

CONSTRUCTING HIGH-DIMENSIONAL NEURAL NETWORK POTENTIALS FOR MOLECULAR DYNAMICS

by

John-Anders Stende

THESIS

for the degree of

MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences
University of Oslo

May 2017

Abstract

This is an abstract text.

To someone

This is a dedication to my cat.

Acknowledgements

I acknowledge my acknowledgements.

Contents

1	Introduction	1
1.1	Goals	1
I	Theory	3
2	Molecular dynamics	7
2.1	Potential energy surfaces	7
2.2	LAMMPS	9
2.2.1	Installing LAMMPS	9
2.2.2	LAMMPS input script	10
2.2.3	LAMMPS structure	14
2.2.4	Extending LAMMPS	14
3	Machine learning	15
3.1	Artificial Neural Networks	16
3.2	Training	19
3.2.1	Backpropagation	19
4	Neural networks in molecular dynamics	25
4.1	High-dimensional NNPs	26
4.1.1	Symmetry functions	28
4.1.2	Symmetry functions and forces	33
4.2	TensorFlow	34
4.2.1	Convolution	34
II	Implementation and validation	35
4.3	Time usage	37
4.4	Training Lennard-Jones potential	39
4.4.1	Many-neighbour Lennard-Jones	41
.1	Appendix	47
.1.1	Symmetry functions derivatives	47

Chapter 1

Introduction

Motivate the reader, outline structure of report and what we have done

1.1 Goals

The main goal of this thesis is to use artificial neural networks to construct many-body potential energy surfaces (PES) to be used in molecular dynamics (MD) simulations. This goal can be split into the following intermediate objectives:

- Train an one-dimensional ANN with the machine learning library TensorFlow (TF) to reproduce the shifted Lennard-Jones (LJ) potential. As a first test case, we will train the NN on random data and investigate how different NN architectures and activation functions influence the quality of the potential fit.
- Train a many-body NN to reproduce the Stillinger-Weber (SW) potential for Si. Training data (configurations of particles and energies) will be sampled from SW LAMMPS simulations. The resulting NNP will then be used to simulate Si with LAMMPS. We therefore need to extend the LAMMPS library with our own neural network potential.
- Construct a many-body NNP that will be able to simulate a system consisting of several atom types. The method will be tested on the Vashishta potential for SiO₂.

Part I

Theory

Theory needed to understand the results and implementations

Chapter 2

Molecular dynamics

Molecular dynamics (MD) is a method to simulate the physical movements of atoms and molecules in gases, liquids and solids. It is thus a type of N-body simulation. The atoms are modelled as point-like particles with interactions described by classical force fields. Their time-evolution is governed by Newton's equations of motion. MD allows one to study the microscopic movement of thousands or millions of atoms, enabling the sampling of macroscopic properties such as temperature, pressure, diffusion, heat capacity and so on.

The dynamics of an ensemble of particles is governed by their interactions. In classical MD, the interactions are described by a classical force field \mathbf{F} , which is defined as the negative gradient of a potential energy surface (PES) E ,

$$\mathbf{F} = -\nabla E \quad (2.1)$$

The PES is represented by a mathematical function, often split into a sum over two- and three-body interactions. Instead of using an analytical functional form to represent the energy, one can calculate the energies and forces for each time step in the MD simulation using an ab initio method. Different quantum mechanical methods like Hartree-Fock (HF) and Density Functional Theory (DFT) are used for this purpose. These approaches to MD are called quantum-classical molecular dynamics (QCMD).

2.1 Potential energy surfaces

The dynamics of an ensemble of particles is governed by their interactions. In classical MD, the interactions are described by a classical force field \mathbf{F} , which is defined as the negative gradient of a potential energy surface (PES) E ,

$$\mathbf{F} = -\nabla E \quad (2.2)$$

The expression for $V(r)$ may be obtained in different ways:

- Empirical potentials: the functional form are chosen and the parameters fitted from empirical knowledge about the system
- Semi-empirical potentials: based on methods from computational QM, but many approximations are made and some parameters are obtained from empirical data
- Quantum mechanical potentials: a computational QM method is used to obtain the energy for different configurations, and then fitted by a functional form

Instead of using a pre-defined functional form to represent the energy, one can calculate the energies and forces for each time step in the MD simulation using an ab initio method. Different quantum mechanical methods like Hartree-Fock (HF) and Density Functional Theory (DFT) are used for this purpose. These approaches to MD are called quantum-classical molecular dynamics (QCMD).

The PES represents the potential energy $V(\mathbf{r})$ of a system as a mathematical function of the configuration of the atomic positions $\mathbf{r} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N)$. To reduce the N -body problem from a fully quantum description to a classical potential we must make the Born-Oppenheimer approximation, which is the assumption that the motion of atomic nuclei and electrons in a molecule can be separated. This enables us to freeze the degrees of freedom of the nuclei and solve the electronic Schrödinger equation with the positions of the nuclei as parameters. Varying the positions of the nuclei in small steps and repeatedly solving the Schrödinger equation, the quantum PES is made. Classically, we can now approximate the atomic nuclei as point-particles that follow Newtonian dynamics, while the positions and velocities of the electrons are baked into the PES, usually representing the ground state. (HA ET EGET AVNSITT OM BORN-OPPENHEIMER I KVANTETEORIDELN KANSKJE?)

As mentioned above, we represent the PES by a pre-defined functional form in classical MD. We can assume that this function can be written as a sum of n -body terms,

$$V(\mathbf{r}) \approx \sum_{i=1}^N V_1(\mathbf{r}_i) + \sum_{i<j}^N V_2(\mathbf{r}_i, \mathbf{r}_j) + \sum_{i<j<k}^N V_3(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k) + \dots \quad (2.3)$$

Each term can be determined by performing an electronic structure calculation for a number of configurations, and fitted by a suitable functional form. Determining how many terms that is sufficient to describe a system adequately is however not trivial. Noble gases like argon have weak interactions and may be well described by a two-body potential V_2 , while molecules with strong angular dependencies on their bond, like silicon (Si), will need at least three-body potentials V_3 . Another challenge is to identify the configuration space of the system. For three-body potentials and more the set of probable configurations is large,

and running an *ab initio* simulation for each configuration is expensive. This problem will be discussed below.

An important step is how to fit the above potential terms obtained from QM calculations to a functional form. There exist a lot of different interpolation techniques; in this theses we will use ANNs to fit the data. (SKRIVE OM ANDRE INTERPOLASJONSTEKNIKKER?)

2.2 LAMMPS

(HAVE TO REFER TO THE MANUAL HERE OR SOMETHING. DOES THAT COUNT AS AN ORDINARY REFERENCE, OR SHOULD IT BE A FOOTNOTE?) LAMMPS (Large-scale Atomic/Molecular Massively Parallel Simulator) is a classical molecular dynamics simulation package developed at Sandia National Laboratories. It is written in highly portable and extendable objected-oriented c++. The package's functionality encompasses a wide variety of potentials, atom types, integrators, thermostats and ensembles and are able to simulate a large number of different systems, including atoms, molecules, coarse-grained particles, polymers, graunular materials and metals. LAMMPS can be run on a single processor or in parallell. There exist several other molecular dynamics packages, like GROMACS, OpenMD, Aber and NAMD that could also have served the purpose of this thesis. We have chosen to work with LAMMPS because it is well documented, easy to expand upon and because the Computational physics group at UiO has a lot of experience with using it for molecular dynamics simulations. In this section we will outline the basic usage of LAMMPS, including a walkthrough of an example input script. We also describe its class hierarchy and how we can add our own functionality to the package.

2.2.1 Installing LAMMPS

We have chosen to install LAMMPS by cloning the Github repository (LINK) and compile the source code by running *make* in the source directory. To compile a serial version of the software, we simply run the command

```
make serial
```

while the corresponding command for the parallel MPI version is

```
make mpi
```

New updates can be downloaded by setting an upstream git remote to point to the LAMMPS GitHub repository and doing a *gitpullupstreammaster*. (TROR IKKE DETTE ER HELT RIKTIG). LAMMPS also have several additional packages that can be installed. This can be done by running

```
make yes-<package name>
```

e.g. *manybody* to install many-body potentials like Stillinger-Weber and Vashishta.

2.2.2 LAMMPS input script

LAMMPS is run by providing an input script as argument to the executable. This input script is read line by line and has its own syntax. A good way to show the basics an input file is to look at a simple example of a script to measure diffusion in liquid argon (CHANGE LC VALUE LATER):

```
#### initialization ####
units      metal
dimension  3
boundary   p p p
atom_style atomic

#### create geometry and atoms ####
lattice    fcc 1.08506
variable   Nc equal 10
region     myRegion block 0 ${Nc} 0 ${Nc} 0 ${Nc}
create_box 1 myRegion
create_atoms 1 box

#### set mass and initial temperature ####
mass       1 1.0
variable   temp equal 300
velocity    all create ${temp} 87287 mom yes

#### compute diffusion ####
compute     displacement all displace/atom

#### potential ####
pair_style  lj/cut 2.5
pair_coeff   1 1 1.0 1.0 2.5
neighbor    0.5 bin
neigh_modify every 20 delay 0 check no

#### integration ####
timestep    0.01
run_        style verlet
fix         integration all nve

#### output ####
thermo      50
thermo_style custom step temp density press ke pe etotal
thermo_modify norm yes
```

```
#### run simulation ####  
run      5000  
dump diffusion all custom 100 diffusion*.dat c_displacement[4]  
run_     10000
```

We will in the following briefly explain what processes that are evoked in LAMMPS when the above commands are read.

```
units metal
```

This command defines the units that are used in the simulation and the output. LAMMPS have eight different unit sets. The *metal* set measures distance in ngstrm, energy in eV and temperature in Kelvin. The choice of units depend on the system that is investigated and the scale we are looking at.

```
boundary p p p
```

We want to measure diffusion in a bulk Argon liquid, thus we want to have periodic boundary conditions in all three dimensions, annotated by a *p*. LAMMPS can also handle stiff (*f*) and adaptive (*s*) non-periodic boundaries. Adaptive means that the position of the face is set so as to encompass the atoms in that dimension.

```
atom_style atomic
```

Different systems need different information to be stored for each atom. For style *atomic*, only the default attributes are associated with each atom, namely coordinates, velocities, atom IDs and types. This is sufficient for pair-interacting, non-bonded systems like Argon.

```
lattice    fcc 5.720
```

Next, we set the initial configuration of the atoms. Here, a fcc lattice with a lattice constant of 5.720 is used. Other types of lattices based on cubic or squared unit cells are also available.

```
variable   Nc equal 10  
region     simBox block 0 ${Nc} 0 ${Nc} 0 ${Nc}
```

The system's geometry is defined with the *region* command. The *block* style is simply a 3-dimensional straight-faced box with a size of $N_c = 10$ unit cells in each dimension. We have labelled the region *simBox*. The number of unit cells N_c is defined as a LAMMPS *variable*. Variables can be referenced elsewhere in the

script by writing *\$variable* to become part of a new input command like above. LAMMPS enables many styles of variables to be defined.

```
create_box 1 simBox
create_atoms 1 box
```

The command *create_box* creates a simulation box based on the specified region, in our case the *block* region defined above. The argument specifies the number of atom types that will be used in the simulation. Next, *create_atoms* with the argument *box* fills the simulation domain with atoms of type 1 on the lattice. LAMMPS also lets you create a random collection of atoms or single atoms at specified coordinates.

```
mass 1 28
variable temp equal 300
velocity all create ${temp} 87287
```

We need to assign mass to the atoms. For metal units, mass is measured in grams/mole. The atoms are also given an initial velocity that corresponds to the given initial temperature.

```
### potential ###
pair_style lj/cut 2.5
pair_coeff 1 1 1.0 1.0 2.5
```

The choice of potential is made with the *pair_style* command. We want to simulate interactions using the cutoff Lennard-Jones potential with a global cutoff of 2.5. The LJ parameters are set with *pair_coeff*, the way this is done depends on the potential being used. For our choice of potential the arguments is as follows: *atom - type - 1 atom - type - 2 sigma epsilon cutoff*. The parameters and cutoff can thus be set separately for each pair of atom types.

```
neighbor 0.5 bin
neighbor_modify every 20 check yes
```

These commands sets parameters that affect the building of neighbor lists. The first argument to *neighbor* is the skin size, while the second selects what algorithm is used to build the lists. The *bin* style creates the lists by binning, which in most cases (including ours) is the fastest method. Further, we can control how often the lists are built with *neighbor_modify*. The above arguments specifies that new neighbour lists are built every 20 steps, but every step LAMMPS checks if any atom has moved more than half the skin distance. If that is the case, new lists are built.

```
timestep 0.01
run_style verlet
fix      integration all nve
```

LAMMPS integrates Newtons' equations of motion with the velocity-Verlet algorithm by default. This is the integrator of choice for most MD applications due to its simplicity and symplectic nature (REF TO SECTION). The rRESPA integrator [11] scheme is also available. LAMMPS does however not integrate and update the positions and velocities of the particles if not explicitly told so. This is done with a *fix*, which is any operation that is applied to the system during timestepping. The above *fix* tells LAMMPS to integrate all atoms in the system so that they follow trajectories consistent with the microcanonical ensemble.

```
thermo 50
thermo_style custom step temp press ke pe etotal
thermo_modify norm yes
```

We can control what thermodynamic properties to calculate and output with *thermo_style*, while *thermo* decides how often they should be computed. We want to output the time step, temperature, pressure, kinetic energy, potential energy and total energy. Also, we want to normalize the extensive quantities (the energies) by the number of atoms.

```
compute displacement all displace/atom
dump diffusion all custom 100 diffusion*.dat
c_displacement[4]
```

To measure diffusion in the Argon liquid, we need to calculate the net displacement of all atoms. This can be done with a *compute*, which defines a computation that is performed on a group of atoms, in this case all atoms. The *displace/atom* compute calculates the current displacement of each atom from its original coordinates, including all effects due to atoms passing through periodic boundaries. For the compute to actually be performed, it needs to be evoked by other LAMMPS commands, like *dump*, which writes a snapshot of atom quantities to one or more files every N timesteps. Computes are referenced via the following notation $c_I D$, where ID is the ID of the compute.

```
run 5000
```

Lastly, we run the simulation for 5000 timesteps. LAMMPS allows several run commands to be issued after one another. This comes in handy if we want to thermalize the system before measuring diffusion. This can be achieved by writing the following commands:

```
run 5000
dump diffusion all custom 100 diffusion*.dat
      c_displacement[4]
run 10000
```

When parsing the second run command, LAMMPS will continue the dynamics while computing eventual new fixes, computes, dumps etc. defined since the last run command.

This section has illuminated some of the basic functionality of LAMMPS through several comments on a simple input script. Now we move on to how class hierarchy of LAMMPS is structured.

2.2.3 LAMMPS structure

LAMMPS is written in C++ in an object-oriented fashion. The class hierarchy is shown schematically in figure (MAKE MY OWN FIGURE IN INKSCAPE). The blue classes are the core classes, these are visible anywhere in LAMMPS. Most of these have subclasses, of which the red ones are called style classes. We recognize many of the input script commands among the style classes; a rule of thumb is that every input script command has a corresponding class and a corresponding file name in the source directory. For instance, the cutoff Lennard-Jones potential used in our input script is a subclass of Pair, and is evoked by the command *pair_style lj/cut* and the source file is named *pair_lj_cut*.

We will not go into detail about the functionality and communication of these classes. The class that is the most relevant for this work is the Pair class, as we have implemented our own NN potential *pair_style*. In the following section we will describe how LAMMPS can be extended with a new pair potential through example.

2.2.4 Extending LAMMPS

To extend LAMMPS with a new potential, we simply add a new *.cpp and *.h file to the source directory and re-compile. The new class will (in theory) work with all other LAMMPS classes without breaking any functionality. A good strategy is to start with an existing potential source file that is somewhat similar to the one we want to make, instead of programming from scratch. We aim to construct a many-body potential, and have chosen to base our implementation on the Stillinger-Weber potential file, *pair_sw.cpp* and *-.h*.

Chapter 3

Machine learning

Machine learning is the science of giving computers the ability to learn without being explicitly programmed. The idea is that there exist generic algorithms which can be used to find patterns in a broad class of data sets without having to write code specifically for each problem. The algorithm will build its own logic based on the data.

There are two main machine learning methods, supervised learning and unsupervised learning. For the former method, you know the answer to a problem and let the computer deduce the logic behind it. When doing the latter, you do not know the answer and are just trying to find patterns and relationship in the data.

Another way to categorize machine learning tasks is to consider the desired output of a system. Some of the most common tasks are:

- **Classification:** Outputs are divided into two or more classes. The goal is to produce a model that assigns inputs into one of these classes. An example is to identify digits based on pictures of hand-written ones. Typically supervised.
- **Regression:** Find a functional relationship between variables, the outputs are thus continuous.
- **Clustering:** Inputs are divided into groups without knowing the different groups, i.e. a form of unsupervised classification.

In this thesis we will use machine learning to do regression, i.e. to find the PES of interacting molecules. One approach to machine learning regression is artificial neural networks.

3.1 Artificial Neural Networks

An Artificial Neural Network (ANN) consists of layers of connected nodes, or neurons. It is supposed to mimic a biological nervous system by letting each neuron interact with other neurons by sending signals between layers. The accumulated signals received by a neuron is used as input to its *activation function*, which is used to send new signals to following layers. This function is supposed to mimic that a biological neuron is activated if the input is above a certain threshold. We will here follow Behler [1] and Dragly [2].

There are many types of ANNs made for different purposes. The most widely used are feed-forward and recurrent ANNs. The feedforward network was the first and arguably the most simple type of ANN. In this network the information or signals only move forward through the layers - from the input layer, through each successive hidden layer and finally to the output layer. On the other hand, recurrent networks have cycles and loops, i.e. bi-directional data flow. Feed-forward neural networks, also called Multilayer Perceptrons (MLP) were shown to be universal function approximators by K. Hornik et al. [3] in 1989. We will use these networks to approximate the PES in thesis.

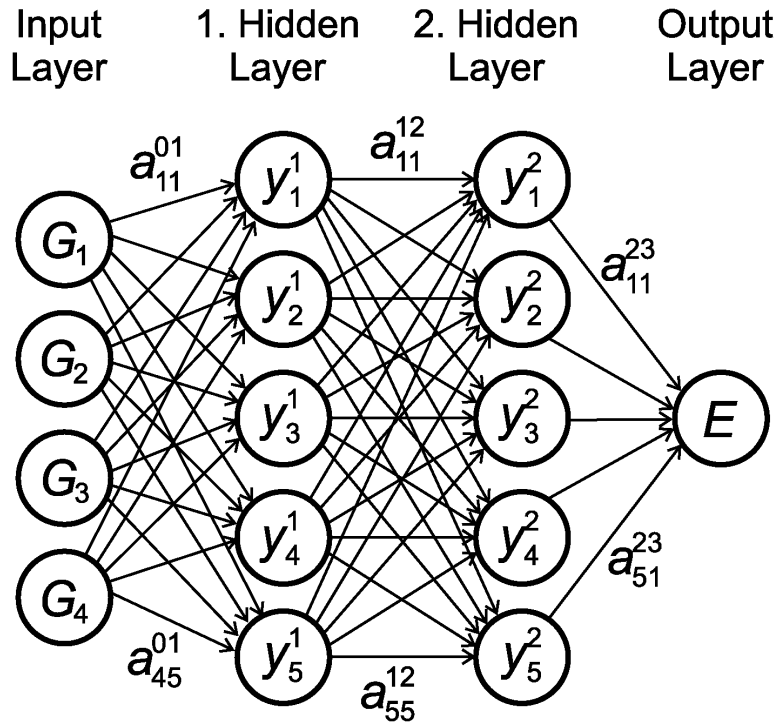


Figure 3.1: Neural network

An MLP networks is shown schematically in Figure 3.1 (STOLEN FROM [3], EDIT TO ONE INPUT LATER).

This network maps input coordinates G_i to a potential energy in the output layer via two hidden layers with an (arbitrary) number of hidden nodes. The hidden neurons have no physical meaning, but have the purpose of defining the functional form of the NN. All nodes in each layer are connected to the nodes in the adjacent layer by *weights*, the fitting parameters of the NN. These are called fully connected layers. We use the notation a_{ij}^{kl} for the weight connecting node number i in layer k with node number k in layer $l = k + 1$. Also, each node in the hidden layers and the output node have a *bias* b_i^j , where i indicates node number and j layer. All weights and biases are real-valued numbers.

The output of this MLP is calculated in the following way. The configuration are specific by the input nodes. Then, a weighted sum x_m^1 of the input coordinates G_i is calculated for each node m in the first hidden layer 1

$$x_m^1 = \sum_{i=1}^4 G_i a_{i,m}^{01} + b_m^1 \quad (3.1)$$

where we also add the bias b_m^1 for each node. This is a linear combination of the input coordinates using the connection weights as coordinates, which is used as input to the activation function f_m^1 of each neuron, thus producing the output y_m^1 of all neurons in layer 1:

$$y_m^1 = f_m^1(x_m^1) \quad (3.2)$$

For an arbitrary node m in a hidden layer n this generalizes to

$$y_m^n = f_m^n(x_m^n) = f_m^n \left(\sum_{i=1}^{N_{n-1}} y_i^{n-1} a_{i,m}^{n-1,n} \right) \quad (3.3)$$

where N_n is the number of nodes in layer n . When the output of all the nodes in the first hidden layer are computed, the values of the subsequent layer can be calculated and so forth until the output is obtained. The complete functional form of an MPL with n hidden layers is thus

$$E = f_1^{n+1} \left(\sum_m^{N_n} f_m^n \left(\sum_l^{N_{n-1}} f_l^{n-1} \left(\cdots \sum_j^{N_1} f_j^1 \left(\sum_{i=1}^{N_0} a_{i,j}^{0,1} G_i + b_j^1 \right) \cdots \right) \right) \right) \quad (3.4)$$

i.e. a nested sum of activation functions. A desired property of activation functions is to converge asymptotically to constant output values for large and small arguments, with a non-linear region in-between. The non-linearity enables the ANN to fit arbitrary functions [3]. Common activation functions [4] for function approximation are the sigmoid function

$$f(x) = \frac{1}{1 + e^{-x}} \quad (3.5)$$

and the hyperbolic tangent

$$f(x) = \tanh(x) \quad (3.6)$$

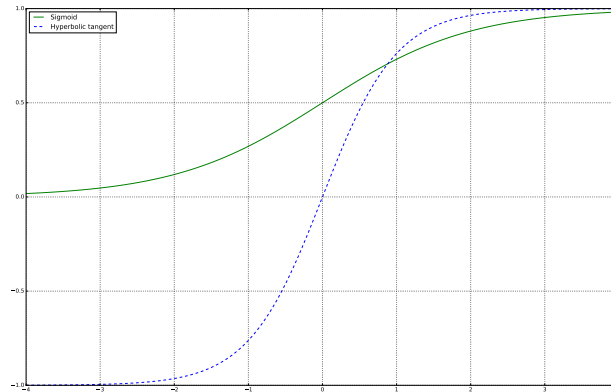


Figure 3.2: Activation functions

Both these functions (Figure 3.2) have the desired properties mentioned above: They are constant (1 or 0) for all values except a small interval around zero, where they have non-linear behaviour. This keeps the outputs of the neurons from blowing up. The sigmoid are more biologically plausible because the output of inactive neurons are zero. Such activation function are called *one-sided*. However, it has been shown [5] that the hyperbolic tangent performs better than the sigmoid for multi-layer networks. Nevertheless, one should assess the problem at hand when deciding what activation function to use; the performance can vary from problem to problem. (VENTE MED DISKUTERE DETTE TIL JEG HAR SNAKKET OM BACKPROPAGATION OSV.) SNAKKE OM SATURATION OSV.

In later years, the rectifier function (Figure ??)

$$f(x) = \max(0, x) \quad (3.7)$$

has become the most popular for deep neural networks [6]. PLOT DENNE FUNKSJONEN. The two former functions are symmetrical around zero osv.... It have been argued to be even more biologically plausible than the sigmoid and also perform better than the hyperbolic tangent for deep NNs. [7].

A linear activation function $f(x) = x$ is often used in the output layer to avoid any constraint in the range of output values. For regression, the most widely used output activation is the unity function $f(x) = 1$, i.e. no activation (HVORFOR?).

As shown in (3.4), the NN consists of several nested terms of the form

$$h(x) = c_1 f(c_2 x + c_3) + c_4 \quad (3.8)$$

where the parameters c_i are weights and biases. By adjusting these parameters, the activation functions can be shifted up and down or left and right, change slope or be rescaled, thus making the NN more flexible.

3.2 Training

The ability of the NN to represent a function accurately depends on the values of the weights. The weights are typically initialized as random numbers, and then optimized iteratively to minimize the error to a set of expected output values. The expected output values are in our case known energies of a set of atomic configurations, which have been obtained in quantum mechanical calculations of electronic structure. This optimization is called *training* and an iteration is called an *epoch*. The fitting process corresponds to the minimization of an error function, also known as the quadratic *cost* function

$$\Gamma = \frac{1}{2N} \sum_{i=1}^N (E_{i,ref} - E_{i,NN})^2 \quad (3.9)$$

$E_{i,ref}$ are the known target energies, while $E_{i,NN}$ are the energies produced by the network. The constant $1/2$ are there to cancel out the exponent when we differentiate later. Other cost functions are used for other ML tasks (SKRIVE OM DISSE. KANSKJE HELE DENNE SEKSJONEN BR SKRIVES MER GENERELL OG IKKE KUN FOR TILFELLET PES?)

3.2.1 Backpropagation

From: <https://www.willamette.edu/gorr/classes/cs449/backprop.html> and <http://neuralnetworksanddeeplearning.com/chap2.html>. Mainly from the first one, but argument for error of neuron from second.

There are a large number of algorithms that can be used to determine the set of weights minimizing the cost function (3.9). Different kinds of gradient descent methods are widely used, but to implement these we need the partial derivatives of the cost function w.r.t. to all weights and biases to determine how we should adjust each weight. The idea behind gradient descent methods is to minimize a function by iteratively taking steps in the direction of the negative of its gradient w.r.t. all its parameters, i.e. the direction of steepest descent towards a minima:

$$w_{ij}^{new} = w_{ij}^{old} - \gamma \frac{\partial \Gamma}{\partial w_{ij}^{old}} \quad (3.10)$$

where we have skipped the layer indicies for clarity. Also, w_{ij} now denotes the weight connecting node j in layer k with node i in the following layer. This change of order is introduced to make it easier to rewrite the below equations for

backpropagation into matrix-vector equations. As seen from (3.10), we need to adjust each weight by an amount

$$\Delta w_{ij} = \frac{\partial \Gamma}{\partial w_{ij}} \quad (3.11)$$

The corresponding expression for the biases is

$$\Delta b_i = \frac{\partial \Gamma}{\partial b_i} \quad (3.12)$$

A common method to obtain these derivatives is backpropagation [8]. In backpropagation, the output of the NN is compared to the desired output, the error is then propagated backwards to adjust the weights. The method is essentially an implementation of the chain rule, and will allow us to calculate the partial derivatives of the cost with respect to all the weights, thereby obtaining the gradient of the network.

To change the value of the cost function (??), we need to change the outputs of the neurons in the network. Changing the input to neuron j by a small amount Δx_j results in the output

$$y_j = f_j(x_j + \Delta x_j) \quad (3.13)$$

This change will propagate through later layers in the network, finally causing the overall cost to change by an amount $\frac{\partial \Gamma}{\partial x_j} \Delta x_j$. If $\partial \Gamma / \partial x_j$ is close to zero, then we can't improve the cost much by perturbing the weighted input x_j ; the neuron is already quite near the optimal value. This is a heuristic argument for $\partial \Gamma / \partial x_j$ to be a measure of the error of the neuron:

$$\delta_j \equiv \frac{\partial \Gamma}{\partial x_j} \quad (3.14)$$

We also define

$$A_i = \{j : w_{ij}\} \quad (3.15)$$

as the set $\{j\}$ of nodes anterior to node i and connected to node i with weights w_{ij} , in addition to

$$P_j = \{i : w_{ij}\} \quad (3.16)$$

as the set $\{i\}$ of nodes posterior to node j and connected to node j with weights w_{ij} .

The weight change (3.11) can be expanded into two factors by use of the chain rule:

$$\Delta w_{ij} = \frac{\partial \Gamma}{\partial x_i} \frac{\partial x_i}{\partial w_{ij}} \quad (3.17)$$

Now we move in the opposite direction compared to the feed-forward function composition stage: First we differentiate the cost w.r.t neuron i 's input, then we

differentiate the input w.r.t. weight w_{ij} connecting neurons j and i . The first term on the right is the error δ_i of node i . The second is

$$\frac{\partial x_i}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_{m \in A_i} w_{im} y_m = y_j \quad (3.18)$$

Putting the two together, we get

$$\Delta w_{ij} = \delta_i y_j \quad (3.19)$$

To compute this quantity, we thus need to know the output and the error for all nodes in the network. The outputs are as we have seen

$$y_i = f_i(x_i) = f_i \left(\sum_{j \in A_i} w_{ij} y_j + b_i \right) \quad (3.20)$$

The error for output neuron o is

$$\delta_o = \frac{\partial \Gamma}{\partial x_o} = \frac{\partial \Gamma}{\partial y_o} \frac{\partial y_o}{\partial x_o} = (d_o - f_o(x_o)) \frac{\partial f_o(x_o)}{\partial x_o} \quad (3.21)$$

When doing regression our output neurons will have the unity activation function, which means that the error is reduced to

$$\delta_o = d_o - y_o \quad (3.22)$$

This error is then propagated backwards through the network. Each hidden neuron will have an error

$$\delta_j = \frac{\partial \Gamma}{\partial x_j} = \sum_{i \in P_j} \frac{\partial \Gamma}{\partial x_i} \frac{\partial x_i}{\partial y_j} \frac{\partial y_j}{\partial x_j} \quad (3.23)$$

The first term on the right is the error of node i . The second is the derivative of the net input of all posterior nodes $\{i\}$ w.r.t. the output of node j :

$$\frac{\partial x_i}{\partial y_j} = \frac{\partial}{\partial y_j} \left(\sum_{m \in A_i} w_{im} y_m + b_i \right) = w_{ij} \quad (3.24)$$

while the third is the derivative of node j 's activation function w.r.t. its net input:

$$\frac{\partial y_j}{\partial x_j} = \frac{\partial f_j(x_j)}{\partial x_j} \equiv f'_j(x_j) \quad (3.25)$$

Putting all the pieces together we obtain

$$\delta_j = f'_j(x_j) \sum_{i \in P_j} \delta_i w_{ij} \quad (3.26)$$

The above expression requires that we know the errors of all the posterior nodes of node j . As long as there are no cycles in the network, there is an ordering of nodes from the output back to the input that respects this condition. It is therefore valid only for feed-forward NNs. The errors are propagated backwards through the whole NN until we reach the input nodes. By propagating the error of only one output neuron, we thus obtain the errors of all the neurons at once. This is the strength of the backpropagation algorithm. For the biases, we have

$$\Delta b_j = \frac{\partial \Gamma}{\partial b_j} = \sum_{i \in P_j} \frac{\partial \Gamma}{\partial x_i} \frac{\partial x_i}{\partial y_j} \frac{\partial y_j}{\partial x_j} \frac{\partial x_j}{\partial b_j} \quad (3.27)$$

The only new term here is the last one:

$$\frac{\partial x_j}{\partial b_j} = \frac{\partial}{\partial b_j} \left(\sum_{m \in A_j} w_{jm} y_m + b_j \right) = 1 \quad (3.28)$$

Consequently, the update rule for the biases is simply the error of each neuron:

$$\Delta b_j = \delta_j \quad (3.29)$$

For fully-connected FFNNs we can easily rewrite these equations in matrix notation. In this notation, the biases, inputs, outputs and errors for all nodes in a layer are combined into vectors, while all weights from one layer to the next form a matrix W_l , where l denotes the layer number. The inputs to all the nodes in layer l can be written

$$\begin{pmatrix} x_1^l \\ x_2^l \\ x_3^l \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{pmatrix} \begin{pmatrix} y_1^k \\ y_2^k \\ y_3^k \end{pmatrix} + \begin{pmatrix} b_1^l \\ b_2^l \\ b_3^l \end{pmatrix}$$

so that equation (3.20) takes the form

$$\mathbf{y}_l = f_l(W_l \mathbf{y}_{l-1} + \mathbf{b}_l) \quad (3.30)$$

Comparing equations (3.20) and (3.26), we observe that

$$\sum_{j \in A_i} w_{ij} = \left(\sum_{i \in A_j} w_{ij} \right)^T \quad (3.31)$$

i.e. the weight matrices used in backpropagation are the transpose of the matrices used in forward activation. The complete backpropagation algorithm including forward activation then looks as follows:

1. Initialize the input layer:

$$\mathbf{y}_0 = \mathbf{x}_0 \quad (3.32)$$

2. Propagate the activity forward: for $l = 1, 2, \dots, L$:

$$\mathbf{y}_l = f_l(W_l \mathbf{y}_{l-1} + \mathbf{b}_l) \quad (3.33)$$

3. Calculate the error in the output layer:

$$\boldsymbol{\delta}_L = \mathbf{t} - \mathbf{y}_L \quad (3.34)$$

4. Calculate the error: for $l = L - 1, L - 2, \dots, 1$:

$$\boldsymbol{\delta}_l = (W_{l+1}^T \boldsymbol{\delta}_{l+1}) \cdot f'_l(\mathbf{x}_l) \quad (3.35)$$

5. Update the weights and biases:

$$\Delta W_l = \boldsymbol{\delta}_l \mathbf{y}_{l-1}^T, \quad \Delta \mathbf{b}_l = \boldsymbol{\delta}_l \quad (3.36)$$

Chapter 4

Neural networks in molecular dynamics

Following [9] and [10]. All the dynamics of a molecular system is determined by the multidimensional potential-energy surface (PES), which in general is a real-valued function depending on all atomic coordinates in the system. Using the Born-Oppenheimer approximation (FORKLARE ELLER REFERERE), various quantum mechanical methods like Hartree-Fock and density functional theory (DFT) are available to directly calculate the PES and forces for a given configuration. Calculating the energies for all relevant configurations of a system is demanding and time-consuming, and only a limited number of energies can be computed and stored. Consequentially, during a *ab initio* molecular dynamics (MD) simulations, the energies and forces for all configurations visited are not available beforehand. In *ab initio* MD, energies and forces are therefore calculated on-the-fly, typically using DFT. This is however a very inefficient approach. A more efficient method is to construct an analytical PES to use in simulations, but this can only be done for very simple systems.

A solution to this problem is the introduction of approximate PESs, and there are two fundamental approaches. The most widely used and conventional method is to replace the solution of the Schrödinger equation by a simplified energy expression based on physical considerations and reasonable approximations. These functions are not obtained by *ab initio* methods, but can in many cases be sufficiently accurate to be used in simulations (BR HA EN EGEN SEKSJON OM DE VANLIGSTE POTENSIALENE).

An alternative approach is to employ machine learning (ML) potentials that have no direct physical meaning. The aim of these purely mathematical functions is to fit an analytic expression to a set of reference data obtained by quantum mechanical calculations. There are various methods (REFERENCES). ML potentials are particularly useful

- if long MD simulations are required

- if many MD trajectories are needed
- if the systems are too large for the application of QM methods

Physical and ML potentials thus enables us to extend the time and length scales of MD simulations beyond the realm of *ab initio* methods.

Feed-forward neural networks (FFNN) have been demonstrated to be useful for the construction of PESs because they are universal function approximators. There is no restriction in the accuracy that can be achieved when constructing neural network potentials (NNP). Unlike physical potentials, they are not restricted by any *ad hoc* functional form. NNPs offer a number of advantages for the construction of PESs:

- Energies can be fitted to very high accuracy, leaving only the underlying error of the reference data
- Evaluation of NNPs require much less CPU time than QM methods
- The NNP expression is unbiased and generally applicable to all types of bonding

Still, there are disadvantages one should be aware of:

- The evaluation of NNPs are notably slower compared to the use of classical force fields
- NNP expressions have no physical interpretation, and have very limited extrapolation capabilities
- NNPs are, for the time being, only applicable to systems containing only a few different chemical elements (but many atoms)

4.1 High-dimensional NNPs

The use of a single FFNN to represent complicated systems containing many atoms of different types is not possible for several reasons. First of all, the degrees of freedom (the number of weight parameters) will be very many, making the training and evaluation of the network slow. Secondly, the NN has a symmetry problem. It does not take into account that exchanging two or more atoms can lead to an energetically equivalent configuration. For example, exchanging the positions of both hydrogen atoms in a water molecule will not alter the configuration energy. However, noting that all weights have numerically different values, changing the order of the input coordinates to the NN will result in a different energy value. This problem can be solved by a different choice of input coordinates, discussed below. The third problem in using a single FFNN is that

it is only applicable to the system size that has been used during the training. If the NN has been trained with 5 inputs, it can not be used for a system containing any other number of atoms.

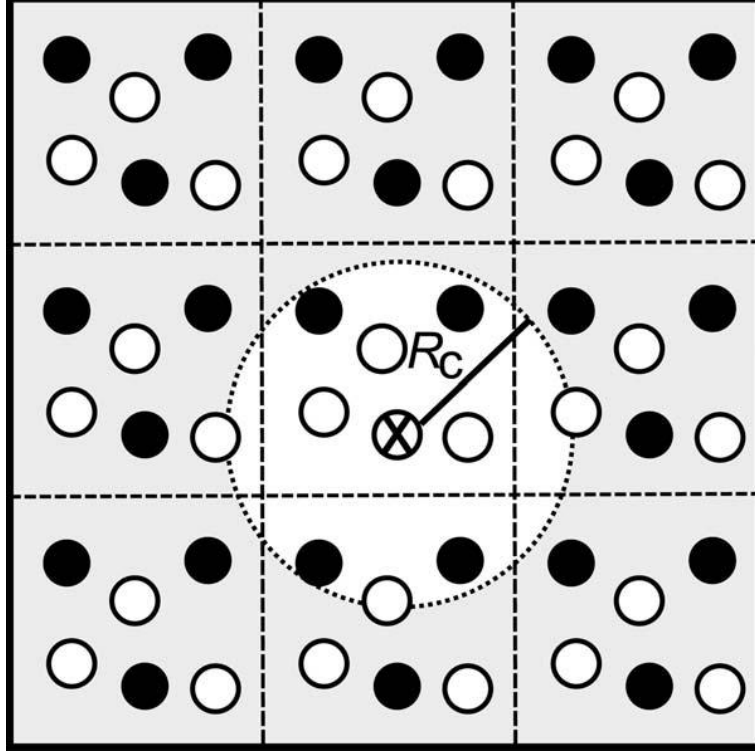


Figure 4.1

A solution to all these problems is to construct the system energy E_s as a sum of N atomic energy contributions E_i , which are provided by a set of individual atomic NNs,

$$E_s = \sum_{i=1}^N E_i \quad (4.1)$$

The atomic energies E_i depend on the local chemical environment up to a cutoff radius R_c as shown in figure 4.1 (STJLET FRA PAPER, M LAGE EGEN?). This is analogous to MD simulations with neighbour lists, thus making it possible to obtain the energies and forces of an atom and all its neighbours simultaneously. The positions of the neighbouring atoms within the cutoff sphere are described by a set of many-body symmetry functions discussed below. As the order of the summation does not change the total energy E_s , the symmetry problem above is solved. Each chemical element have separate types of NNs with their own architecture and weights, but for a given element all atomic NNs are equal. Such

a system of NNs are also applicable to different system sizes: To add an atom of a given type, we simply extend the set of NNs by another network identical to the other NNs of the same type, and if an atom is removed, we delete the respective NN.

- The introduction of the cutoff R_c reduces the effective dimensionality of the problem, which allows to use NNs of tractable size
- The total energy is invariant with respect to the order of the atoms
- The NNP can be used for systems of different sizes
- We

We know that the energy of a molecule does not change under a translation or a rotation. Thus our symmetry functions should be invariant with respect to these operations.

4.1.1 Symmetry functions

To construct adequate symmetry functions, we need to have a cutoff function that defines the atomic environments. One such function is [9],

$$f_c(R_{ij}) = \begin{cases} 0.5 \cdot \left[\cos\left(\frac{\pi R_{ij}}{R_c}\right) + 1 \right], & R_{ij} \leq R_c \\ 0, & R_{ij} > R_c \end{cases} \quad (4.2)$$

which is the monotonically decreasing part of a cosine function on $R_{ij} \in [0, R_c]$, seen in Figure 4.2.

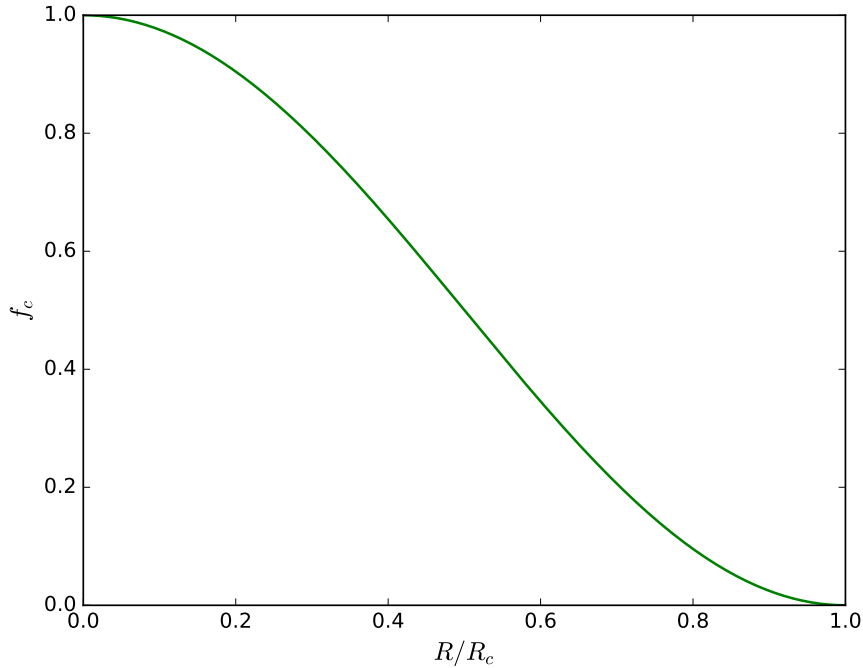


Figure 4.2: Plot of the cutoff function (4.2) applied in this thesis. This function is used to define a chemical environment around a central atom: only the atoms within the cutoff radius R_c contribute to its energy. These are called neighbouring atoms. The closer a neighbouring atom is, the larger the energy contribution, as is the case for most physical systems.

This function has the desirable property that it decreases with increasing distance R_{ij} between the central atom i and its neighbour j . At the cutoff radius R_c , it has zero value and slope, which is important to avoid discontinuities when computing energies and forces. Atoms beyond the cutoff radius are not a part of the central atom's chemical environment and therefore do not contribute to its energy.

Several types of many-body symmetry functions can be constructed based on these cutoff functions. They can be divided into two classes: *radial* symmetry functions, describing the radial distribution of neighbours up to the cutoff radius, and *angular* symmetry functions, specifying their angular arrangement. All symmetry functions depend on the positions of all the atoms inside the cutoff spheres. They make it possible to obtain a constant number of function values independent of the number of neighbours, which can change during MD simulations.

The most basic radial symmetry function is simply the sum of the cutoff

functions for all the neighbours j to atom i ,

$$G_i^1 = \sum_{j=1}^N f_c(R_{ij}) \quad (4.3)$$

We need a set of such functions with different cutoff radii to describe the radial arrangement of the neighbouring atoms. A better alternative is to use a sum of products of Gaussians and the cutoff function,

$$G_i^2 = \sum_{j=1}^N \exp[-\eta(R_{ij} - R_s)^2] \cdot f_c(R_{ij}) \quad (4.4)$$

We now have two parameters that can be adjusted to probe different radii. The width parameter η determines the radial extension of the symmetry functions, while the shifting parameter R_s displaces the Gaussians to improve the sensitivity of the symmetry functions at specific radii. A third option is

$$G_i^3 = \sum_{j=1}^N \cos(\kappa R_{ij}) \cdot f_c(R_{ij}) \quad (4.5)$$

which are damped cosine functions with a period length adjusted by parameter κ . We will however not use this function because of the existence of negative function values which can lead to atoms canceling each other's contribution to the sum. In Figure 4.3 we see the radial symmetry functions for several different parameters. It is clear that a set of such functions have a large flexibility when we tune the parameters in different ways. This is what makes them able to represent the radial distribution of neighbours around a central atom.

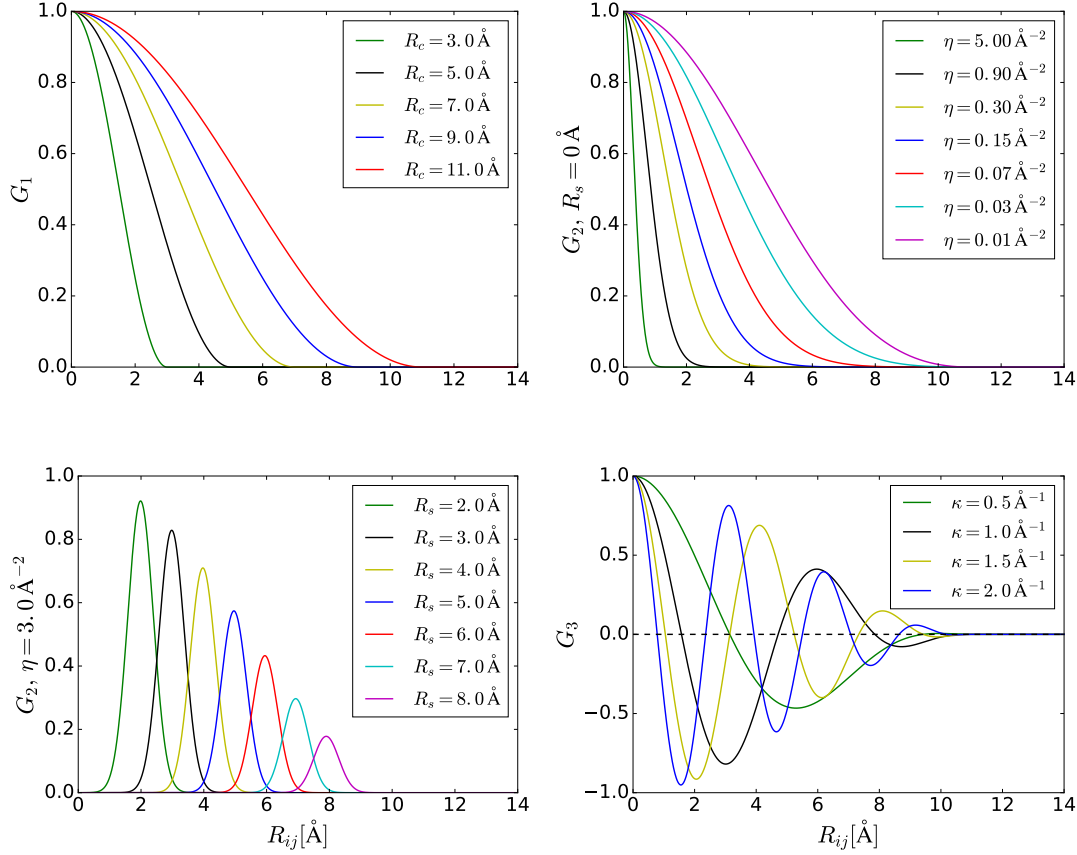


Figure 4.3: Radial symmetry functions G^1 , G^2 and G^3 for an atom with one neighbour only. A set of such functions represents the radial distribution of neighbours around a central atom placed at the origin. For G^2 and G^3 a cutoff $R_c = 11.0 \text{ \AA}$ has been used.

To obtain a suitable structural fingerprint of the atomic environments, we also need the angular distribution of neighbouring atoms. This can be achieved by using functions depending on θ_{ijk} , which is the angle formed by the central atom i and the two interatomic distances R_{ij} and R_{ik} . The potential is periodic with respect to this angle, so we can use the cosine of θ_{ijk} instead. We thus define an angular symmetry function as a sum over all cosines with respect to any possible pair of neighbours j and k , multiplied by Gaussians of the three interatomic distances and the respective cutoff functions,

$$G_i^4 = 2^{1-\zeta} \sum_{j \neq i} \sum_{k > j} [(1 + \lambda \cos \theta_{ijk})^\zeta \cdot \exp(-\eta(R_{ij}^2 + R_{ik}^2 + R_{jk}^2))]. \quad (4.6)$$

$$f_c(R_{ij})f_c(R_{ik})f_c(R_{jk})] \quad (4.7)$$

This function becomes zero if any of the interatomic distances is larger than R_c .

The parameter η takes here into account that the angular contribution depends on the atomic separations. The angular arrangement can be investigated by using different values for ζ while the normalization factor $2^{1-\zeta}$ ensures that the range of values is independent of the choice of ζ . The parameter $\lambda \in -1, 1$ can be used to invert the shape of the cosine function: for $\lambda = +1$ the maxima of the cosine terms are at $\theta_{ijk} = 0^\circ$, while for $\lambda = -1$ they are located at $\theta_{ijk} = 180^\circ$. The cutoff function R_{jk} is included to ensure that only triplets where all three inter-atomic distances are within the cutoff radius. Another function can be defined that has no constraint on R_{jk} ,

$$G_i^5 = 2^{1-\zeta} \sum_{j \neq i} \sum_{k > j} [(1 + \lambda \cos \theta_{ijk})^\zeta \cdot \exp(-\eta(R_{ij}^2 + R_{ik}^2)) \cdot \quad (4.8)$$

$$f_c(R_{ij})f_c(R_{ik})] \quad (4.9)$$

(4.9) will generally have larger function values than (4.7) because the lack of constraint on R_{jk} results in a larger number of non-zero terms in the summation. The angular part of G^4 and G^5 is identical, shown in Figure 4.4 for different values of ζ .

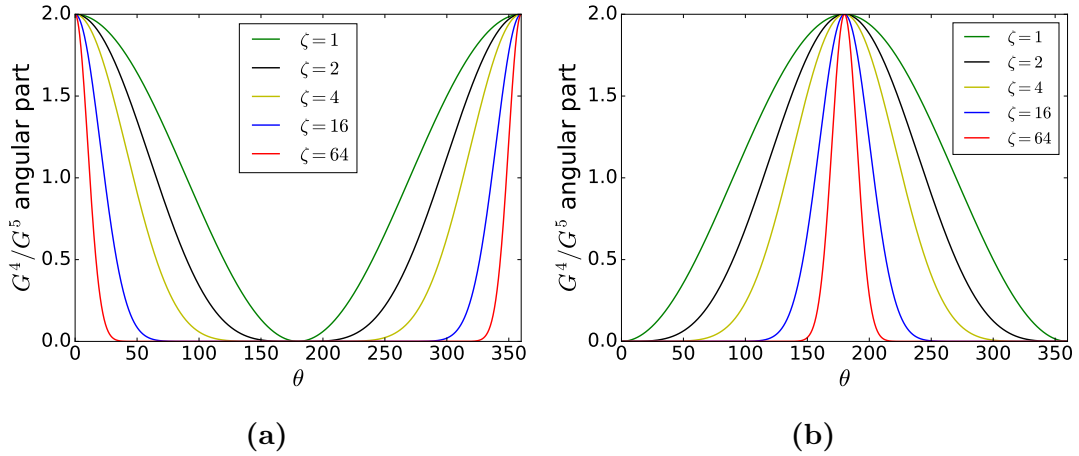


Figure 4.4: Angular part of symmetry functions G^4 and G^5 for an atom with one neighbour only. A set of such functions represents the angular distribution of neighbours around a central atom placed at the origin. $\lambda = +1$ for Figure 4.4a, $\lambda = -1$ for Figure 4.4b.

The parameter values R_c, η, R_s, ζ and λ are not automatically optimized, like the weights. Once a set of symmetry functions has been determined, they remain fixed during the training of the NN. They are therefore an integral part of the NN, and need to be evaluated together with the NN itself when applied in simulations. The symmetry function set have to be customized for different systems, but in our experience the quality of the fit is not very sensitive to the choice of

parameters. An adequate set will shorten the convergence time, but it will not necessarily result in a lower RMSE. However, it is important to use at least as many symmetry functions as there are degrees of freedom in the system. If this is not the case, the NN will receive insufficient information and will not be able to fit the data properly.

4.1.2 Symmetry functions and forces

To do MD simulations we need need the forces (2.2), i.e. the gradient of the PES. NNs have as we have seen well-defined functional forms, and analytical derivatives are therefore readily available. The pair-force on the central atom i from neighbour j w.r.t some atomic coordinate α is

$$F_{ij,\alpha} = -\frac{\partial E}{\partial \alpha_{ij}} = \sum_{\mu=1}^M \frac{\partial E}{\partial G_{\mu}} \frac{\partial G_{\mu}}{\partial \alpha_{ij}} \quad (4.10)$$

where $\alpha_{ij} = \alpha_i - \alpha_j = \alpha_j$ and M is the number of symmetry functions, i.e. the number of inputs to the NN. As seen by the r.h.s of this equation, we need to apply the chain rule: The first term $\partial E / \partial G_{\mu}$ is the derivative of the output (E) w.r.t. symmetry function number μ , the second term is the derivative of G_{μ} w.r.t. the atomic coordinate α . The former is obtained with a slightly modified form of backpropagation.

1. Instead of backpropagating the error, we backpropagate the derivative of the output neuron's activation function
2. The derivative is propagated all the way back to the input nodes. When training, the propagation stops at the weights of the first hidden layer.

This clarifies how a fully-connected feed-forward NN is built: To get the derivative of the output neuron w.r.t. the input neurons, we have to also compute the derivative of all the neurons and weights in-between because all nodes in each layer are connected to all nodes in the following layer. The term $\partial G_{\mu} / \partial \alpha_{ij}$ depends only on the form of the symmetry functions. See appendix for derivatives of these. Since all symmetry functions G_{μ} depend on the atomic coordinates of all the neighbours, we have to sum over all the input nodes to get the total force from neighbour j on central atom i .

The total force on atom i is obtained by calculating (4.10) for all neighbours and adding them,

$$F_{i,\alpha}^{TOT} = -\sum_{j=1}^N \sum_{\mu=1}^M \frac{\partial E}{\partial G_{\mu}} \frac{\partial G_{\mu}}{\partial \alpha_{ij}} \quad (4.11)$$

where N is the number of neighbours. TODO: THIS EXPRESSION HAS TO BE MODIFIED A BIT IF I INCLUDE MORE THAN ONE ATOM TYPE. THEN

THE TOTAL ENERGY E WILL BE DIFFERENT DEPENDING ON WHICH ATOM TYPE EACH NEIGHBOUR IS AND SO ON...

4.2 TensorFlow

TensorFlow is an open source software library for numerical computation using data flow graphs (DFG). A DFG describe mathematical computations with a directed graph of nodes and edges. Nodes typically represent mathematical operations, while the graph edges describe the input/output relationship between nodes. In TensorFlow, these edges carry dynamically-sized tensors, i.e. there is a flow of tensors between nodes in the graph.

We train a neural network by introducing different variables at the nodes. We want to minimize the cross-entropy (error) with respect to these variables. The derivatives of the output w.r.t. these are calculated using backpropagation. These derivatives are then used to minimize the cross-entropy using the gradient descent method or other optimization algorithms.

TensorFlow uses ReLU neurons, i.e. the activation function is

$$f(x) = \max(0, x) \quad (4.12)$$

where x is the input to the neuron. This is also known as a ramp function, or a rectifier (in analogy to half-wave rectification in electrical engineering). The above function has lately been argued to be more biologically plausible (REFERENCE) than the widely used sigmoid function and hyperbolic tangent. A smooth approximation to the rectifier is the analytic function

$$f(x) = \ln(1 + e^x) \quad (4.13)$$

which is called the softplus function. Using this kind of activation function prevents dead neurons. In practice we give initialize the neurons with a slightly positive initial bias.

4.2.1 Convolution

ANNs are supposed to mimic biological systems. We know that neurons in the visual cortex of animals have localized receptive fields, i.e. they respond only to stimuli in a certain location of the visual field. These regions are overlapping and covers the entire visual field. This can be exploited in ANNs by making the hidden layers only connect to a small contiguous region of the input layer. This also speeds up computations. These are called locally connected networks, unlike fully connected networks, where all the hidden layers are connected to the complete input layer.

Using this approach, we can aggregate statistics of these features at various locations. This is called pooling.

Part II

Implementation and validation

4.3 Time usage

The workflow when training a NN and using it in MD is as follows:

1. Generate training data
2. Train a NN to fit the data
3. Use the analytical expression for the trained NN as a potential in MD

The main application programming interface (API) for TensorFlow (TF) is Python. The NN is therefore trained using Python scripts, while C++ is utilized to run the MD simulations for speed purposes. Thus we need a way to transfer the NN from Python to C++.

There are built-in functions in TF that automatically saves computational graphs (architecture of network plus nodes for training) and variables (weights and activation functions) as binary files. These files can then be read and the NN evaluated using the TF C++ API. This API is however quite under-developed and the documentation is very sparse. Another alternative is to manually save the weights and biases as a text file, along with information on which activation function(s) that has been used. This file can then be used to reconstruct the NN in the same way that it is made in Python: By representing the connection weights between layers as matrices and the biases as vectors. The activation functions must be coded manually. A linear algebra library should be utilized for speed and simplicity; we have chosen to use Armadillo.

To find out which alternative is the most efficient for evaluating a NN, we compare the time usage in three cases: with TF in Python, with TF in C++ and with Armadillo. The results are shown in Figure 4.5.

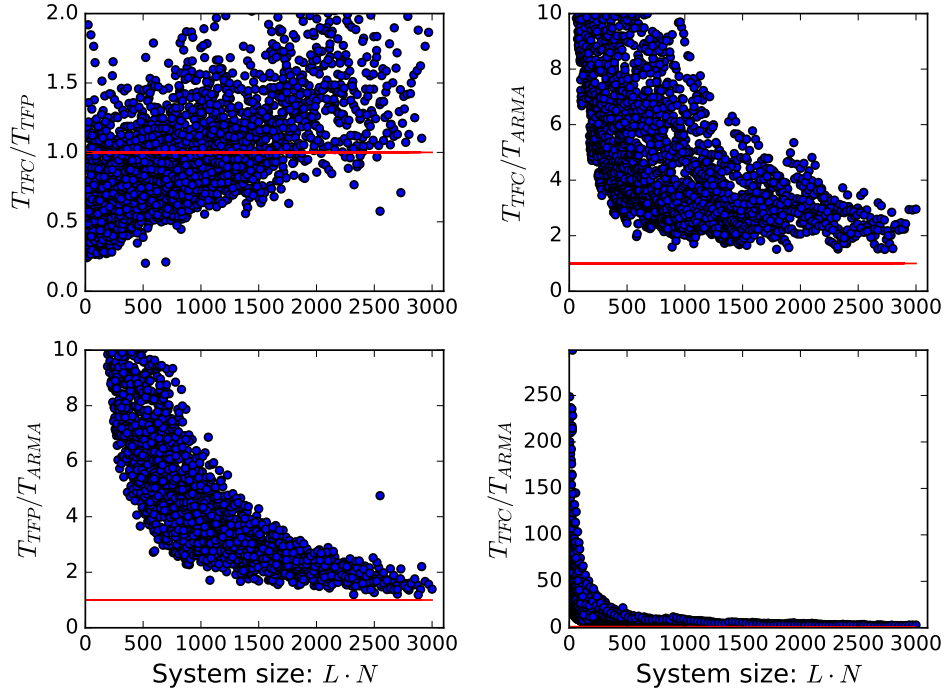


Figure 4.5: Scatter plot of CPU time when evaluating untrained NNs with random weights and sigmoid activation functions using the TF Python API (T_{TFP}), the TF C++ API (T_{TFC}) and Armadillo (T_{ARMA}). L is the number of layers, N is the number of nodes in each hidden layer. All the NNs have one input and one output. The time has been found by averaging over 50 evaluations for each NN architecture. T_{TFC}/T_{ARMA} is also shown in an uncut version to demonstrate how large the time difference is for very small NNs.

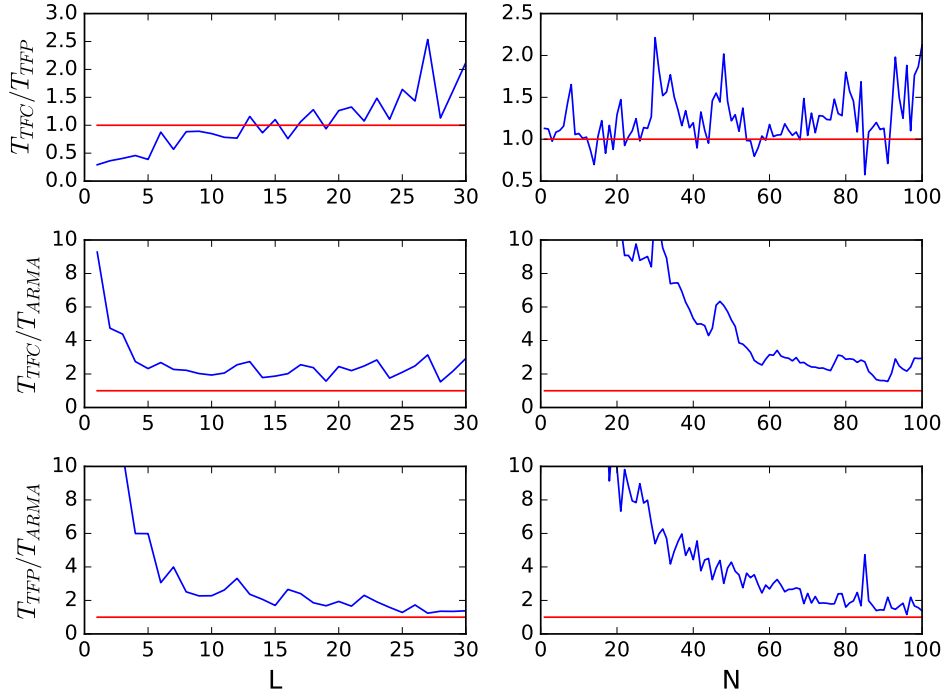


Figure 4.6: Left column: CPU time comparison for $N = 100$. Right column: CPU time comparison for $L = 30$.

As we can see, evaluation of an NN with the built-in TF methods are notably slower than with Armadillo for small NNs, but seems to converge to 1 for large networks. The convergence rate is however surprisingly low. One reason for this can be that the *Session* environment in TF has lot of overhead. However, the main computational strength of TF resides in GPUs, which will be tested out later on.

4.4 Training Lennard-Jones potential

To verify that the implementation of the NN and the backpropagation algorithm with Armadillo in our MD code is correct, we train a NN to reproduce the shifted Lennard Jones potential (REF TIL LIKNING). The error of the trained network on the training data interval is evaluated both in Python and C++. The two error plots are shown in Figure 4.7.

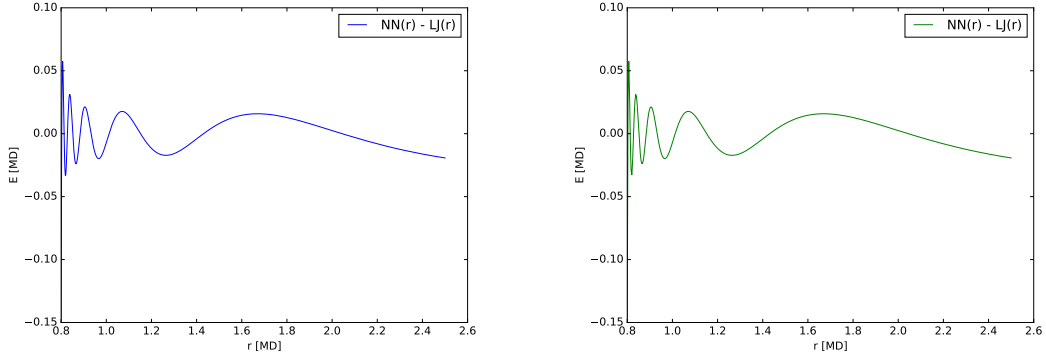


Figure 4.7: Error of a network trained to reproduce the shifted Lennard-Jones potential. The NN is trained in Python and the error on the training data interval is calculated, shown on the left. The NN is also written to file for evaluation in C++, shown on the right.

We observe that the error has exactly the same shape and value on the whole interval, we can therefore confirm that the NN is implemented correctly.

To test the implementation of the backpropagation algorithm, the same NN is differentiated in C++ and compared to the analytical derivative of the L-J potential. The gradient of a NN is defined as the derivative of one of the outputs with respect to the input(s). There are two differences when using backpropagation to find the gradient of a network compared to the use of backpropagation in training:

1. The derivative of the output neuron's activation function is backpropagated instead of the output error
2. We backpropagate all the way down to the input layer, during training we stop at the first hidden layer

As mentioned before, the backpropagation algorithm is simply an application of the chain rule: To obtain the derivative of the output node with respect to the input nodes we need to differentiate through all the nodes in-between because we have fully connected layers. The gradient error is plotted in Figure 4.8.

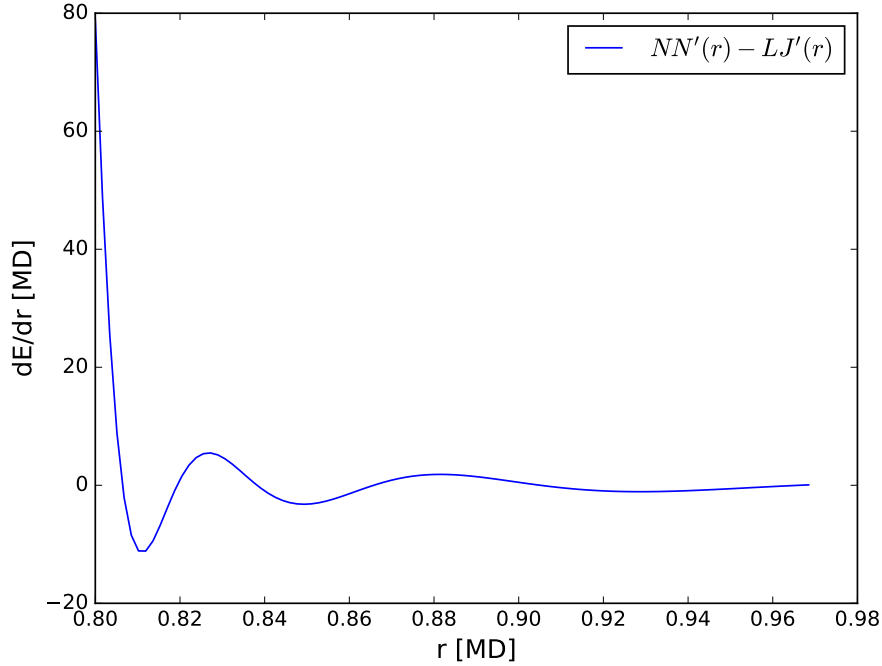
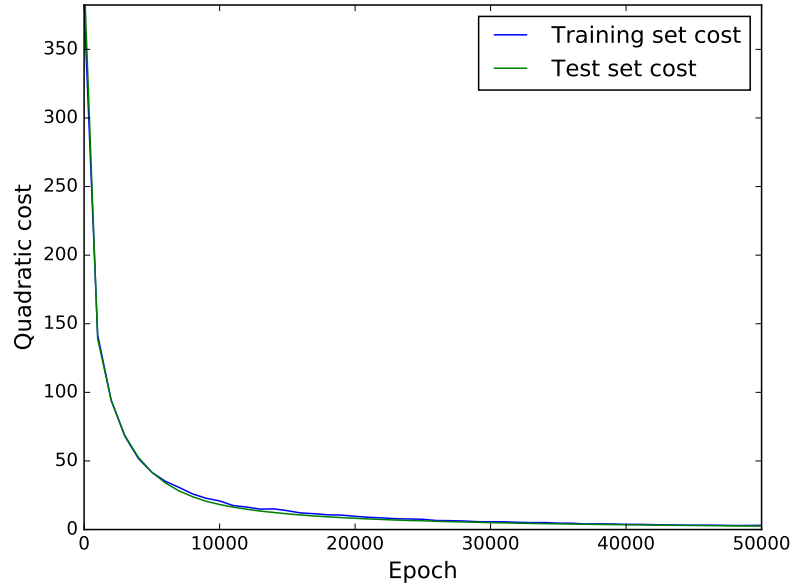


Figure 4.8: Error of the gradient of a NN trained to reproduce the shifted Lennard-Jones potential. The NN is trained in Python and evaluated and differentiated in C++ using the backpropagation algorithm. The result is compared to the analytical derivative of the LJ potential. Only a part of the training interval is shown, the graph is essentially flat after this point

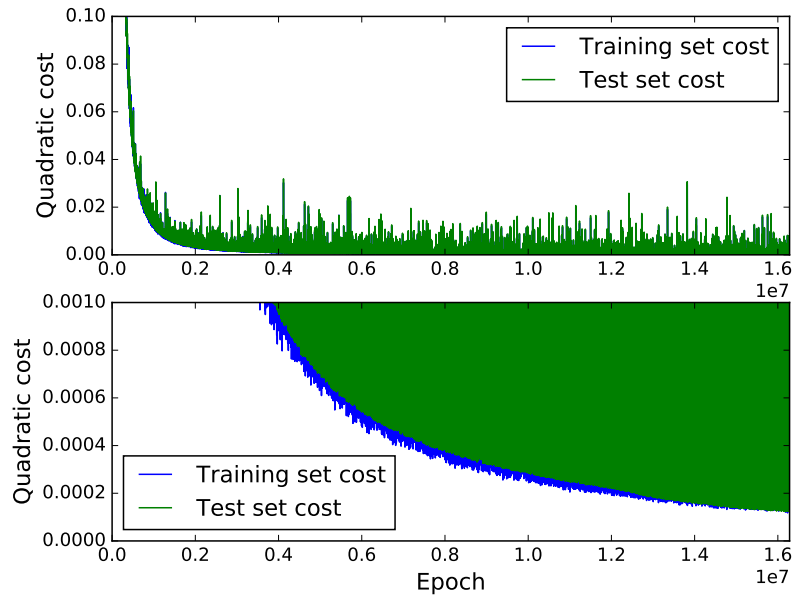
Again we get an oscillating error over the training interval. The oscillations near the left end of the interval is much larger than for the energy, but they get damped more quickly. The large error near the left end of the interval is a natural consequence of the shape of the LJ potential for these values: the $1/r^{13}$ term dominates with its very steep slope, the NN is not fed a sufficient amount of data in this region to exactly reproduce the rapid changing derivative. This tells us that it is important to train the NN on a larger data interval than needed in simulations or possibly feed the NN more data in critical regions during training.

4.4.1 Many-neighbour Lennard-Jones

A NN can also be trained to obtain the energy of several neighbouring atoms at once. The training procedure is exactly the same, but now we train the network so that the output is the sum of N LJ-potentials for N different distances r . The NN will have N inputs and 1 output. To get an error of the same order as for 1 neighbour, we need to have a lot more nodes. The result of a training session with $N = 20$ inputs is shown in figures Figure 4.9. (Training session 02.12-19.36.50)



(a) Zoomed out



(b) Zoomed in

Figure 4.9: Quadratic cost (3.9) of training set and test set for NN trained to yield the sum of the shifted LJ-potential of 20 neighbouring atoms. The NN has 20 inputs, one hidden layer with 200 neurons and one output node. All nodes in the hidden layer have sigmoid activation functions (3.5). Figure 4.9a shows the cost for the first 50 000 epochs, while Figure 4.9b displays the cost for all epochs, but on a smaller scale.

The cost rapidly decreases during the first 10 000 epochs, then it starts to flatten out. This behaviour seems to be general when doing regression with ANNs. There is no sign of overfitting: we do not observe an increase in the test cost relative to the training cost. This is (PROBABLY) due to the fact that our training set and test set are random numbers, i.e. there will be no specific patterns or configurations present in the test set that is not covered by the training set. When we start doing training on data produced by MD trajectories or QM calculations, the situation will be different.

The NN above outputs the total energy on an atom that is surrounded by N neighbouring atoms, but it does not provide the forces. These must be obtained by calculating the analytical derivative of the NN with backpropagation as described above for each time step. An alternative approach that may speed up the MD simulation is to have the NN also output the total force on the central atom. This can be done giving the cartesian coordinates and distance of N atoms as input to the NN so that each input vector looks like

$$(x_1, y_1, z_1, r_1, x_2, y_2, z_2, r_2, \dots, x_N, y_N, z_N, r_N) \quad (4.14)$$

which yields an output vector

$$(F_x^{TOT}, F_y^{TOT}, F_z^{TOT}, E^{TOT}) \quad (4.15)$$

We have to make sure that the input coordinates are in the range $x, y, z \in [-r_{cut}, r_{cut}]$ to cover the whole configuration space. The distances however, are still in the range $r \in [a, r_{cut}]$ where a is the minimum distance that two atoms can have in a given MD simulation (see function *energyAndForceCoordinates* in *generateData*). As the NN now have 4 inputs per neighbour and 4 outputs, we need to have more nodes to get a sufficiently low error. We run a training session with 5 neighbours (20 inputs) (session 13.12-12.59.51). The cost has the same shape as before, but the error is much larger to begin with and also converges slower (SKAL JEG GIDDE HA ENDA ET ERROR PLOT HER?).

One way to test if this works is to move only one of the input atoms while freezing the others and see if we reproduce the LJ-potential form. We load the above training session and test with the command `python approximateFunction2.py -load TrainingData/13.12-12.59.51/Checkpoints/ckpt-999 -test` with this code:

```
if testFlag:
    # pick an input vector
    coordinates = xTrain[0]
    coordinates = coordinates.reshape([1, self.inputs])
    neighbours = self.inputs/4
    xNN = np.zeros(neighbours)
    yNN = np.zeros(neighbours)
    zNN = np.zeros(neighbours)
```

```

rNN = np.zeros(neighbours)
# extract coordinates and distances
for i in range(neighbours):
    xNN[i] = coordinates[0,i*4]
    yNN[i] = coordinates[0,i*4 + 1]
    zNN[i] = coordinates[0,i*4 + 2]
    rNN[i] = coordinates[0,i*4 + 3]

# vary coordinates of only one atom and see
# if the resulting potential is similar to LJ
N = 500
r = np.linspace(0.8, 2.5, N)
energyNN = []
energyLJ = []
forceNN = []
forceLJ = []
xyz = np.zeros(3)
for i in range(N):
    r2 = r[i]**2
    xyz[0] = np.random.uniform(0, r2)
    xyz[1] = np.random.uniform(0, r2-xyz[0])
    xyz[2] = r2 - xyz[0] - xyz[1]
    #np.random.shuffle(xyz)
    x = np.sqrt(xyz[0])# * np.random.choice([-1,1])
    y = np.sqrt(xyz[1])# * np.random.choice([-1,1])
    z = np.sqrt(xyz[2])# * np.random.choice([-1,1])
    coordinates[0][0] = x; coordinates[0][1] = y;
    coordinates[0][2] = z
    coordinates[0][3] = r[i]
    energyAndForce = sess.run(prediction,
        feed_dict={self.x: coordinates})
    energyNN.append(energyAndForce[0][3])
    rNN[0] = r[i]
    energyLJ.append(np.sum(self.function(rNN)))
    forceNN.append(energyAndForce[0][0])
    xNN[0] = x
    forceLJ.append(np.sum(self.functionDerivative(rNN)*xNN/rNN))

# convert to arrays
energyNN = np.array(energyNN); energyLJ = np.array(energyLJ)
forceNN = np.array(forceNN); forceLJ = np.array(forceLJ)

# plot error
plt.plot(r, energyNN - energyLJ)
plt.show()
plt.plot(r, forceNN - forceLJ)
plt.show()
print 'Cost: ', (np.sum((energyNN - energyLJ)**2 + (forceNN
    - forceLJ)**2))/N

# see if the energy is zero when all neighbours is at

```

```

    cutoff distance
    inputz = np.array([1.87, 1.32, 1.006,
        2.5]*neighbours).reshape([1, self.inputs])
    r = np.array([2.5]*neighbours)
    energyLJ = sum(self.function(r))
    ef = sess.run(prediction, feed_dict={self.x: inputz})

    print 'Approximate energy: ', ef[0,3]
    print 'LJ energy: ', energyLJ

    numberOfEpochs = 0

```

which produces the following output:

```

Model restored
Cost: 0.132622806039
NN energy at cutoff: 0.125104
LJ energy at cutoff: 0.0
Time elapsed: 1.2e-05

```

We see that the NN energy at cutoff is not zero, because of the cost (IS THIS CORRECT? THE COST AT CUTOFF IS PRACTICALLY ZERO). The error of the NN energy and forces over the training data interval is shown in

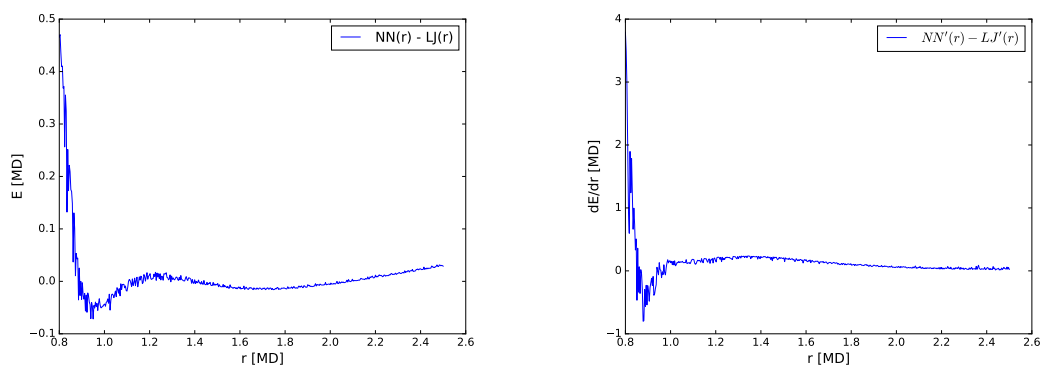


Figure 4.10: Error of a network trained to yield the total energy and force on a central atom from 5 neighbouring atoms. The energy contribution from each neighbour is a shifted LJ potential. Trained for 1e6 epochs, 100 nodes in 2 layers. The energy error is shown on the left, force error on the right.

The error oscillates in exactly the same way as for the 1-input NN above, suggesting that the training is done correctly. However, it is very important that the energies and forces in an MD simulation is consistent, and this is not ensured here, where the two are obtained in a somewhat independent way. Any inconsistency between the energies and forces will magnify for each time step, resulting

in inaccurate statistical properties like temperature and pressure. Therefore it is a better idea to have the NN only output the total energy, and find the forces by analytically differentiating the NN for each time step. 4 Å

.1 Appendix

.1.1 Symmetry functions derivatives

We list here the derivatives of the symmetry functions used for training neural networks. These have to be known to compute forces for MD simulations. In some cases we only show the derivative with respect to $R = \sqrt{x^2 + y^2 + z^2}$. Each component α can be calculated by

$$\frac{\partial R}{\partial \alpha} = \frac{\alpha}{R} \quad (16)$$

The following notation applies in the following,

$$R_{ij} = \left(x_{ij}^2 + y_{ij}^2 + z_{ij}^2 \right)^{1/2} = \left((x_j - x_i)^2 + (y_j - y_i)^2 + (z_j - z_i)^2 \right)^{1/2} \quad (17)$$

Cutoff function (4.2):

$$\frac{\partial f_c(R)}{\partial R} = -\frac{1}{2} \frac{\pi}{R_c} \sin \left(\frac{\pi R}{R_c} \right) = M(R) \quad (18)$$

This expression is valid for all combinations of indices. The derivatives w.r.t. to individual coordinates are

$$\frac{\partial f_c(R_{ij})}{\partial x_{ij}} = M(R_{ij}) \frac{x_{ij}}{R_{ij}} \quad (19)$$

and

$$\frac{\partial f_c(R_{ik})}{\partial x_{ik}} = M(R_{ik}) \frac{x_{ik}}{R_{ik}} \quad (20)$$

while

$$\frac{\partial f_c(R_{ij})}{\partial x_{ik}} = \frac{\partial f_c(R_{ik})}{\partial x_{ij}} = 0 \quad (21)$$

because atom i is always in origo, which means we can do the substitutions $x_{ij} \rightarrow x_j$ and $x_{ik} \rightarrow x_k$. SHOULD I DO THIS CONSEQUENTLY BELOW?????. Also,

$$\frac{\partial f_c(R_{jk})}{\partial x_{ij}} = -M(R_{jk}) \frac{x_{jk}}{R_{jk}} \quad (22)$$

and

$$\frac{\partial f_c(R_{jk})}{\partial x_{ij}} = M(R_{jk}) \frac{x_{jk}}{R_{jk}} \quad (23)$$

The derivative of each term in G_i^1 (4.3) is simply (18).

G_i^2 (4.4):

$$\frac{\partial G_i^2}{\partial R_{ij}} = \exp(-\eta(R_{ij} - R_s)^2) \left[2\eta(R_s - R_{ij}) + \frac{\partial f_c}{\partial R_{ij}} \right] \quad (24)$$

G_i^3 has not been used in this theses.

G_i^4 (4.7):

$$G_i^4 = F_1(\theta)F_2(R_{ij}, R_{ik}, R_{jk})F_3(R_{ij}, R_{ij}, R_{jk}) \quad (25)$$

where

$$F_1(\theta) = 2^{1-\zeta}(1 + \lambda \cos \theta_{ijk})^\zeta \quad (26)$$

$$F_2(R_{ij}, R_{ik}, R_{jk}) = \exp[-\eta(R_{ij}^2 + R_{ik}^2 + R_{jk}^2)] \quad (27)$$

$$F_3(R_{ij}, R_{ij}, R_{jk}) = f_c(R_{ij})f_c(R_{ik})f_c(R_{jk}) \quad (28)$$

where

$$\cos \theta = \frac{x_{ij}x_{ik} + y_{ij}y_{ik} + z_{ij}z_{ik}}{R_{ij}R_{ik}} \quad (29)$$

Using the product rule:

$$\frac{\partial G_i^4}{\partial x_{ij}} = \frac{\partial F_1}{\partial x_{ij}}F_2F_3 + F_1 \left(\frac{\partial F_2}{\partial x_{ij}}F_3 + F_2 \frac{\partial F_3}{\partial x_{ij}} \right) \quad (30)$$

$$\frac{\partial F_1}{\partial x_{ij}}F_2F_3 + F_1 \frac{\partial F_2}{\partial x_{ij}}F_3 + F_1F_2 \frac{\partial F_3}{\partial x_{ij}} \quad (31)$$

We have

$$\frac{\partial F_1}{\partial x_{ij}} = \frac{\partial F_1}{\partial \cos \theta} \frac{\partial \cos \theta}{\partial x_{ij}} \quad (32)$$

where

$$\frac{\partial F_1}{\partial \cos \theta} = \lambda \zeta (1 + \cos \theta)^{\zeta-1} = K \quad (33)$$

and

$$\frac{\partial \cos \theta}{\partial x_{ij}} = \frac{x_{ik}}{r_{ij}r_{ik}} - \frac{x_{ij} \cos \theta}{r_{ij}^2} \quad (34)$$

To get the corresponding expression for x_{ik} , simply substitute $x_{ij} \leftrightarrow x_{ik}$ and $r_{ij} \leftrightarrow r_{ik}$. Further,

$$\frac{\partial F_2}{\partial x_{ij}} = \frac{\partial F_2}{\partial r_{ij}} \frac{\partial r_{ij}}{\partial x_{ij}} = (-2\eta r_{ij} \frac{x_{ij}}{r_{ij}} - 2\eta r_{jk} \frac{x_{jk}}{r_{jk}}) \exp[-\eta(R_{ij}^2 + R_{ik}^2 + R_{jk}^2)] \quad (35)$$

$$= -2\eta F_2 x_{ij} - 2\eta F_2 x_{jk} \quad (36)$$

The corresponding expression for x_{ik} is

$$\frac{\partial F_2}{\partial x_{ik}} = -2\eta F_2 x_{ik} + 2\eta F_2 x_{jk} \quad (37)$$

Lastly,

$$\frac{\partial F_3}{\partial x_{ij}} = \frac{\partial f_c(r_{ij})}{\partial x_{ij}} f_c(r_{ik})f_c(r_{jk}) + f_c(r_{ij})f_c(r_{ik}) \frac{\partial f_c(r_{jk})}{\partial x_{ij}} \quad (38)$$

$$= f_c(r_{ik}) \left(M(r_{ij}) \frac{x_{ij}}{r_{ij}} f_c(r_{jk}) - f_c(r_{ij}) M(r_{jk}) \frac{x_{jk}}{r_{jk}} \right) \quad (39)$$

The corresponding expression for x_{ik} is

$$\frac{\partial F_3}{\partial x_{ik}} = f_c(r_{ij}) \left(\frac{\partial f_c(r_{ik})}{\partial x_{ik}} f_c(r_{jk}) + f_c(r_{ik}) \frac{\partial f_c(r_{jk})}{\partial x_{ik}} \right) \quad (40)$$

$$= f_c(r_{ij}) \left(M(r_{ik}) \frac{x_{ik}}{r_{ik}} f_c(r_{jk}) + f_c(r_{ik}) M(r_{jk}) \frac{x_{jk}}{r_{jk}} \right) \quad (41)$$

$$(42)$$

Then, taking (31) into account:

$$\frac{\partial G_i^4}{\partial x_{ij}} = \left(\frac{x_{ik}}{r_{ij}r_{ik}} - \frac{x_{ij} \cos \theta}{r_{ij}^2} \right) K F_2 F_3 - \quad (43)$$

$$2\eta F_1 F_2 F_3 x_{jk} - 2\eta F_1 F_2 F_3 x_{ij} + \quad (44)$$

$$F_1 F_2 f_c(r_{ik}) \left(M(r_{ij}) \frac{x_{ij}}{r_{ij}} f_c(r_{jk}) + f_c(r_{ij}) M(r_{jk}) \frac{x_{jk}}{r_{jk}} \right) \quad (45)$$

and for x_{ik} :

$$\frac{\partial G_i^4}{\partial x_{ik}} = \left(\frac{x_{ij}}{r_{ij}r_{ik}} - \frac{x_{ik} \cos \theta}{r_{ik}^2} \right) K F_2 F_3 + \quad (46)$$

$$2\eta F_1 F_2 F_3 x_{jk} - 2\eta F_1 F_2 F_3 x_{ik} + \quad (47)$$

$$F_1 F_2 f_c(r_{ij}) \left(M(r_{ik}) \frac{x_{ik}}{r_{ik}} f_c(r_{jk}) - f_c(r_{ik}) M(r_{jk}) \frac{x_{jk}}{r_{jk}} \right) \quad (48)$$

or in terms of x_{ij} and x_{ik} and x_{jk} :

$$\frac{\partial G_i^4}{\partial x_{ij}} = x_{ij} \left(- \frac{\cos \theta}{r_{ij}^2} K F_2 F_3 - 2\eta F_1 F_2 F_3 + \quad (49)$$

$$F_1 F_2 M(r_{ij}) f_c(r_{ik}) f_c(r_{jk}) \frac{1}{r_{ij}} \right) + \quad (50)$$

$$x_{ik} \frac{K F_2 F_3}{r_{ij}r_{ik}} - x_{jk} \left(F_1 F_2 M(r_{jk}) f_c(r_{ik}) f_c(r_{ij}) + 2\eta F_1 F_2 F_3 \right) \frac{1}{r_{jk}} \quad (51)$$

and

$$\frac{\partial G_i^4}{\partial x_{ik}} = \frac{x_{ik}}{r_{ik}} \left(- \frac{\cos \theta}{r_{ik}^2} K F_2 F_3 - 2\eta F_1 F_2 F_3 + \quad (52)$$

$$F_1 F_2 M(r_{ik}) f_c(r_{ij}) f_c(r_{jk}) \frac{1}{r_{ik}} \right) + \quad (53)$$

$$x_{ij} \frac{K F_2 F_3}{r_{ij}r_{ik}} + x_{jk} \left(F_1 F_2 M(r_{jk}) f_c(r_{ij}) f_c(r_{ik}) + 2\eta F_1 F_2 F_3 \right) \frac{1}{r_{jk}} \quad (54)$$

The derivative of G_i^5 (4.9) is found in a similar way,

$$G_i^5 = F_1(\theta)F_2(R_{ij}, R_{ik})F_3(R_{ij}, R_{ij}) \quad (55)$$

where

$$F_1(\theta) = 2^{1-\zeta}(1 + \lambda \cos \theta_{ijk})^\zeta \quad (56)$$

$$F_2(R_{ij}, R_{ik}) = \exp[-\eta(R_{ij}^2 + R_{ik}^2)] \quad (57)$$

$$F_3(R_{ij}, R_{ij}) = f_c(R_{ij})f_c(R_{ik}) \quad (58)$$

The derivative of F_1 is the same, while for F_2 we obtain

$$\frac{\partial F_2}{\partial x_{ij}} = -2\eta F_2 x_{ij} \quad (59)$$

and

$$\frac{\partial F_2}{\partial x_{ik}} = -2\eta F_2 x_{ik} \quad (60)$$

For F_3 ,

$$\frac{\partial F_3}{\partial x_{ij}} = \frac{\partial f_c(r_{ij})}{\partial x_{ij}} \frac{x_{ij}}{r_{ij}} f_c(r_{ik}) = M(r_{ij})f_c(r_{ik}) \frac{x_{ij}}{r_{ij}} \quad (61)$$

and

$$\frac{\partial F_3}{\partial x_{ik}} = f_c(r_{ij}) \frac{\partial f_c(r_{ik})}{\partial x_{ij}} \frac{x_{ik}}{r_{ik}} = f_c(r_{ij})M(r_{ik}) \frac{x_{ik}}{r_{ik}} \quad (62)$$

so that

$$\frac{\partial G_i^5}{\partial x_{ij}} = \left(\frac{x_{ik}}{r_{ij}r_{ik}} - \frac{x_{ij} \cos \theta}{r_{ij}^2} \right) K F_2 F_3 - 2\eta F_1 F_2 F_3 x_{ij} + \quad (63)$$

$$F_1 F_2 M(r_{ij}) f_c(r_{ik}) \frac{x_{ij}}{r_{ij}} \quad (64)$$

and

$$\frac{\partial G_i^5}{\partial x_{ik}} = \left(\frac{x_{ij}}{r_{ij}r_{ik}} - \frac{x_{ij} \cos \theta}{r_{ik}^2} \right) K F_2 F_3 - 2\eta F_1 F_2 F_3 x_{ik} + \quad (65)$$

$$F_1 F_2 M(r_{ik}) f_c(r_{ij}) \frac{x_{ik}}{r_{ik}} \quad (66)$$

In terms of x_{ij} and x_{ik} :

$$\frac{\partial G_i^5}{\partial x_{ij}} = x_{ij} \left(-\frac{\cos \theta}{r_{ij}^2} K F_2 F_3 - 2\eta F_1 F_2 F_3 + F_1 F_2 M(r_{ij}) f_c(r_{ik}) \frac{1}{r_{ij}} \right) + \quad (67)$$

$$x_{ik} \frac{K F_2 F_3}{r_{ij}r_{ik}} \quad (68)$$

and

$$\frac{\partial G_i^5}{\partial x_{ik}} = x_{ik} \left(-\frac{\cos \theta}{r_{ik}^2} K F_2 F_3 - 2\eta F_1 F_2 F_3 + F_1 F_2 M(r_{ik}) f_c(r_{ij}) \frac{1}{r_{ik}} \right) + \quad (69)$$

$$x_{ij} \frac{K F_2 F_3}{r_{ij} r_{ik}} \quad (70)$$

Bibliography

- [1] J. Behler *Neural network potential-energy surfaces in chemistry: a tool for large-scale simulations* Phys. Chem. Chem. Phys., 2011, **13**, 17930-17955 (2011)
- [2] Dragly, Sverre Arne *Bridging quantum mechanics and molecular dynamics with artificial neural networks* Thesis for the degree of Master of Science, University of Oslo (2014)
- [3] K. Hornik, M. Stinchcombe and H. White *Multilayer Feedforward Networks are Universal Approximators* Neural Networks **2**, 359 (1989).
- [4] R. Rojas. *Neural Networks: A Systematic Introduction*. Springer, 1996.
- [5] B. Karlik and A. V. Olgac *Performance Analysis of Various Activation Functions in Generalized MLP Architectures of Neural Networks*. IJAE **1** 111-122 (2011)
- [6] Y. LeCun, Y. Bengio and G. Hinton *Deep learning* Nature **521**, 436-444 (2015)
- [7] X. Glorot, A. Bordes, Y. Bengio *Deep Sparse Rectifier Neural Networks* JMLR **15**, 315-323 (2011)
- [8] D. E. Rumelhart, G. E. Hinton and R. J. Williams *Learning representations by back-propagating errors* Nature **323**, 533-536 (1986)
- [9] J. Behler *Constructing High-Dimensional Neural Network Potentials: A Tutorial Review* International Journal of Quantum Chemistry **115**, 1032-1050 (2015)
- [10] J. Behler *Neural network potential-energy surfaces in chemistry: a tool for large-scale simulations* Phys. Chem. Chem. Phys. **13**, 17930-17955 (2011)
- [11] Tuckerman, Berne and Martyna *Reversible multiscale molecular dynamics* J. Chem Phys, **97**, 1990 (1992).