means that the wave function should also have a cusp, known as the Coulomb cusp condition, Kato theorem or just the cusp condition.

### 2.2.5 Basis set

To go further, we need to define a basis set,  $\phi_n(\mathbf{r})$  which should be chosen carefully based on the system.

### 2.2.6 Jastrow factor

The Jastrow factor is introduced in order to capture interaction properties, ...

We will use the Padé-Jastrow factor, which can be expressed as

$$J(\mathbf{r}; \beta, \gamma) = \exp \left( \sum_{i=1}^N \sum_{j=1}^i \frac{\beta(i, j) r_{ij}}{1 + \gamma r_{ij}} \right) \quad (2.15)$$

which gets small when  $r_{ij}$  gets small. If we recall that the probability distribution is given by the wave function squared, we see that the Padé-Jastrow factor gives a lower probability of finding particles close to each other, which is what we want.



# Chapter 3

## Systems

We must be clear that when it comes to atoms, language can be used only as in poetry.

---

Niels Bohr

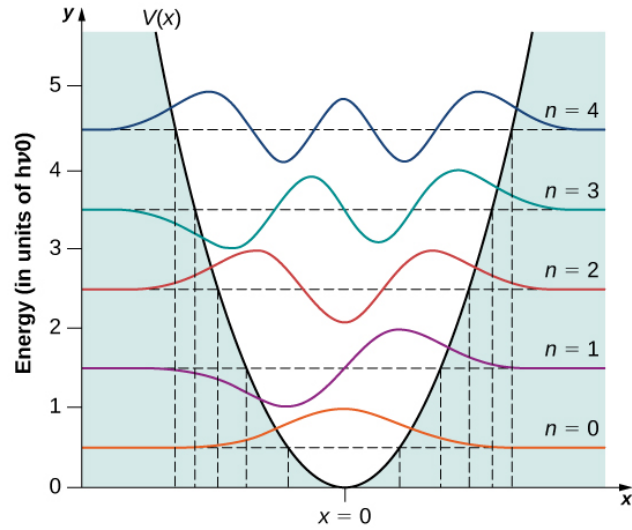


Figure 3.1: The quantum harmonic oscillator, with the Hermite functions represented up to 4th order. As in classical mechanics, the harmonic oscillator can describe various quantum systems, such as lattice vibration (phonons) and quantum fields.

When defining a system, we also need to specify the basis set to be used. The single particle functions are often known, and they are well-suited as a basis for the total

### 3.1 Quantum dots

Quantum dots are very small particles, and contain fermions or bosons hold together by an external potential. Since these particles have discrete electronic states like an atom, they are often called artificial atoms.

In this thesis we will study electrons trapped in harmonic oscillators, which give an

external potential affecting particle  $i$  is given by

$$u_i = \frac{1}{2}m\omega^2 r_i^2. \quad (3.1)$$

Using natural units as described in Appendix B, we can write the Hamiltonian as

$$\hat{\mathcal{H}} = \sum_{i=1}^P \left( -\frac{1}{2}\nabla_i^2 + \frac{1}{2}r_i^2 \right) + \sum_{i<j} \frac{1}{r_{ij}} \quad (3.2)$$

where the energy is scaled with respect to atomic units and lengths are scaled with respect to the Bohr radius.

The exact solutions of the non-interacting Hamiltonian are the Hermite functions,

$$f_n(x) = H_n(x) \exp(-x^2/2) \quad (3.3)$$

which we will use as our basis. The first four Hermite functions are illustrated in figure (??).

## 3.2 Atomic systems

We will also investigate real atoms, which have Hamiltonians defined by the Born-Oppenheimer approximation.

The External potential affecting particle  $i$  is therefore

$$u_i = -\frac{1}{2} \frac{Z}{r_i} \quad (3.4)$$

$$\hat{\mathcal{H}} = \sum_{i=1}^P \left( -\frac{1}{2}\nabla_i^2 - \frac{1}{2} \frac{Z}{r_i} \right) + \sum_{i<j} \frac{1}{r_{ij}}, \quad (3.5)$$

which applies for an atom with stationary nucleus at origin. The first part is related to the kinetic energy, while the second is

”Colloquially, we call such solutions and derived properties as electronic structure.”

## 3.3 Molecular systems

## 3.4 Bose-Einstein Condensation

## 3.5 Electron gas

## 3.6 Helium gas

He<sup>3</sup> and He<sup>4</sup>

# Chapter 4

## Quantum Monte-Carlo Methods

Some great methods are developed over the past decade. We will focus on the Hartree-Fock and variational Monte Carlo methods and give a detailed explanation of those methods. Additionally, the well-known methods configuration interaction and coupled cluster will be described briefly for some kind of completeness.

Monte Carlo methods in quantum mechanics are a bunch of methods that are built on diffusion processes, and includes Variational Monte Carlo (VMC), Diffusion Monte Carlo (DMC) and others. The common denominator is that we move particles in order to find the optimal configuration, usually where the energy is minimized. The particles can be moved isotropic, i.e., uniformly in all directions, or they can be affected by a drift force which makes the process anisotropic.

### 4.0.1 Isotropic processes

In isotropic processes, we have random walks where the particles move randomly in space which falls under the category Markov chains. If we assume constant timestep, the step index can be considered the time, thus we have a time-dependent probability density  $P(\mathbf{x}, t)$ . This probability density needs to satisfy the isotropic diffusion equation,

$$\frac{\partial P(\mathbf{x}, t)}{\partial t} = D \frac{\partial^2}{\partial \mathbf{x}^2} P(\mathbf{x}, t). \quad (4.1)$$

### 4.0.2 Anisotropic processes

For anisotropic processes, we have a drift and the moves are no longer considered random but falls still under the category Markov chains. Because of the drift, we need to rewrite the diffusion equation and we end up with the Fokker-Planck equation,

$$\frac{\partial P(\mathbf{x}, t)}{\partial t} = D \frac{\partial}{\partial \mathbf{x}} \left( \frac{\partial}{\partial \mathbf{x}} - F \right) P(\mathbf{x}, t) \quad (4.2)$$

which needs to be satisfied. The new positions in coordinate space are given as solution of the Langevin equation

$$\frac{\partial \mathbf{x}(t)}{\partial t} = D \mathbf{F}(\mathbf{x}(t)) + \eta \quad (4.3)$$

### 4.0.3 Variational Monte Carlo

The Variational Monte Carlo (hereafter, VMC) method is today widely used when it comes to the study of ground state properties of quantum mechanical systems. It is

a Monte Carlo method which makes use of Metropolis sampling, and has been used in studies of fermionic systems since the 1970's. [Deb14] If we go back to the variational principle in equation (2.9), we see that by choosing a wave function which satisfies the criteria, we will get an energy larger or equal to the ground state energy.

There are two main problems we need to solve

1. We seldomly know the correct wave function
2. The integral we need to find the energy is hard or impossible to solve

Let us take the last problem first. The solution is to approximate the integral with a sum,

$$\begin{aligned}
 E &\leq \frac{\int \Psi_T(\mathbf{r})^* \hat{H} \Psi_T(\mathbf{r}) d\mathbf{r}}{\int \Psi_T(\mathbf{r})^* \Psi_T(\mathbf{r}) d\mathbf{r}} \\
 &= \int P(\mathbf{r}) E_L(\mathbf{r}) d\mathbf{r} \\
 &\approx \frac{1}{N} \sum_{i=1}^N E_L(\mathbf{r}_i)
 \end{aligned} \tag{4.4}$$

where the local energy is defined as

$$E_L(\mathbf{r}) \equiv \frac{1}{\Psi_T(\mathbf{r})} \hat{H} \Psi_T(\mathbf{r}) \tag{4.5}$$

and the  $\mathbf{r}_i$  is withdrawn from the probability distribution  $P(\mathbf{r})$ , which is given by

$$P(\mathbf{r}) = \frac{|\Psi_T(\mathbf{r})|^2}{\int |\Psi_T(\mathbf{r})|^2 d\mathbf{r}}. \tag{4.6}$$

The energy found from equation (4.4) is an expectation value, and we therefore know that the true energy lies within the standard error. When increasing the number of Monte-Carlo cycles, the standard error decreases and we get a more accurate energy. In the limit  $N$  goes to infinity, the variance should approach zero,

$$\langle E \rangle = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N E_L(\mathbf{r}_i). \tag{4.7}$$

For more statistical details, see [Deb14].

Given a wave function we are now able to estimate the corresponding energy, but how do we find a good wave function? In VMC, we define a wave function with variational parameters, which are adjusted in order to minimize the energy for every iteration. For every iteration, we run  $N$  Monte-Carlo cycles where we withdraw a new position  $\mathbf{r}_i$ . Whether or not the proposed move should be accepted is determined by the Metropolis algorithm.

#### 4.0.3.1 The Metropolis Algorithm

The genius of the metropolis algorithm, is that the acceptance of a move is not based on the probabilities themselves, but the ratio between the new and the old probability. In that way, we avoid calculating the sum over all probabilities, which often is expensive or even impossible to calculate.

In its simplest form, the move is proposed randomly, and it is accepted if the ratio is larger than a random number between 0 and 1. However, with this approach a lot of moves will be rejected, which wastes CPU time. A better method is **importance sampling**, which makes a educated guess of the best way to move based on diffusion processes, and move the particle in that direction.

A time-dependent probability density needs to satisfy the Fokker-Planck equation



# Chapter 5

## Ab Initio Quantum Methods

The term *ab initio* means from first principles, implying that only physical constants are put into the methods. The Monte-Carlo methods are not considered *ab initio* as nonphysical hyper parameters are required. Some example methods are Hartree-Fock methods and post-Hartree-Fock methods like Configuration Interaction and Coupled Cluster, which will be discussed in the following.

### 5.1 Hartree-Fock

Describe this detailed

### 5.2 Configuration Interaction

The completeness relation

### 5.3 Coupled Cluster

The coupled cluster method is the *de facto* standard wave function-based method for electronic structure calculations. [Pal05] The method approximates the wave function with an exponential expansion,

$$|\Psi_{CC}\rangle = e^{\hat{T}} |\Phi_0\rangle \quad (5.1)$$

where  $\hat{T}$  is the cluster operator, entirely given by  $\hat{T} = \hat{T}_1 + \hat{T}_2 + \hat{T}_3 + \dots$  with

$$\hat{T}_n = \left(\frac{1}{n!}\right)^2 \sum_{abc\dots} \sum_{ijk\dots} t_{ijk\dots}^{abc\dots} a_a^\dagger a_b^\dagger a_c^\dagger \cdots a_k a_j a_i. \quad (5.2)$$

We again want to solve the Schrödinger equation,

$$\hat{H} |\Psi\rangle = \hat{H} e^{\hat{T}} |\Phi_0\rangle = \epsilon e^{\hat{T}} |\Phi_0\rangle, \quad (5.3)$$

which can be simplified by multiplying with  $e^{-\hat{T}}$  from the left. This introduces us to the **similarity transformed Hamiltonian**

$$\bar{H} = e^{-\hat{T}} \hat{H} e^{\hat{T}}. \quad (5.4)$$

If we on one hand now multiply with the reference bra on the left hand side, we easily observe that

$$\langle \Phi_0 | \bar{H} | \Phi_0 \rangle = \epsilon \quad (5.5)$$

which is the coupled cluster energy equation. On the other hand, we can multiply with an excited bra on left hand side, and find that

$$\langle \Phi_{ijk\dots}^{abc\dots} | \bar{H} | \Phi_0 \rangle = 0 \quad (5.6)$$

which are the coupled cluster amplitude equations. The similarity transformed Hamiltonian can be rewritten using the Baker-Campbell-Hausdorff expansion

$$\begin{aligned} \bar{H} &= \hat{H} + [\hat{H}, \hat{T}] \\ &\quad + \frac{1}{2} [[\hat{H}, \hat{T}], \hat{T}] \\ &\quad + \frac{1}{6} [[[ \hat{H}, \hat{T} ], \hat{T}], \hat{T}] \\ &\quad + \frac{1}{24} [[[[ \hat{H}, \hat{T} ], \hat{T}], \hat{T}], \hat{T}] \\ &\quad + \dots \end{aligned} \quad (5.7)$$

and we are in principle set to solve the amplitude equations with respect to the amplitudes  $t_{ijk\dots}^{abc\dots}$  and then find the energy. The expansion is able to reproduce the true wave function exactly using a satisfying number of terms and an infinite basis. This is, of course, not possible, but even by limiting us to the first few coupled cluster operators, the results are often good compared to other methods. [DCFSI07]



# Chapter 6

## Machine Learning

The use of the term *machine learning* has exploded over the past years, and sometimes it sounds like it is a totally new field. However, the truth is that many of the methods we use are quite old, where for instance *linear regression* was known early in the 19th century. [Leg][Gau] Those methods have just recently been taken under the *machine learning* umbrella, which is one of the reasons why the term is mentioned so often. As the UiO professor Anne Solberg pointed out during one of her lectures

*"In the early 1990's I was working a lot with machine learning, but at that time we called it pattern recognition and regression..."*

Another reason why machine learning has an increasingly popularity is that some of the algorithms have been significantly improved. *Convolutional Neural Networks* (CNNs) are now as good as humans when it comes to object recognition in images, and *Long Short-Term Memory Recurrent Neural Networks* (LSTM RNNs) has improved voice recognition.

We can conclude that Machine learning includes all methods where we want to fit models to data sets, and pattern recognition and regression are clearly some of those methods. In both cases we have concrete targets for the training, and the training is therefore called supervised training. In other cases we do not have targets, but want to train our model based on a probability distribution or so, which is called unsupervised learning.

### 6.1 Supervised Learning

In supervised learning methods, the input data has known targets such that we can fit a model to give the correct outputs. Linear regression is perhaps the most intuitive example on this, where we want to find the line that fits some data points in the best possible way. This corresponds to a neural network without any hidden layer, but when we add more layers the model is no longer linear and it all gets more complex. For simple classification tasks, logistic regression can be used.

#### 6.1.1 Linear regression

In linear regression, the dependent variable  $y_i$  is a linear combination of the parameters, and for a dependent variable this can be written as

$$y_i = \sum_j X_{ij} \beta_j \tag{6.1}$$

where  $\beta_j$ 's are the unknown parameters to be found. In principle,  $X_{ij}$  can be an arbitrary function of the arguments  $x_i$ , but often one wants a polynomial model which corresponds to  $X_{ij} = x_i^j$ .

The three most commonly used linear regression methods are Ordinary Least Square (OLS) regression, Ridge regression and Lasso regression, where the former has the cost function

$$c(\vec{\beta}) = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2 \quad \text{OLS,} \quad (6.2)$$

which is minimized when

$$\vec{\beta} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \vec{y}. \quad (6.3)$$

Similarly, the Ridge cost function is

$$c(\vec{\beta}) = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad \text{Ridge} \quad (6.4)$$

where  $\lambda$  is called the penalty. This is minimized when

$$\vec{\beta} = (\hat{X}^T \hat{X} + \lambda I)^{-1} \hat{X}^T \vec{y}, \quad (6.5)$$

and finally the Lasso cost function is given by

$$c(\vec{\beta}) = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j \quad \text{Lasso.} \quad (6.6)$$

### 6.1.2 Logistic regression

Despite its name, logistic regression is not a fitting tool, but rather a classification tool. Traditionally, the perceptron model was used for 'hard classification', which sets the outputs directly to binary values. However, often we are interested in the probability of a given category, which means that we need a continuous *activation function*. Logistic regression can, like linear regression, be considered as a function where coefficients are adjusted with the intention to minimize the error. Here, the coefficients are called *weights*. The process goes like this: The inputs are multiplied by several weights, and by adjusting those weights the model can classify every *linear classification problem*. A drawing of the perceptron is found in figure (6.1).

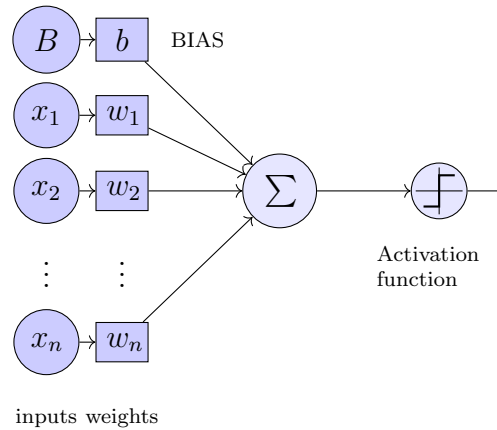


Figure 6.1: Logistic regression model with  $n$  inputs.

In logistic regression, we usually have one binary output node for each class, but for two categories one output node is sufficient, which can be fired or not fired.

Initially, one needs to train the perceptron such that it knows which outputs are correct, and for that one needs to know the outputs that correspond to the inputs. Every time the network is trained, the weights are adjusted such that the error is minimized.

The very first step is to calculate the initial outputs (forward phase), where the weights usually are set to small random numbers. Then the error is calculated, and the weights are updated to minimize the error (backward phase). So far so good.

### 6.1.2.1 Forward phase

Let us look at it from a mathematical perspective, and calculate the net output. The net output seen from an output node is simply the sum of all the "arrows" that point towards the node, see figure (6.1), where each "arrow" is given by the left-hand node multiplied with its respective weight. For example, the contribution from input node 2 to the output node follows from  $X_2 \cdot w_2$ , and the total net output to the output  $O$  is therefore

$$net = \sum_{i=1}^I x_i \cdot w_i + b \cdot 1. \quad (6.7)$$

Just some notation remarks:  $x_i$  is the value of input node  $i$  and  $w_i$  is the weight which connects input  $i$  to the output.  $b$  is the bias weight, which we will discuss later.

You might wonder why we talk about the net output all the time, do we have other outputs? If we look at the network mathematically, what we talk about as the net output should be our only output. Anyway, it turns out to be convenient mapping the net output to a final output using an activation function, which is explained further in section 6.1.2.5. The activation function,  $f$ , takes in the net output and gives the output,

$$out = f(net). \quad (6.8)$$

If not everything is clear right now, it is fine. We will discuss the most important concepts before we dive into the maths.

### 6.1.2.2 BIAS

As mentioned above, we use biases when calculating the outputs. The nodes, with value  $B$ , are called the bias nodes, and the weights,  $b$ , are called the bias weights. But why do we need those?

Suppose we have two inputs of value zero, and one output which should not be zero (this could for instance be a NOR gate). Without the bias, we will not be able to get any other output than zero, and in fact the network would struggle to find the right weights even if the output had been zero.

The bias value  $B$  does not really matter since the network will adjust the bias weights with respect to it, and is usually set to 1 and ignored in the calculations. [2]

### 6.1.2.3 Learning rate

In principle, the weights could be updated without adding any learning rate ( $\eta = 1$ ), but this can in practice be problematic. It is easy to imagine that the outputs can be fluctuating around the targets without decreasing the error, which is not ideal, and a

learning rate can be added to avoid this. The downside is that with a low learning rate the network needs more training to obtain the correct results (and it takes more time), so we need to find a balance.

#### 6.1.2.4 Cost function

The cost function is what defines the error, and in logistic regression the cross-entropy function is a naturally choice. [3] It reads

$$c(\mathbf{W}) = - \sum_{i=1}^n \left[ y_i \log f(\mathbf{x}_i^T \mathbf{W}) + (1 - y_i) \log[1 - f(\mathbf{x}_i^T \mathbf{W})] \right] \quad (6.9)$$

where  $\mathbf{W}$  contains all weights, included the bias weight ( $\mathbf{W} \equiv [b, \mathbf{W}]$ ), and similarly does  $\mathbf{x}$  include the bias node, which is 1;  $\mathbf{x} \equiv [1, \mathbf{x}]$ . Further, the  $f(x)$  is the activation function discussed in next section.

The cross-entropy function is derived from likelyhood function, which famously reads

$$p(y|x) = \hat{y}^y \cdot (1 - \hat{y})^{1-y}. \quad (6.10)$$

Working in the log space, we can define a log likelyhood function

$$\log [p(y|x)] = \log [\hat{y}^y \cdot (1 - \hat{y})^{1-y}] \quad (6.11)$$

$$= y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad (6.12)$$

which gives the log of the probability of obtaining  $y$  given  $x$ . We want this quantity to increase then the cost function is decreased, so we define our cost function as the negative log likelyhood function. [7]

Additionally, including a regularization parameter  $\lambda$  inspired by Ridge regression is often convenient, such that the cost function is

$$c(\mathbf{W})^+ = c(\mathbf{W}) + \lambda \|\mathbf{W}\|_2^2. \quad (6.13)$$

We will later study how this regularization affects the classification accuracy.

#### 6.1.2.5 Activation function

Above, we were talking about the activation function, which is used to activate the net output. In binary models, this is often just a step function firing when the net output is above a threshold. For continuous outputs, the logistic function given by

$$f(x) = \frac{1}{1 + e^{-x}}. \quad (6.14)$$

is usually used in logistic regression to return a probability instead of a binary value. This function has a simple derivative, which is advantageous when calculating a large number of derivatives. As shown in section ??, the derivative is simply

$$\frac{df(x)}{dx} = x(1 - x). \quad (6.15)$$

$\tanh(x)$  is another popular activation function in logistic regression, which more or less holds the same properties as the logistic function.

### 6.1.2.6 Backward phase

Now all the tools for finding the outputs are in place, and we can calculate the error. If the outputs are larger than the targets (which are the exact results), the weights need to be reduced, and if the error is large the weights need to be adjusted a lot. The weight adjustment can be done by any minimization method, and we will look at a couple of gradient methods. To illustrate the point, we will stick to the **gradient descent** (GD) method in the calculations, even though other methods will be used later. The principle of GD is easy: each weight is "moved" in the direction of steepest slope,

$$\mathbf{w}^+ = \mathbf{w} - \eta \cdot \frac{\partial c(\mathbf{w})}{\partial \mathbf{w}}, \quad (6.16)$$

where  $\eta$  is the learning rate and  $c(\mathbf{w})$  is the cost function. We use the chain rule to simplify the derivative

$$\frac{\partial c(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial c(\mathbf{w})}{\partial out} \cdot \frac{\partial out}{\partial net} \cdot \frac{\partial net}{\partial \mathbf{w}} \quad (6.17)$$

where the first is the derivative of the cost function with respect to the output. For the cross-entropy function, this is

$$\frac{\partial c(\mathbf{w})}{\partial out} = -\frac{y}{out} + \frac{1-y}{1-out}. \quad (6.18)$$

Further, the second derivative is the derivative of the activation function with respect to the output, which is given in (6.15)

$$\frac{\partial out}{\partial net} = out(1-out). \quad (6.19)$$

The latter derivative is the derivative of the net output with respect to the weights, which is simply

$$\frac{\partial net}{\partial \mathbf{w}} = \mathbf{x}. \quad (6.20)$$

If we now recall that  $out = f(\mathbf{x}^T \mathbf{w})$ , we can write

$$\frac{\partial c(\mathbf{w})}{\partial \mathbf{w}} = [f(\mathbf{x}^T \mathbf{w}) - \mathbf{y}] \mathbf{x} \quad (6.21)$$

and obtain a weight update algorithm

$$\mathbf{w}^+ = \mathbf{w} - \eta \cdot [f(\mathbf{x}^T \mathbf{w}) - \mathbf{y}]^T \mathbf{x}. \quad (6.22)$$

where the bias weight is included implicitly in  $\mathbf{w}$  and the same applies for  $\mathbf{x}$ .

## 6.1.3 Neural network

If you have understood logistic regression, understanding a neural network should not be a difficult task. According to **the universal approximation theorem**, a neural network with only one hidden layer can approximate any continuous function. [8] However, often multiple layers are used since this tends to give fewer nodes in total.

In figure (6.2), a two-layer neural network (one hidden layer) is illustrated. It has some similarities with the logistic regression model in figure (6.1), but a hidden layer and

multiple outputs are added. In addition, the output is no longer probabilities and can take any number, which means that we do not need to use the logistic function on the outputs anymore.

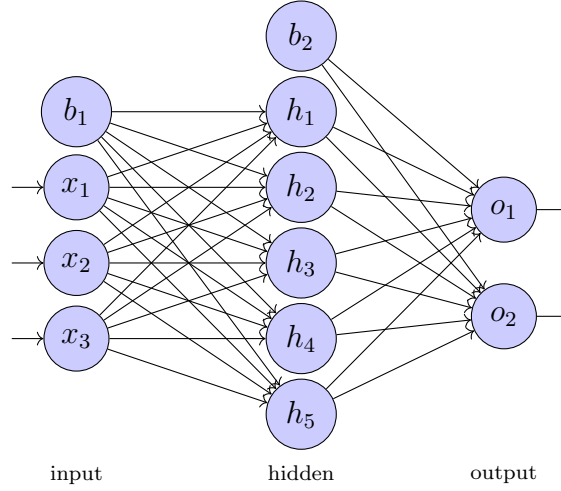


Figure 6.2: Neural network with 3 input nodes, 5 hidden nodes and 2 output nodes, in addition to the bias nodes.

Without a hidden layer, we have seen that the update of weights is quite straight forward. For a neural network consisting of multiple layers, the question is: how do we update the weights when we do not know the values of the hidden nodes? And how do we know which layer causing the error? This will be explained in section 6.1.3.3, where one of the most popular techniques for that is discussed. Before that we will generalize the forward phase presented in logistic regression.

### 6.1.3.1 Forward phase

In section 6.1.2.1, we saw how the output is found for a single perceptron. Since we only had one output node, the weights could be stored in an array. Generally, it is more practical to store the weights in matrices, since they will have indices related to both the node on left-hand side and the node on the right-hand side. For instance, the weight between input node  $x_3$  and hidden node  $h_5$  in figure (6.2) is usually labeled as  $w_{35}$ . Since we have two layers, we also need to denote which weight set it belongs to, which we will do by a superscript ( $w_{35} \Rightarrow w_{35}^{(1)}$ ). In the same way,  $\mathbf{W}^1$  is the matrix containing all  $w_{ij}^{(1)}$ ,  $\mathbf{x}$  is the vector containing all  $x_i$ 's and so on. We then find the net outputs at the hidden layer to be

$$net_{h,j} = \sum_{i=1}^I x_i \cdot w_{ij}^{(1)} = \mathbf{x}^T \mathbf{W}_j^{(1)} \quad (6.23)$$

where the  $\mathbf{x}$  and  $\mathbf{W}^{(1)}$  again are understood to take the biases. This will be the case henceforth. The real output to the hidden nodes will be

$$h_j = f(net_{h,j}). \quad (6.24)$$

Further, we need to find the net output to the output nodes, which is obviously just

$$net_{o,j} = \sum_{i=1}^H h_i \cdot w_{ij}^{(2)} = \mathbf{h}^T \mathbf{W}_j^{(2)} \quad (6.25)$$

We can easily generalize this. Looking at the net output to a hidden layer  $l$ , we get

$$\mathbf{net}_{h_l} = \mathbf{h}^{(l-1)T} \mathbf{W}^{(l)}. \quad (6.26)$$

### 6.1.3.2 Activation function

Before 2012, the logistic, the tanh and the pure linear functions were the standard activation functions, but then Alex Krizhevsky published an article where he introduced a new activation function called *Rectified Linear Units (ReLU)* which outperformed the classical activation functions. [4] The network he used is now known as AlexNet, and helped to revolutionize the field of computer vision. [5] After that, the ReLU activation function has been modified several times (avoiding zero derivative among others), and examples of innovations are *leaky ReLU* and *Exponential Linear Unit (ELU)*. All those networks are linear for positive numbers, and small for negative numbers. Often, especially in the output layer, a straight linear function is used as well.

In figure (??), *standard RELU*, *leaky RELU* and *ELU* are plotted along with the logistic function.

### 6.1.3.3 Backward Propagation

Backward propagation is probably the most used technique for updating the weights, and is actually again based on equation (6.16). What differs, is the differentiation of the net input with respect to the weight, which gets more complex as we add more layers. For one hidden layer, we have two sets of weights, where the last layer is updated in a similar way as for a network without hidden layer, but the inputs are replaced with the values of the hidden nodes:

$$w_{ij}^{(2)+} = w_{ij}^{(2)} - \eta \cdot [f(h_i^T w_{ij}) - y_j]^T h_i. \quad (6.27)$$

We recognize the first part as  $\delta_{ok}$ , such that

$$w_{ij}^{(1)+} = w_{ij}^{(1)} - \eta \cdot \sum_{k=1}^O \delta_{ok} \cdot w_{jk}^{(2)} \cdot f'(out_{hj}) \cdot x_i \quad (6.28)$$

where we recall  $\delta_{ok}$  as

$$\delta_{ok} = -(t_{ok} - out_{ok}) \cdot f'(out_{ok}).$$

For more layers, the procedure is the same, but we keep on inserting the obtained outputs from various layers.

### 6.1.3.4 Summary

Since it will be quite a lot calculations, I will just express the results here, and move the calculations to Appendix B. The forward phase in a three-layer perceptron is

$$\begin{aligned}
 net_{hi} &= \sum_j w_{ji}^{(1)} \cdot x_j \\
 out_{hi} &= f(net_{hi}) \\
 net_{ki} &= \sum_j w_{ji}^{(2)} \cdot out_{hj} \\
 out_{ki} &= f(net_{ki}) \\
 net_{oi} &= \sum_j w_{ji}^{(3)} \cdot out_{kj} \\
 out_{oi} &= f(net_{oi})
 \end{aligned} \tag{6.29}$$

which can easily be turned into vector form. The backward propagation follows from the two-layer example, and we get

$$\begin{aligned}
 w_{ij}^{(3)} &= w_{ij}^{(3)} - \eta \cdot \delta_{oj} \cdot out_{ki} \\
 w_{ij}^{(2)} &= w_{ij}^{(2)} - \eta \sum_{k=1}^O \delta_{ok} \cdot w_{jk}^{(3)} \cdot f'(out_{kj}) \cdot out_{hi} \\
 w_{ij}^{(1)} &= w_{ij}^{(1)} - \eta \sum_{k=1}^O \sum_{l=1}^K \delta_{ok} \cdot w_{lk}^{(3)} \cdot f'(out_{kl}) \cdot w_{jl}^{(2)} f'(out_{hj}) \cdot x_i
 \end{aligned}$$

where we again use the short hand

$$\delta_{oj} = (t_j - out_{oj}) \cdot f'(out_{oj}).$$

If we compare with the weight update formulas for the two-layer case, we recognize some obvious connections, and it is easy to imagine how we can construct a general weight update algorithm, no matter how many layers we have.

Now over to the problem we want to solve using neural networks.

## 6.2 Unsupervised Learning

How can we train using unsupervised learning?

### 6.2.1 Statistical foundation

Bayesian statistics



### 6.2.2 Boltzmann Machines

Boltzmann Machines are based on the more primitive Hopfield network, where a system of nodes is set up which defines the system energy. Inspired by statistical mechanics, the probability of finding the system in a state of energy  $E$  is given by

$$P = \exp(-E/kT) \quad (6.30)$$

which is the Boltzmann distribution, hence the name Boltzmann machine.  $k$  is known as Boltzmann's constant and  $T$  is the system temperature, but henceforth they both will be avoided by scaling  $E' = E/kT$ .

In the most general form, all nodes are connected to all other nodes, that is an unrestricted Boltzmann machine, see figure .. for illustration.

INSERT FIGURE

In the same manner as in a Feed-forward Neural Network, we can directly multiply each node  $s_i$  with its respective inner weights  $w_{ij}$  and then with the opposite node  $s_j$ , which defines the system energy together with the bias contributions. This gives the energy

$$E = - \sum_{i=1}^N \sum_{j=i}^N s_i w_{ij} s_j - \sum_{i=1}^N s_i b_i \quad (6.31)$$

for a system of  $N$  nodes, which is the so-called binary-binary network and the most basic architecture. During training, the weights are adjusted in order to maximize the probability...

### 6.2.3 Restricted Boltzmann Machines

When there is an unrestricted guy, a restricted guy must exist as well.



# Chapter 7

## Optimization

### 7.1 Optimization of Trial Wave Function

#### 7.1.1 Energy Minimization

#### 7.1.2 Variance Minimization

### 7.2 Minimization Algorithms

Gradient methods

#### 7.2.1 Gradient Descent

#### 7.2.2 Conjugate Gradient



# Chapter 8

## Scientific Programming

Since this thesis is much about writing code, it is natural including a few words about it.

The computer's language itself is binary, and is the lowest level. To translate commands to this language, we need a "translator", which is a language that fills the gap between the binary language and human commands. This language is categorized in levels based on how similar they are to the binary language. Low-level languages are similar to the binary language, which means fast but complicated. High-level languages are easy to work with, but are not as fast as low-level languages. Might mention grammar etc

One can either do *procedural programming* or *object oriented programming*. The former means that the code is written in the same order as the program flow goes, while one in the latter defines objects.

### 8.1 Object Orientated Programming

In the everyday life, we are surrounded by objects all the time which we can place in different categories. For instance, a *cat* is an object with a name, race, age and so on, and can be placed in the class *animals*. In object oriented programming, the class could be implemented as

```
class Animal:
    def __init__(self, animal, name, race, age):
        self.animal = animal
        self.name = name
        self.race = race
        self.age = age

    def __call__(self):
        return "%s is a %s that's %d years old and of the race %s"%(self.name, self.animal,
            ↪ self.age, self.race)
```

where...

The next step is to define an object, which in our case is the cat with name "Schrodinger":

```
Alma = Animal("cat", "Schroedeger", "Ragdoll", 4)
print(Alma())
```

This implies that "Schrodinger" is a cat of race "Ragdoll" and of age 4. When calling this class from "Schrodinger", the class returns

```
$python3 simple_class.py
>>> Schrodinger is a cat that's 4 years old and of the race Ragdoll
```

You might wonder how this is related to scientific programming. The answer is that it is often convenient to define various parts of the code as objects to increase the liability

and maximize reuse of code. For example, we use various Hamiltonians, where each can be defined as a subclass of the Hamiltonian superclass.

The code above is written in Python, but the exact task could be performed in C++.

INCLUDE C++ IMPLEMENTATION

As one can see,

In C++ one needs to define constructors and destructors. All variables used inside the functions are defines in the constructor, while they are removed in the destructor to free up memory. In Python, the constructors are called `__init__` by default, and memory is handled automatically.

### 8.1.1 Inheritance

This is also called parent and child, respectively.

Parent and child Polymorphy: Child inherit from the parents. Virtual functions to achieve runtime polymorphism Should define virtual destructor as well

1. Single inheritance
2. Multiple inheritances

Python and C++ support multiple inheritances. Multilevel inheritance: Got child and grand child. Hierarchical inheritance: Parent got several children. <https://www.geeksforgeeks.org/inheritance-in-python/>

### 8.1.2 Pointers

Sometimes we do not want to send the object itself, but either its address, such that..

### 8.1.3 Virtual Functions

Often one wants to define a template of objects... where the super class defines which functions its objects should have. In C++, this can be achieved by virtual functions, functions with arguments specified but task undefined. Those functions are overwritten by the corresponding functions in the object (hence virtual),

# Chapter 9

## Implementation

In many cases, planning is half the job, and so is true for a good VMC implementation. In fact, the program was restructured three times before we landed on this final version. It is based on Morten Ledum's VMC framework [<https://github.com/mortele/variational-monte-carlo-fys4411>], which was meant as an example implementation in the course *FYS4411 - Computational Physics II: Quantum Mechanical Systems*.

During the implementations we had two main areas of focus

- Efficient code
- High legibility, tidy code

Unlike many other VMC codes, our code was developed flexible with respect to the wave functions. This means that one can combine various wave function elements, where each element is implemented separately. For instance, the simple Gaussian function and the Padé-Jastrow factor were implemented separately, but they can easily be combined. The way one does this in practice, is to append multiple wave function elements to the vector `WaveFunctionElements` in main. If one for instance want to combine the simple Gaussian with the Padé-Jastrow factor, this can be done by

```
System* system = new System();
std::vector<class WaveFunction*> WaveFunctionElements;
WaveFunctionElements.reserve(2);
WaveFunctionElements.push_back(new class SimpleGaussian(system, alpha));
WaveFunctionElements.push_back(new class PadeJastrow(system, beta, Gamma));
system->setWaveFunction(WaveFunctionElements);
```

The big advantage of this implementation is that we do not need to hard code every possible combination of wave function elements, which reduces the number of code lines significantly. This also eases adding new elements, since we only need to calculate the derivatives of the particular element (do not need to worry about cross terms). Exactly how this is done can be read in Appendix D.

The con is that the program will be slightly slower.

### 9.1 Foundation

### 9.2 Structure

Add structure chart

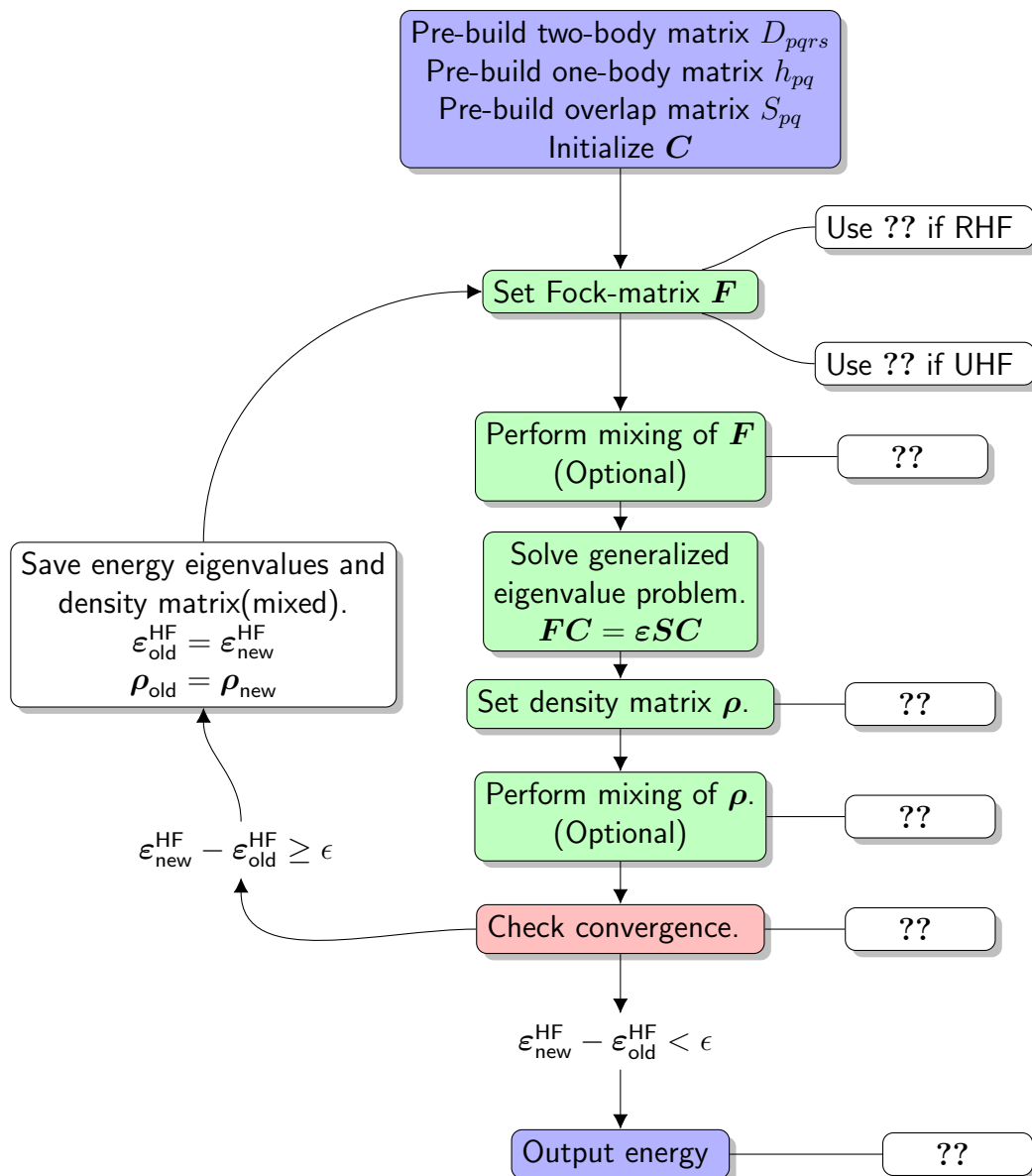


Figure 9.1: Hartree-Fock Algorithm.

## 9.3 Optimization algorithms

### 9.3.1 Stochastic Gradient Descent

### 9.3.2 ADAM



# Chapter 10

## Results

### 10.1 Comparison of Wave Functions

Compare machine learning wave functions to standard VMC wave functions

### 10.2 Ground state energies



# Chapter 11

## Conclusion and future work

See if machine learning is able to describe the three-body interaction, with nuclear physics applications.



# Appendix A

## Dirac notation

The Dirac notation, also called bracket notation, was suggested by Paul Dirac in a 1939 paper with the purpose of improving the reading ease. [Dir39]



# Appendix B

## Scaling

### B.0.1 Harmonic Oscillator - Natural units

Hamiltonian is in one dimension given by

$$H = -\frac{\hbar^2}{2m} \frac{\partial^2}{\partial x^2} + \frac{1}{2} m \omega^2 x^2 \quad (\text{B.1})$$

which has corresponding wavefunctions

$$\psi_n(x) = \frac{1}{\sqrt{2^n n!}} \cdot \left( \frac{m\omega}{\pi\hbar} \right)^{1/4} \exp\left(-\frac{m\omega}{2\hbar} x^2\right) H_n\left(\sqrt{\frac{m\omega}{\hbar}} x\right). \quad (\text{B.2})$$

We want to get rid of  $\hbar$  and  $m$  in equation (B.1), and we start with scaling  $H' \equiv H/\hbar$ , such that the Hamiltonian reduces to

$$H' = -\frac{\hbar}{2m} \frac{\partial^2}{\partial x^2} + \frac{1}{2} \frac{m\omega^2}{\hbar} x^2 \quad (\text{B.3})$$

One can now observe that the fraction  $\hbar/m$  comes in both terms, which we can take out by scaling the length  $x' \equiv x \cdot \sqrt{m/\hbar}$ . The final Hamiltonian is

$$H = \frac{1}{2} \frac{\partial^2}{\partial x^2} + \frac{1}{2} \omega^2 x^2 \quad (\text{B.4})$$

which corresponds to setting  $\hbar = m = 1$ . In natural units, one often sets  $\omega = 1$  as well by scaling  $H' = H/\hbar\omega$ , but since we want to keep the  $\omega$ -dependency, we do it slightly different. This means that the exact wavefunctions for the one-particle-one-dimension case is

$$\psi_n(x) = \exp\left(-\frac{\omega}{2} x^2\right) H_n(\sqrt{\omega} x) \quad (\text{B.5})$$

where we take advantage of the Metropolis algorithm and ignore the normalization constant.

### B.0.2 Atomic systems - Atomic units





# Appendix C

## Machine Learning Wave Function

We have seen that the probability of having a set of positions  $\mathbf{X}$  with a set of hidden nodes  $\mathbf{h}$  is given by

$$F(\mathbf{X}, \mathbf{H}) = \frac{1}{Z} e^{-\beta E(\mathbf{X}, \mathbf{H})} \quad (\text{C.1})$$

where we set  $\beta = 1/kT = 1$ ,  $Z$  is the partition function and  $E(\mathbf{X}, \mathbf{h})$  is the system energy

$$E(\mathbf{X}, \mathbf{H}) = \sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2} - \sum_{j=1}^N b_j H_j - \sum_{i,j=1}^{M,N} \frac{X_i W_{ij} H_j}{\sigma_i^2} \quad (\text{C.2})$$

such that

$$F_{\text{RBM}}(\mathbf{X}, \mathbf{h}) = e^{\sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2}} e^{\sum_{j=1}^N \left( b_j h_j + \sum_{i=1}^M \frac{X_i W_{ij} h_j}{\sigma_i^2} \right)}. \quad (\text{C.3})$$

We skip the partition function because it will not affect the results (it is just a normalization constant). The probability of a set of positions only is therefore the sum over all sets of  $\mathbf{h}$ ,  $\{\mathbf{h}\}$ :

$$\begin{aligned} F_{\text{RBM}}(\mathbf{X}) &= \sum_{\{\mathbf{h}\}} e^{\sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2}} \prod_{j=1}^N e^{b_j h_j + \sum_{i=1}^M \frac{X_i W_{ij} h_j}{\sigma_i^2}} \\ &= \sum_{h_1} \sum_{h_2} \dots \sum_{h_N} e^{\sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2}} \times \\ &\quad e^{b_1 h_1 + \sum_{i=1}^M \frac{X_i W_{i1} h_1}{\sigma_i^2}} e^{b_2 h_2 + \sum_{i=1}^M \frac{X_i W_{i2} h_2}{\sigma_i^2}} \dots e^{b_N h_N + \sum_{i=1}^M \frac{X_i W_{iN} h_N}{\sigma_i^2}} \\ &= e^{\sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2}} \prod_{j=1}^N \sum_{h_j=0}^1 e^{b_j h_j + \sum_{i=1}^M \frac{X_i W_{ij} h_j}{\sigma_i^2}} \\ &= e^{\sum_{i=1}^M \frac{(X_i - a_i)^2}{2\sigma_i^2}} \prod_{j=1}^N \left( 1 + e^{b_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}} \right) \end{aligned} \quad (\text{C.4})$$



# Appendix D

## Wave Function Elements

### D.1 Kinetic Energy Calculations

The local energy, defined in equation (4.5), is

$$E_L = \frac{1}{\Psi_T} \hat{H} \Psi_T \quad (\text{D.1})$$

$$= \sum_{k=1}^M \left[ -\frac{1}{2\Psi_T} \nabla_k^2 \Psi_T + U_k + V_k \right]. \quad (\text{D.2})$$

The first term, which is the kinetic energy term, is the only wave function-dependent one. It will in this appendix be evaluated for various wave function elements. From the definition of differentiation of a logarithm, we have that

$$\frac{1}{\Psi_T} \nabla_k \Psi_T = \nabla_k \ln \Psi_T, \quad (\text{D.3})$$

which provides the following useful relation

$$\frac{1}{\Psi_T} \nabla_k^2 \Psi_T = \nabla_k^2 \ln \Psi_T + (\nabla_k \ln \Psi_T)^2. \quad (\text{D.4})$$

Consider a trial wave function,  $\Psi_T$ , consisting of a product of  $p$  wave function elements,  $\{\phi_1, \phi_2 \dots \phi_p\}$ ,

$$\Psi_T = \prod_{i=1}^p \phi_i. \quad (\text{D.5})$$

The kinetic energy related to this trial wave function is then computed by

$$\frac{1}{\Psi_T} \nabla_k^2 \Psi_T = \sum_{i=1}^p \nabla_k^2 \ln \phi_i + \left( \sum_{i=1}^p \nabla_k \ln \phi_i \right)^2, \quad (\text{D.6})$$

which can be computed given all local derivatives  $\nabla_k^2 \ln \phi_i$  and  $\nabla_k \ln \phi_i$ . For each wave function element given below, those local derivatives will be evaluated. In addition, we need to know the derivative of the local energy with respect to the variational parameters in order to update the parameters correctly.

## D.2 Parameter Update

In gradient based optimization methods, as we use, one needs to know the gradient of the local energy with respect to all variational parameters  $\alpha_i$ ,

$$\partial_{\alpha_i} E_L \equiv \frac{\partial E_L(\alpha_i)}{\partial \alpha_i}. \quad (\text{D.7})$$

If we assume that each parameter,  $\alpha$ , only exists in a wave function element,

$$\Psi_T(\alpha) = \phi_1(\alpha) \prod_{i=2}^p \phi_i \quad (\text{D.8})$$

the derivative of the entire local energy is reduced to the derivative of the kinetic energy term given the wave function element,

$$\partial_{\alpha} E_L = \partial_{\alpha} \left( \nabla_k^2 \ln \phi_1(\alpha) + \sum_{i=2}^p \nabla_k^2 \ln \phi_i + \left( \nabla_k \ln \phi_1(\alpha) + \sum_{i=2}^p \nabla_k \ln \phi_i \right)^2 \right) \quad (\text{D.9})$$

$$= \partial_{\alpha} \nabla_k^2 \ln \phi_1(\alpha) + 2 \left( \sum_{i=1}^p \nabla_k \ln \phi_i \right) \cdot \partial_{\alpha} \nabla_k \ln \phi_1(\alpha). \quad (\text{D.10})$$

The sum is already evaluated in the kinetic energy calculations, which means that to calculate the gradients, it is sufficient to calculate  $\partial_{\alpha} \nabla_k \ln \phi_1(\alpha)$  and  $\partial_{\alpha} \nabla_k^2 \ln \phi_1(\alpha)$  for all wave function elements with respect to their variational parameters.

## D.3 Derivatives

### D.3.1 Simple Gaussian

A natural starting point is the Gaussian function, since it appears in all harmonic oscillator calculations. For  $N$  particles in  $D$  dimensions, the function is given by

$$\Psi(\alpha, \mathbf{r}) = \exp \left[ -\frac{1}{2} \alpha \sum_{i=1}^N r_i^2 \right] \quad (\text{D.11})$$

where the derivative with respect to coordinate  $r_k$  is

$$\nabla_k \ln \Psi(\alpha) = -\alpha r_k \quad (\text{D.12})$$

and the second derivative is

$$\nabla_k^2 \ln \Psi(\alpha) = -\alpha N D. \quad (\text{D.13})$$

The gradients for those derivatives are

$$\partial_{\alpha} \nabla_k \ln \Psi(\alpha) = -r_k \quad (\text{D.14})$$

and

$$\partial_{\alpha} \nabla_k^2 \ln \Psi(\alpha) = -N D \quad (\text{D.15})$$

respectively.

### D.3.2 Padé-Jastrow Factor

The Padé-Jastrow factor is introduced in order to take care of the correlations. It is specified in equation (2.15),

$$J(\mathbf{r}; \beta, \gamma) = \exp \left( \sum_{i=1}^N \sum_{j=1}^i \frac{\beta(i, j) r_{ij}}{1 + \gamma r_{ij}} \right) \quad (\text{D.16})$$

, which gives the first and second derivatives

$$\nabla_k \ln J = \sum_{j=1}^k \left( \frac{\beta_{kj}}{(1 + \gamma r_{kj})^2} \right) \quad (\text{D.17})$$

and

$$\nabla_k^2 \ln J = \sum_{j=1}^k \left( \frac{\beta_{kj}}{(1 + \gamma r_{kj})^2} \left[ \frac{1}{r_k} - \frac{\gamma}{1 + \gamma r_{kj}} \right] \right) \quad (\text{D.18})$$

respectively.

The derivative of those again, with respect to the parameters  $\gamma$  and  $\beta_{kl}$  are

$$\partial_\gamma \nabla_k \ln J = -2 \sum_{j=1}^k \left( \frac{\beta_{kj}}{(1 + \gamma r_{kj})^3} \right) \quad (\text{D.19})$$

$$\partial_\gamma \nabla_k^2 \ln J = -2 \sum_{j=1}^k \left( \frac{\beta_{kj}}{(1 + \gamma r_{kj})^3} \left[ \frac{1}{r_{kj}} + \frac{2}{r_k} - \frac{3\gamma}{1 + \gamma r_{kj}} \right] \right) \quad (\text{D.20})$$

$$\partial_{\beta_{kl}} \nabla_k \ln J = \frac{1}{(1 + \gamma r_{kl})^2} \quad (\text{D.21})$$

$$\partial_{\beta_{kl}} \nabla_k^2 \ln J = \frac{2}{(1 + \gamma r_{kl})^2} \left[ \frac{1}{r_l} - \frac{\gamma}{1 + \gamma r_{kl}} \right] \quad (\text{D.22})$$

### D.3.3 Slater Determinant

We need to introduce the Slater Determinant to study more fermions.

### D.3.4 NQS-Gaussian

Now over to the real deal; the machine learning inspired wave function elements.

### D.3.5 NQS-Jastrow Factor

### D.3.6 Hydrogen-Like Orbitals

$$\Psi(\alpha, \mathbf{r}) = \exp \left[ -\frac{1}{2} \alpha \sum_{i=1}^N r_i \right] \quad (\text{D.23})$$

where the derivative with respect to coordinate  $r_k$  is

$$\nabla_k \ln \Psi(\alpha) = -\alpha \quad (\text{D.24})$$

and the second derivative is

$$\nabla_k^2 \ln \Psi(\alpha) = 0 \tag{D.25}$$

The gradients for those derivatives are

$$\partial_\alpha \nabla_k \ln \Psi(\alpha) = -1 \tag{D.26}$$

and

$$\partial_\alpha \nabla_k^2 \ln \Psi(\alpha) = 0 \tag{D.27}$$

respectively.

# Bibliography

- [Dir39] P. A. M. Dirac. “A new notation for quantum mechanics”. In: *Mathematical Proceedings of the Cambridge Philosophical Society* 35.3 (1939), 416–418. DOI: 10.1017/S0305004100021162.
- [Gri05] D.J. Griffiths. *Introduction to quantum mechanics*. 2nd Edition. Pearson PH, 2005. ISBN: 0131911759.
- [Pal05] Josef Paldus. *The beginnings of coupled-cluster theory: An eyewitness account*. Amsterdam: Elsevier, 2005, pp. 115 –147. ISBN: 9780444517197. DOI: 10.1016/B978-044451719-7/50050-0.
- [DCFSI07] T Daniel Crawford and Henry F. Schaefer III. “An Introduction to Coupled Cluster Theory for Computational Chemists”. In: vol. 14. Jan. 2007, pp. 33 –136. ISBN: 9780470125915. DOI: 10.1002/9780470125915.ch2.
- [Deb14] Sukanta Deb. “Variational Monte Carlo technique”. In: *Resonance* 19.8 (2014), pp. 713–739. ISSN: 0973-712X. DOI: 10.1007/s12045-014-0079-x. URL: <https://doi.org/10.1007/s12045-014-0079-x>.
- [Gau] C.F. Gauss. “Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum”. In: ().
- [Leg] A.M. Legendre. “Nouvelles méthodes pour la détermination des orbites des comètes”. In: ().