

Boltzmann machines and deep learning and discussions of project 2

Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no^{1,2}

Department of Physics and Center for Computing in Science Education, University
of Oslo, Oslo, Norway¹

Department of Physics and Astronomy and Facility for Rare Ion Beams, Michigan
State University, East Lansing, Michigan, USA²

April 12

© 1999-2024, Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no. Released under CC

Attribution-NonCommercial 4.0 license

Plans for the week of April 8-12, 2024

1. Neural Networks and Boltzmann Machines
 2. If time, start discussion on how to implement Slater determinants, see last part of slides. This topic will be continued next week.
-
- ▶ Video of lecture
 - ▶ Handwritten notes

Energy models

Last week we defined a domain \mathbf{X} of stochastic variables $\mathbf{X} = \{x_0, x_1, \dots, x_{n-1}\}$ with a pertinent probability distribution

$$p(\mathbf{X}) = \prod_{x_i \in \mathbf{X}} p(x_i),$$

where we have assumed that the random variables x_i are all independent and identically distributed (iid).

We will now assume that we can define this function in terms of optimization parameters Θ , which could be the biases and weights of deep network, and a set of hidden variables we also assume to be random variables which also are iid. The domain of these variables is $\mathbf{H} = \{h_0, h_1, \dots, h_{m-1}\}$.

Probability model

We define a probability

$$p(x_i, h_j; \Theta) = \frac{f(x_i, h_j; \Theta)}{Z(\Theta)},$$

where $f(x_i, h_j; \Theta)$ is a function which we assume is larger or equal than zero and obeys all properties required for a probability distribution and $Z(\Theta)$ is a normalization constant. Inspired by statistical mechanics, we call it often for the partition function. It is defined as (assuming that we have discrete probability distributions)

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta).$$

Marginal and conditional probabilities

We can in turn define the marginal probabilities

$$p(x_i; \Theta) = \frac{\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)}{Z(\Theta)},$$

and

$$p(h_j; \Theta) = \frac{\sum_{x_i \in \mathbf{X}} f(x_i, h_j; \Theta)}{Z(\Theta)}.$$

Change of notation

Note the change to a vector notation. A variable like \mathbf{x} represents now a specific **configuration**. We can generate an infinity of such configurations. The final partition function is then the sum over all such possible configurations, that is

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta),$$

changes to

$$Z(\Theta) = \sum_{\mathbf{x}} \sum_{\mathbf{h}} f(\mathbf{x}, \mathbf{h}; \Theta).$$

If we have a binary set of variable x_i and h_j and M values of x_i and N values of h_j we have in total 2^M and 2^N possible \mathbf{x} and \mathbf{h} configurations, respectively.

We see that even for the modest binary case, we can easily approach a number of configuration which is not possible to deal with.

Optimization problem

At the end, we are not interested in the probabilities of the hidden variables. The probability we thus want to optimize is

$$p(\mathbf{X}; \Theta) = \prod_{x_i \in \mathbf{X}} p(x_i; \Theta) = \prod_{x_i \in \mathbf{X}} \left(\frac{\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)}{Z(\Theta)} \right),$$

which we rewrite as

$$p(\mathbf{X}; \Theta) = \frac{1}{Z(\Theta)} \prod_{x_i \in \mathbf{X}} \left(\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta) \right).$$

Further simplifications

We simplify further by rewriting it as

$$p(\mathbf{X}; \Theta) = \frac{1}{Z(\Theta)} \prod_{x_i \in \mathbf{X}} f(x_i; \Theta),$$

where we used $p(x_i; \Theta) = \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)$. The optimization problem is then

$$\arg \max_{\Theta \in \mathbb{R}^p} p(\mathbf{X}; \Theta).$$

Optimizing the logarithm instead

Computing the derivatives with respect to the parameters Θ is easier (and equivalent) if we compute the logarithm of the probability. We will thus optimize

$$\arg \max_{\Theta \in \mathbb{R}^p} \log p(\mathbf{X}; \Theta),$$

which leads to

$$\nabla_{\Theta} \log p(\mathbf{X}; \Theta) = 0.$$

Expression for the gradients

This leads to the following equation

$$\nabla_{\Theta} \log p(\mathbf{X}; \Theta) = \nabla_{\Theta} \left(\sum_{x_i \in \mathbf{X}} \log f(x_i; \Theta) \right) - \nabla_{\Theta} \log Z(\Theta) = 0.$$

The first term is called the positive phase and we assume that we have a model for the function f from which we can sample values. Below we will develop an explicit model for this. The second term is called the negative phase and is the one which leads to more difficulties.

The derivative of the partition function

The partition function, defined above as

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta),$$

is in general the most problematic term. In principle both x and h can span large degrees of freedom, if not even infinitely many ones, and computing the partition function itself is often not desirable or even feasible. The above derivative of the partition function can however be written in terms of an expectation value which is in turn evaluated using Monte Carlo sampling and the theory of Markov chains, popularly shortened to MCMC (or just MC²).

Explicit expression for the derivative

We can rewrite

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\nabla_{\Theta} Z(\Theta)}{Z(\Theta)},$$

which reads in more detail

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\nabla_{\Theta} \sum_{x_i \in \mathbf{x}} f(x_i; \Theta)}{Z(\Theta)}.$$

We can rewrite the function f (we have assumed that is larger or equal than zero) as $f = \exp \log f$. We can then rewrite the last equation as

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\sum_{x_i \in \mathbf{x}} \nabla_{\Theta} \exp \log f(x_i; \Theta)}{Z(\Theta)}.$$

Final expression

Taking the derivative gives us

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\sum_{x_i \in \mathbf{X}} f(x_i; \Theta) \nabla_{\Theta} \log f(x_i; \Theta)}{Z(\Theta)},$$

which is the expectation value of $\log f$

$$\nabla_{\Theta} \log Z(\Theta) = \sum_{x_i \in \mathbf{X}} p(x_i; \Theta) \nabla_{\Theta} \log f(x_i; \Theta),$$

that is

$$\nabla_{\Theta} \log Z(\Theta) = \mathbb{E}(\log f(x_i; \Theta)).$$

This quantity is evaluated using Monte Carlo sampling, with Gibbs sampling as the standard sampling rule. Before we discuss the explicit algorithms, we need to remind ourselves about Markov chains and sampling rules like the Metropolis-Hastings algorithm and Gibbs sampling.

Introducing the energy model

As we will see below, a typical Boltzmann machines employs a probability distribution

$$p(\mathbf{x}, \mathbf{h}; \Theta) = \frac{f(\mathbf{x}, \mathbf{h}; \Theta)}{Z(\Theta)},$$

where $f(\mathbf{x}, \mathbf{h}; \Theta)$ is given by a so-called energy model. If we assume that the random variables x_i and h_j take binary values only, for example $x_i, h_j = \{0, 1\}$, we have a so-called binary-binary model where

$$f(\mathbf{x}, \mathbf{h}; \Theta) = -E(\mathbf{x}, \mathbf{h}; \Theta) = \sum_{x_i \in \mathbf{X}} x_i a_i + \sum_{h_j \in \mathbf{H}} b_j h_j + \sum_{x_i \in \mathbf{X}, h_j \in \mathbf{H}} x_i w_{ij} h_j,$$

where the set of parameters are given by the biases and weights $\Theta = \{\mathbf{a}, \mathbf{b}, \mathbf{W}\}$. **Note the vector notation** instead of x_i and h_j for f . The vectors \mathbf{x} and \mathbf{h} represent a specific instance of stochastic variables x_i and h_j . These arrangements of \mathbf{x} and \mathbf{h} lead to a specific energy configuration.

More compact notation

With the above definition we can write the probability as

$$p(\mathbf{x}, \mathbf{h}; \Theta) = \frac{\exp(\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{x}^T \mathbf{W} \mathbf{h})}{Z(\Theta)},$$

where the biases \mathbf{a} and \mathbf{h} and the weights defined by the matrix \mathbf{W} are the parameters we need to optimize.

Anticipating results to be derived

Since the binary-binary energy model is linear in the parameters a_i , b_j and w_{ij} , it is easy to see that the derivatives with respect to the various optimization parameters yield expressions used in the evaluation of gradients like

$$\frac{\partial E(\mathbf{x}, \mathbf{h}; \Theta)}{\partial w_{ij}} = -x_i h_j,$$

and

$$\frac{\partial E(\mathbf{x}, \mathbf{h}; \Theta)}{\partial a_i} = -x_i,$$

and

$$\frac{\partial E(\mathbf{x}, \mathbf{h}; \Theta)}{\partial b_j} = -h_j.$$

Boltzmann Machines, marginal and conditional probabilities

A generative model can learn to represent and sample from a probability distribution. The core idea is to learn a parametric model of the probability distribution from which the training data was drawn. As an example

1. A model for images could learn to draw new examples of cats and dogs, given a training dataset of images of cats and dogs.
2. Generate a sample of an ordered or disordered Ising model phase, having been given samples of such phases.
3. Model the trial function for Monte Carlo calculations

Generative and discriminative models

Generative and discriminative models use both gradient-descent based learning procedures for minimizing cost functions. However, in energy based models we don't use backpropagation and automatic differentiation for computing gradients, instead we turn to Markov Chain Monte Carlo methods.

A typical deep neural network has several hidden layers. A restricted Boltzmann machine has normally one hidden layer, however several RBMs can be stacked to make up deep Belief Networks, of which they constitute the building blocks.

Basics of the Boltzmann machine

A BM is what we would call an undirected probabilistic graphical model with stochastic continuous or discrete units.

It is interpreted as a stochastic recurrent neural network where the state of each unit(neurons/nodes) depends on the units it is connected to. The weights in the network represent thus the strength of the interaction between various units/nodes.

More about the basics

A standard BM network is divided into a set of observable and visible units \mathbf{x} and a set of unknown hidden units/nodes \mathbf{h} .

Additionally there can be bias nodes for the hidden and visible layers. These biases are normally set to 1.

BMs are stackable, meaning they can be trained a BM which serves as input to another BM. We can construct deep networks for learning complex PDFs. The layers can be trained one after another, a feature which makes them popular in deep learning

Difficult to train

However, they are often hard to train. This leads to the introduction of so-called restricted BMs, or RBMS. Here we take away all lateral connections between nodes in the visible layer as well as connections between nodes in the hidden layer.

The network layers

1. A function \mathbf{x} that represents the visible layer, a vector of M elements (nodes). This layer represents both what the RBM might be given as training input, and what we want it to be able to reconstruct. This might for example be the pixels of an image, the spin values of the Ising model, or coefficients representing speech.
2. The function \mathbf{h} represents the hidden, or latent, layer. A vector of N elements (nodes). Also called "feature detectors".

Goal of hidden layer

The goal of the hidden layer is to increase the model's expressive power. We encode complex interactions between visible variables by introducing additional, hidden variables that interact with visible degrees of freedom in a simple manner, yet still reproduce the complex correlations between visible degrees in the data once marginalized over (integrated out).

The parameters

The network parameters, to be optimized/learned:

1. \mathbf{a} represents the visible bias, a vector of same length M as \mathbf{x} .
2. \mathbf{b} represents the hidden bias, a vector of same length N as \mathbf{h} .
3. \mathbf{W} represents the interaction weights, a matrix of size $M \times N$.

Note that we have specified the lengths of \mathbf{b} and \mathbf{h} . These lengths define the number of visible and hidden units, respectively.

Joint distribution

The restricted Boltzmann machine is described by a Boltzmann distribution

$$P_{rbm}(\mathbf{x}, \mathbf{h}, \Theta) = \frac{1}{Z(\Theta)} \exp -(E(\mathbf{x}, \mathbf{h}, \Theta)),$$

where Z is the normalization constant or partition function discussed earlier and defined as

$$Z(\Theta) = \int \int \exp -E(\mathbf{x}, \mathbf{h}, \Theta) d\mathbf{x} d\mathbf{h}.$$

It is common to set the temperature T to one. It is omitted in the equations above. The energy is thus a dimensionless function.

Network Elements, the energy function

The function $E(\mathbf{x}, \mathbf{h}, \Theta)$ gives the **energy** of a configuration (pair of vectors) (\mathbf{x}, \mathbf{h}) . The lower the energy of a configuration, the higher the probability of it. This function also depends on the parameters \mathbf{a} , \mathbf{b} and W . Thus, when we adjust them during the learning procedure, we are adjusting the energy function to best fit our problem.

Defining different types of RBMs

There are different variants of RBMs, and the differences lie in the types of visible and hidden units we choose as well as in the implementation of the energy function $E(\mathbf{x}, \mathbf{h}, \Theta)$. The connection between the nodes in the two layers is given by the weights w_{ij} .

Binary-Binary RBM:

RBMs were first developed using binary units in both the visible and hidden layer. The corresponding energy function is defined as follows:

$$E(\mathbf{x}, \mathbf{h}, \Theta) = - \sum_i^M x_i a_i - \sum_j^N b_j h_j - \sum_{i,j}^{M,N} x_i w_{ij} h_j,$$

where the binary values taken on by the nodes are most commonly 0 and 1.

Gaussian-binary RBM

Another variant is the RBM where the visible units are Gaussian while the hidden units remain binary:

$$E(\mathbf{x}, \mathbf{h}, \Theta) = \sum_i^M \frac{(x_i - a_i)^2}{2\sigma_i^2} - \sum_j^N b_j h_j - \sum_{i,j}^{M,N} \frac{x_i w_{ij} h_j}{\sigma_i^2}.$$

This type of RBMs are useful when we model continuous data (i.e., we wish \mathbf{x} to be continuous). The parameter σ_i^2 is meant to represent a variance and is often just set to one.

Cost function

When working with a training dataset, the most common training approach is maximizing the log-likelihood of the training data. The log likelihood characterizes the log-probability of generating the observed data using our generative model. Using this method our cost function is chosen as the negative log-likelihood. The learning then consists of trying to find parameters that maximize the probability of the dataset, and is known as Maximum Likelihood Estimation (MLE).

Denoting the parameters as $\Theta = a_1, \dots, a_M, b_1, \dots, b_N, w_{11}, \dots, w_{MN}$, the log-likelihood is given by

$$\begin{aligned}\mathcal{L}(\{\Theta_i\}) &= \langle \log P_{\theta}(\mathbf{x}) \rangle_{data} \\ &= -\langle E(\mathbf{x}; \{\Theta_i\}) \rangle_{data} - \log Z(\{\Theta_i\}),\end{aligned}$$

where we used that the normalization constant does not depend on the data, $\langle \log Z(\{\Theta_i\}) \rangle = \log Z(\{\Theta_i\})$. Our cost function is the negative log-likelihood, $\mathcal{C}(\{\Theta_i\}) = -\mathcal{L}(\{\Theta_i\})$

Optimization / Training

The training procedure of choice often is Stochastic Gradient Descent (SGD). It consists of a series of iterations where we update the parameters according to the equation

$$\Theta_{k+1} = \Theta_k - \eta \nabla \mathcal{C}(\Theta_k)$$

at each k -th iteration. There are a range of variants of the algorithm which aim at making the learning rate η more adaptive so the method might be more efficient while remaining stable.

Gradients

We now need the gradient of the cost function in order to minimize it. We find that

$$\begin{aligned}\frac{\partial \mathcal{C}(\{\Theta_i\})}{\partial \Theta_i} &= \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{data} + \frac{\partial \log Z(\{\Theta_i\})}{\partial \Theta_i} \\ &= \langle O_i(\mathbf{x}) \rangle_{data} - \langle O_i(\mathbf{x}) \rangle_{model}.\end{aligned}$$

Simplifications

In order to simplify notation we defined the "operator"

$$O_i(\mathbf{x}) = \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i},$$

and used the statistical mechanics relationship between expectation values and the log-partition function:

$$\langle O_i(\mathbf{x}) \rangle_{model} = \text{Tr} P_{\Theta}(\mathbf{x}) O_i(\mathbf{x}) = - \frac{\partial \log Z(\{\Theta_i\})}{\partial \Theta_i}.$$

Positive and negative phases

As discussed earlier, the data-dependent term in the gradient is known as the positive phase of the gradient, while the model-dependent term is known as the negative phase of the gradient. The aim of the training is to lower the energy of configurations that are near observed data points (increasing their probability), and raising the energy of configurations that are far from observed data points (decreasing their probability).

Gradient examples

The gradient of the negative log-likelihood cost function of a Binary-Binary RBM is then

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial w_{ij}} = \langle x_i h_j \rangle_{data} - \langle x_i h_j \rangle_{model}$$

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial a_{ij}} = \langle x_i \rangle_{data} - \langle x_i \rangle_{model}$$

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial b_{ij}} = \langle h_i \rangle_{data} - \langle h_i \rangle_{model}.$$

To get the expectation values with respect to the *data*, we set the visible units to each of the observed samples in the training data, then update the hidden units according to the conditional probability found before. We then average over all samples in the training data to calculate expectation values with respect to the data.

Kullback-Leibler relative entropy

When the goal of the training is to approximate a probability distribution, as it is in generative modeling, another relevant measure is the **Kullback-Leibler divergence**, also known as the relative entropy or Shannon entropy. It is a non-symmetric measure of the dissimilarity between two probability density functions p and q . If p is the unknown probability which we approximate with q , we can measure the difference by

$$\text{KL}(p||q) = \int_{-\infty}^{\infty} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}.$$

Kullback-Leibler divergence

Thus, the Kullback-Leibler divergence between the distribution of the training data $f(\mathbf{x})$ and the model distribution $p(\mathbf{x}|\Theta)$ is

$$\begin{aligned}\text{KL}(f(\mathbf{x})||p(\mathbf{x}|\Theta)) &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log \frac{f(\mathbf{x})}{p(\mathbf{x}|\Theta)} d\mathbf{x} \\ &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} - \int_{-\infty}^{\infty} f(\mathbf{x}) \log p(\mathbf{x}|\Theta) d\mathbf{x} \\ &= \langle \log f(\mathbf{x}) \rangle_{f(\mathbf{x})} - \langle \log p(\mathbf{x}|\Theta) \rangle_{f(\mathbf{x})} \\ &= \langle \log f(\mathbf{x}) \rangle_{data} + \langle E(\mathbf{x}) \rangle_{data} + \log Z \\ &= \langle \log f(\mathbf{x}) \rangle_{data} + \mathcal{C}_{LL}.\end{aligned}$$

Maximizing log-likelihood

The first term is constant with respect to Θ since $f(\mathbf{x})$ is independent of Θ . Thus the Kullback-Leibler Divergence is minimal when the second term is minimal. The second term is the log-likelihood cost function, hence minimizing the Kullback-Leibler divergence is equivalent to maximizing the log-likelihood.

To further understand generative models it is useful to study the gradient of the cost function which is needed in order to minimize it using methods like stochastic gradient descent.

More on the partition function

The partition function is the generating function of expectation values, in particular there are mathematical relationships between expectation values and the log-partition function. In this case we have

$$\left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{model} = \int p(\mathbf{x}|\Theta) \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} d\mathbf{x} = -\frac{\partial \log Z(\Theta_i)}{\partial \Theta_i}.$$

Here $\langle \cdot \rangle_{model}$ is the expectation value over the model probability distribution $p(\mathbf{x}|\Theta)$.

Setting up for gradient descent calculations

Using the previous relationship we can express the gradient of the cost function as

$$\begin{aligned}\frac{\partial \mathcal{C}_{LL}}{\partial \Theta_i} &= \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{data} + \frac{\partial \log Z(\Theta_i)}{\partial \Theta_i} \\ &= \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{data} - \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{model}\end{aligned}$$

Difference of moments

This expression shows that the gradient of the log-likelihood cost function is a **difference of moments**, with one calculated from the data and one calculated from the model. The data-dependent term is called the **positive phase** and the model-dependent term is called the **negative phase** of the gradient. We see now that minimizing the cost function results in lowering the energy of configurations \mathbf{x} near points in the training data and increasing the energy of configurations not observed in the training data. That means we increase the model's probability of configurations similar to those in the training data.

More observations

The gradient of the cost function also demonstrates why gradients of unsupervised, generative models must be computed differently from those of, for example, FNNs. While the data-dependent expectation value is easily calculated based on the samples \mathbf{x}_i in the training data, we must sample from the model in order to generate samples from which to calculate the model-dependent term. We sample from the model by using MCMC-based methods. We cannot sample from the model directly because the partition function Z is generally intractable.

Adding hyperparameters

As in supervised machine learning problems, the goal is also here to perform well on **unseen** data, that is to have good generalization from the training data. The distribution $f(x)$ we approximate is not the **true** distribution we wish to estimate, it is limited to the training data. Hence, in unsupervised training as well it is important to prevent overfitting to the training data. Thus it is common to add regularizers to the cost function in the same manner as we discussed for say linear regression.

Mathematical details

Because we are restricted to potential functions which are positive it is convenient to express them as exponentials.

The original RBM had binary visible and hidden nodes. They were shown to be universal approximators of discrete distributions. It was also shown that adding hidden units yields strictly improved modelling power.

Binary-binary (BB) RBMs

The common choice of binary values are 0 and 1. However, in some physics applications, -1 and 1 might be a more natural choice. We will here use 0 and 1. We have the energy function

$$E_{BB}(\mathbf{x}, \mathbf{h}, \Theta) = - \sum_i^M x_i a_i - \sum_j^N b_j h_j - \sum_{i,j}^{M,N} x_i w_{ij} h_j.$$

Marginal probability

We have the binary-binary marginal probability defined as

$$\begin{aligned} p_{BB}(\mathbf{x}, \mathbf{h}, \Theta) &= \frac{1}{Z_{BB}(\Theta)} e^{\sum_i^M a_i x_i + \sum_j^N b_j h_j + \sum_{ij}^{M,N} x_i w_{ij} h_j} \\ &= \frac{1}{Z_{BB}(\Theta)} e^{\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{x}^T \mathbf{W} \mathbf{h}} \end{aligned}$$

with the partition function

$$Z_{BB}(\Theta) = \sum_{\mathbf{x}, \mathbf{h}} e^{\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{x}^T \mathbf{W} \mathbf{h}}.$$

Marginal Probability Density Function for the visible units

In order to find the probability of any configuration of the visible units we derive the marginal probability density function.

$$\begin{aligned} p_{BB}(\mathbf{x}, \Theta) &= \sum_{\mathbf{h}} p_{BB}(\mathbf{x}, \mathbf{h}, \Theta) \\ &= \frac{1}{Z_{BB}} \sum_{\mathbf{h}} e^{\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{x}^T \mathbf{W} \mathbf{h}} \\ &= \frac{1}{Z_{BB}} e^{\mathbf{a}^T \mathbf{x}} \sum_{\mathbf{h}} e^{\sum_j^N (b_j + \mathbf{x}^T \mathbf{w}_{*j}) h_j} \\ &= \frac{1}{Z_{BB}} e^{\mathbf{a}^T \mathbf{x}} \sum_{\mathbf{h}} \prod_j^N e^{(b_j + \mathbf{x}^T \mathbf{w}_{*j}) h_j} \\ &= \frac{1}{Z_{BB}} e^{\mathbf{a}^T \mathbf{x}} \left(\sum_{h_1} e^{(b_1 + \mathbf{x}^T \mathbf{w}_{*1}) h_1} \times \sum_{h_2} e^{(b_2 + \mathbf{x}^T \mathbf{w}_{*2}) h_2} \times \right. \\ &\quad \left. \dots \times \sum_{h_N} e^{(b_N + \mathbf{x}^T \mathbf{w}_{*N}) h_N} \right) \end{aligned}$$

Marginal probability for hidden units

A similar derivation yields the marginal probability of the hidden units

$$p_{BB}(\mathbf{h}, \Theta) = \frac{1}{Z_{BB}(\Theta)} e^{\mathbf{b}^T \mathbf{h}} \prod_i^M (1 + e^{a_i + \mathbf{w}_{i*}^T \mathbf{h}}).$$

Conditional Probability Density Functions

We derive the probability of the hidden units given the visible units using Bayes' rule (we drop the explicit Θ dependence)

$$\begin{aligned} p_{BB}(\mathbf{h}|\mathbf{x}) &= \frac{p_{BB}(\mathbf{x}, \mathbf{h})}{p_{BB}(\mathbf{x})} \\ &= \frac{\frac{1}{Z_{BB}} e^{\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{h} + \mathbf{x}^T \mathbf{W} \mathbf{h}}}{\frac{1}{Z_{BB}} e^{\mathbf{a}^T \mathbf{x}} \prod_j^N (1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}})} \\ &= \frac{e^{\mathbf{a}^T \mathbf{x}} e^{\sum_j^N (b_j + \mathbf{x}^T \mathbf{w}_{*j}) h_j}}{e^{\mathbf{a}^T \mathbf{x}} \prod_j^N (1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}})} \\ &= \prod_j^N \frac{e^{(b_j + \mathbf{x}^T \mathbf{w}_{*j}) h_j}}{1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}}} \\ &= \prod_j^N p_{BB}(h_j|\mathbf{x}). \end{aligned}$$

On and off probabilities

From this we find the probability of a hidden unit being "on" or "off":

$$\begin{aligned} p_{BB}(h_j = 1|\mathbf{x}) &= \frac{e^{(b_j + \mathbf{x}^T \mathbf{w}_{*j})h_j}}{1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}}} \\ &= \frac{e^{(b_j + \mathbf{x}^T \mathbf{w}_{*j})}}{1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}}} \\ &= \frac{1}{1 + e^{-(b_j + \mathbf{x}^T \mathbf{w}_{*j})}}, \end{aligned}$$

and

$$p_{BB}(h_j = 0|\mathbf{x}) = \frac{1}{1 + e^{b_j + \mathbf{x}^T \mathbf{w}_{*j}}}.$$

Conditional probability for visible units

Similarly we have that the conditional probability of the visible units given the hidden are

$$\begin{aligned} p_{BB}(\mathbf{x}|\mathbf{h}) &= \prod_i^M \frac{e^{(a_i + \mathbf{w}_{i*}^T \mathbf{h})x_i}}{1 + e^{a_i + \mathbf{w}_{i*}^T \mathbf{h}}} \\ &= \prod_i^M p_{BB}(x_i|\mathbf{h}). \end{aligned}$$

We have

$$\begin{aligned} p_{BB}(x_i = 1|\mathbf{h}) &= \frac{1}{1 + e^{-(a_i + \mathbf{w}_{i*}^T \mathbf{h})}} \\ p_{BB}(x_i = 0|\mathbf{h}) &= \frac{1}{1 + e^{a_i + \mathbf{w}_{i*}^T \mathbf{h}}}. \end{aligned}$$

Gaussian-Binary Restricted Boltzmann Machines

Inserting into the expression for $E_{RBM}(\mathbf{x}, \mathbf{h}, \Theta)$ in equation results in the energy

$$\begin{aligned} E_{GB}(\mathbf{x}, \mathbf{h}, \Theta) &= \sum_i^M \frac{(x_i - a_i)^2}{2\sigma_i^2} - \sum_j^N b_j h_j - \sum_{ij}^{M,N} \frac{x_i w_{ij} h_j}{\sigma_i^2} \\ &= \left\| \frac{\mathbf{x} - \mathbf{a}}{2\sigma} \right\|^2 - \mathbf{b}^T \mathbf{h} - \left(\frac{\mathbf{x}}{\sigma^2} \right)^T \mathbf{W} \mathbf{h}. \end{aligned}$$

Joint Probability Density Function

$$\begin{aligned} p_{GB}(\mathbf{x}, \mathbf{h}, \Theta) &= \frac{1}{Z_{GB}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \mathbf{h} + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{W} \mathbf{h}} \\ &= \frac{1}{Z_{GB}} e^{-\sum_i^M \frac{(x_i - a_i)^2}{2\sigma_i^2} + \sum_j^N b_j h_j + \sum_{ij}^{M,N} \frac{x_i w_{ij} h_j}{\sigma_i^2}} \\ &= \frac{1}{Z_{GB}} \prod_{ij}^{M,N} e^{-\frac{(x_i - a_i)^2}{2\sigma_i^2} + b_j h_j + \frac{x_i w_{ij} h_j}{\sigma_i^2}}. \end{aligned}$$

Partition function

The partition function is given by

$$Z_{GB} = \int \sum_{\tilde{\mathbf{h}}}^{\tilde{\mathbf{H}}} e^{-\|\frac{\tilde{\mathbf{x}} - \mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \tilde{\mathbf{h}} + (\frac{\tilde{\mathbf{x}}}{\sigma^2})^T \mathbf{W} \tilde{\mathbf{h}}} d\tilde{\mathbf{x}}.$$

Marginal Probability Density Functions

We proceed to find the marginal probability densities of the Gaussian-binary RBM. We first marginalize over the binary hidden units to find $p_{GB}(\mathbf{x})$

$$\begin{aligned} p_{GB}(\mathbf{x}) &= \sum_{\tilde{\mathbf{h}}} p_{GB}(\mathbf{x}, \tilde{\mathbf{h}}) \\ &= \frac{1}{Z_{GB}} \sum_{\tilde{\mathbf{h}}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \tilde{\mathbf{h}} + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{W} \tilde{\mathbf{h}}} \\ &= \frac{1}{Z_{GB}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2} \prod_j^N (1 + e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}). \end{aligned}$$

Then the visible units

We next marginalize over the visible units. This is the first time we marginalize over continuous values. We rewrite the exponential factor dependent on \mathbf{x} as a Gaussian function before we integrate in the last step.

$$\begin{aligned} p_{GB}(\mathbf{h}) &= \int p_{GB}(\tilde{\mathbf{x}}, \mathbf{h}) d\tilde{\mathbf{x}} \\ &= \frac{1}{Z_{GB}} \int e^{-\|\frac{\tilde{\mathbf{x}} - \mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \mathbf{h} + (\frac{\tilde{\mathbf{x}}}{\sigma^2})^T \mathbf{W} \mathbf{h}} d\tilde{\mathbf{x}} \\ &= \frac{1}{Z_{GB}} e^{\mathbf{b}^T \mathbf{h}} \int \prod_i^M e^{-\frac{(\tilde{x}_i - a_i)^2}{2\sigma_i^2} + \frac{\tilde{x}_i \mathbf{w}_{i*}^T \mathbf{h}}{\sigma_i^2}} d\tilde{\mathbf{x}} \\ &= \frac{1}{Z_{GB}} e^{\mathbf{b}^T \mathbf{h}} \left(\int e^{-\frac{(\tilde{x}_1 - a_1)^2}{2\sigma_1^2} + \frac{\tilde{x}_1 \mathbf{w}_{1*}^T \mathbf{h}}{\sigma_1^2}} d\tilde{x}_1 \right. \\ &\quad \times \int e^{-\frac{(\tilde{x}_2 - a_2)^2}{2\sigma_2^2} + \frac{\tilde{x}_2 \mathbf{w}_{2*}^T \mathbf{h}}{\sigma_2^2}} d\tilde{x}_2 \\ &\quad \times \dots \end{aligned}$$

Conditional Probability Density Functions

We finish by deriving the conditional probabilities.

$$\begin{aligned} p_{GB}(\mathbf{h}|\mathbf{x}) &= \frac{p_{GB}(\mathbf{x}, \mathbf{h})}{p_{GB}(\mathbf{x})} \\ &= \frac{\frac{1}{Z_{GB}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \mathbf{h} + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{W} \mathbf{h}}}{\frac{1}{Z_{GB}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2} \prod_j^N (1 + e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}})} \\ &= \prod_j^N \frac{e^{(b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}) h_j}}{1 + e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}} \\ &= \prod_j^N p_{GB}(h_j|\mathbf{x}). \end{aligned}$$

Hidden units

The conditional probability of a binary hidden unit h_j being on or off again takes the form of a sigmoid function

$$\begin{aligned} p_{GB}(h_j = 1|\mathbf{x}) &= \frac{e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}}{1 + e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}} \\ &= \frac{1}{1 + e^{-b_j - (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}} \\ p_{GB}(h_j = 0|\mathbf{x}) &= \frac{1}{1 + e^{b_j + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{w}_{*j}}}. \end{aligned}$$

Visible units

The conditional probability of the continuous \mathbf{x} now has another form, however.

$$\begin{aligned}
 p_{GB}(\mathbf{x}|\mathbf{h}) &= \frac{p_{GB}(\mathbf{x}, \mathbf{h})}{p_{GB}(\mathbf{h})} \\
 &= \frac{\frac{1}{Z_{GB}} e^{-\|\frac{\mathbf{x}-\mathbf{a}}{2\sigma}\|^2 + \mathbf{b}^T \mathbf{h} + (\frac{\mathbf{x}}{\sigma^2})^T \mathbf{W} \mathbf{h}}}{\frac{1}{Z_{GB}} e^{\mathbf{b}^T \mathbf{h}} \prod_i^M \sqrt{2\pi\sigma_i^2} e^{\frac{2a_i \mathbf{w}_{i*}^T \mathbf{h} + (\mathbf{w}_{i*}^T \mathbf{h})^2}{2\sigma_i^2}}} \\
 &= \prod_i^M \frac{1}{\sqrt{2\pi\sigma_i^2}} \frac{e^{-\frac{(x_i - a_i)^2}{2\sigma_i^2} + \frac{x_i \mathbf{w}_{i*}^T \mathbf{h}}{2\sigma_i^2}}}{e^{\frac{2a_i \mathbf{w}_{i*}^T \mathbf{h} + (\mathbf{w}_{i*}^T \mathbf{h})^2}{2\sigma_i^2}}} \\
 &= \prod_i^M \frac{1}{\sqrt{2\pi\sigma_i^2}} \frac{e^{-\frac{x_i^2 - 2a_i x_i + a_i^2 - 2x_i \mathbf{w}_{i*}^T \mathbf{h}}{2\sigma_i^2}}}{e^{\frac{2a_i \mathbf{w}_{i*}^T \mathbf{h} + (\mathbf{w}_{i*}^T \mathbf{h})^2}{2\sigma_i^2}}} \\
 &= \prod_i^M \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{x_i^2 - 2a_i x_i + a_i^2 - 2x_i \mathbf{w}_{i*}^T \mathbf{h} + 2a_i \mathbf{w}_{i*}^T \mathbf{h} + (\mathbf{w}_{i*}^T \mathbf{h})^2}{2\sigma_i^2}}
 \end{aligned}$$

Comments

The form of these conditional probabilities explains the name "Gaussian" and the form of the Gaussian-binary energy function. We see that the conditional probability of x_i given \mathbf{h} is a normal distribution with mean $b_i + \mathbf{w}_{i*}^T \mathbf{h}$ and variance σ_i^2 . For the quantum mechanical calculations however, there are several ingredients which simplify our calculations.

Neural Quantum States

The wavefunction should be a probability amplitude depending on \mathbf{x} . The RBM model is given by the joint distribution of \mathbf{x} and \mathbf{h}

$$F_{rbm}(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} e^{-\frac{1}{T_0} E(\mathbf{x}, \mathbf{h})}$$

To find the marginal distribution of \mathbf{x} we set:

$$\begin{aligned} F_{rbm}(\mathbf{x}) &= \sum_{\mathbf{h}} F_{rbm}(\mathbf{x}, \mathbf{h}) \\ &= \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} \end{aligned}$$

Model for the trial wave function

Now this is what we use to represent the wave function, calling it a neural-network quantum state (NQS)

$$\begin{aligned}\Psi(X) &= F_{rbm}(x) \\ &= \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(x, \mathbf{h})} \\ &= \frac{1}{Z} \sum_{\{h_j\}} e^{-\sum_i^M \frac{(x_i - a_i)^2}{2\sigma^2} + \sum_j^N b_j h_j + \sum_{i,j}^{M,N} \frac{x_i w_{ij} h_j}{\sigma^2}} \\ &= \frac{1}{Z} e^{-\sum_i^M \frac{(x_i - a_i)^2}{2\sigma^2}} \prod_j^N (1 + e^{b_j + \sum_i^M \frac{x_i w_{ij}}{\sigma^2}})\end{aligned}$$

Allowing for complex valued functions

The above wavefunction is the most general one because it allows for complex valued wavefunctions. However it fundamentally changes the probabilistic foundation of the RBM, because what is usually a probability in the RBM framework is now an amplitude. This means that a lot of the theoretical framework usually used to interpret the model, i.e. graphical models, conditional probabilities, and Markov random fields, breaks down.

Squared wave function

If we assume the wavefunction to be positive definite, however, we can use the RBM to represent the squared wavefunction, and thereby a probability. This also makes it possible to sample from the model using Gibbs sampling, because we can obtain the conditional probabilities.

$$\begin{aligned} |\Psi(X)|^2 &= F_{rbm}(X) \\ \Rightarrow \Psi(X) &= \sqrt{F_{rbm}(X)} \\ &= \frac{1}{\sqrt{Z}} \sqrt{\sum_{\{h_j\}} e^{-E(X,h)}} \\ &= \frac{1}{\sqrt{Z}} \sqrt{\sum_{\{h_j\}} e^{-\sum_i^M \frac{(X_i - a_i)^2}{2\sigma^2} + \sum_j^N b_j h_j + \sum_{i,j}^{M,N} \frac{X_i w_{ij} h_j}{\sigma^2}}} \\ &= \frac{1}{\sqrt{Z}} e^{-\sum_i^M \frac{(X_i - a_i)^2}{4\sigma^2}} \sqrt{\sum_{\{h_j\}} \prod_j^N e^{b_j h_j + \sum_i^M \frac{X_i w_{ij} h_j}{\sigma^2}}} \end{aligned}$$

Cost function

This is where we deviate from what is common in machine learning. Rather than defining a cost function based on some dataset, our cost function is the energy of the quantum mechanical system. From the variational principle we know that minimizing this energy should lead to the ground state wavefunction. As stated previously the local energy is given by

$$E_L = \frac{1}{\psi} \hat{H} \psi.$$

And the gradient

$$G_i = \frac{\partial \langle E_L \rangle}{\partial \alpha_i} = 2(\langle E_L \frac{1}{\Psi} \frac{\partial \Psi}{\partial \alpha_i} \rangle - \langle E_L \rangle \langle \frac{1}{\Psi} \frac{\partial \Psi}{\partial \alpha_i} \rangle),$$

where $\alpha_i = a_1, \dots, a_M, b_1, \dots, b_N, w_{11}, \dots, w_{MN}$.

Additional details

We use that $\frac{1}{\Psi} \frac{\partial \Psi}{\partial \alpha_i} = \frac{\partial \ln \Psi}{\partial \alpha_i}$, and find

$$\ln \Psi(\mathbf{X}) = -\ln Z - \sum_m^M \frac{(X_m - a_m)^2}{2\sigma^2} + \sum_n^N \ln(1 + e^{b_n + \sum_i^M \frac{X_i w_{in}}{\sigma^2}}).$$

Final equation

This gives

$$\frac{\partial}{\partial a_m} \ln \Psi = \frac{1}{\sigma^2} (X_m - a_m)$$

$$\frac{\partial}{\partial b_n} \ln \Psi = \frac{1}{e^{-b_n - \frac{1}{\sigma^2} \sum_i^M X_i w_{in}} + 1}$$

$$\frac{\partial}{\partial w_{mn}} \ln \Psi = \frac{X_m}{\sigma^2 (e^{-b_n - \frac{1}{\sigma^2} \sum_i^M X_i w_{in}} + 1)}.$$

Code example

This part is best seen using the jupyter-notebook

```
# 2-electron VMC code for 2dim quantum dot with importance sampling  
# Using gaussian rng for new positions and Metropolis- Hastings  
# Added restricted boltzmann machine method for dealing with the wavef  
# RBM code based heavily off of:  
# https://github.com/CompPhysics/ComputationalPhysics2/tree/gh-pages/d
```

```
from math import exp, sqrt  
from random import random, seed, normalvariate  
import numpy as np  
import matplotlib.pyplot as plt  
from mpl_toolkits.mplot3d import Axes3D  
from matplotlib import cm  
from matplotlib.ticker import LinearLocator, FormatStrFormatter  
import sys
```

```
# Trial wave function for the 2-electron quantum dot in two dims
```

```
def WaveFunction(r,a,b,w):  
    sigma=1.0  
    sig2 = sigma**2  
    Psi1 = 0.0  
    Psi2 = 1.0  
    Q = Qfac(r,b,w)  
  
    for iq in range(NumberParticles):  
        for ix in range(Dimension):  
            Psi1 += (r[iq[ix]]-a[iq[ix]])**2
```

Project 2, VMC for fermions: Efficient calculation of Slater determinants

The potentially most time-consuming part is the evaluation of the gradient and the Laplacian of an N -particle Slater determinant. We have to differentiate the determinant with respect to all spatial coordinates of all particles. A brute force differentiation would involve $N \cdot d$ evaluations of the entire determinant which would even worsen the already undesirable time scaling, making it $Nd \cdot O(N^3) \sim O(d \cdot N^4)$.

This poses serious hindrances to the overall efficiency of our code.

Matrix elements of Slater determinants

The efficiency can be improved however if we move only one electron at the time. The Slater determinant matrix \hat{D} is defined by the matrix elements

$$d_{ij} = \phi_j(x_i)$$

where $\phi_j(r_i)$ is a single particle wave function. The columns correspond to the position of a given particle while the rows stand for the various quantum numbers.

Efficient calculation of Slater determinants

What we need to realize is that when differentiating a Slater determinant with respect to some given coordinate, only one row of the corresponding Slater matrix is changed.

Therefore, by recalculating the whole determinant we risk producing redundant information. The solution turns out to be an algorithm that requires to keep track of the *inverse* of the Slater matrix.

Efficient calculation of Slater determinants

Let the current position in phase space be represented by the $(N \cdot d)$ -element vector r^{old} and the new suggested position by the vector r^{new} .

The inverse of \hat{D} can be expressed in terms of its cofactors C_{ij} and its determinant (this our notation for a determinant) $|\hat{D}|$:

$$d_{ij}^{-1} = \frac{C_{ji}}{|\hat{D}|} \quad (1)$$

Notice that the interchanged indices indicate that the matrix of cofactors is to be transposed.

Efficient calculation of Slater determinants

If \hat{D} is invertible, then we must obviously have $\hat{D}^{-1}\hat{D} = 1$, or explicitly in terms of the individual elements of \hat{D} and \hat{D}^{-1} :

$$\sum_{k=1}^N d_{ik} d_{kj}^{-1} = \delta_{ij} \quad (2)$$

Efficient calculation of Slater determinants

Consider the ratio, which we shall call R , between $|\hat{D}(\mathbf{r}^{\text{new}})|$ and $|\hat{D}(\mathbf{r}^{\text{old}})|$. By definition, each of these determinants can individually be expressed in terms of the i -th row of its cofactor matrix

$$R \equiv \frac{|\hat{D}(\mathbf{r}^{\text{new}})|}{|\hat{D}(\mathbf{r}^{\text{old}})|} = \frac{\sum_{j=1}^N d_{ij}(\mathbf{r}^{\text{new}}) C_{ij}(\mathbf{r}^{\text{new}})}{\sum_{j=1}^N d_{ij}(\mathbf{r}^{\text{old}}) C_{ij}(\mathbf{r}^{\text{old}})} \quad (3)$$

Efficient calculation of Slater determinants

Suppose now that we move only one particle at a time, meaning that \mathbf{r}^{new} differs from \mathbf{r}^{old} by the position of only one, say the i -th, particle. This means that $\hat{D}(\mathbf{r}^{\text{new}})$ and $\hat{D}(\mathbf{r}^{\text{old}})$ differ only by the entries of the i -th row. Recall also that the i -th row of a cofactor matrix \hat{C} is independent of the entries of the i -th row of its corresponding matrix \hat{D} . In this particular case we therefore get that the i -th row of $\hat{C}(\mathbf{r}^{\text{new}})$ and $\hat{C}(\mathbf{r}^{\text{old}})$ must be equal. Explicitly, we have:

$$C_{ij}(\mathbf{r}^{\text{new}}) = C_{ij}(\mathbf{r}^{\text{old}}) \quad \forall j \in \{1, \dots, N\} \quad (4)$$

Efficient calculation of Slater determinants

Inserting this into the numerator of eq. (3) and using eq. (1) to substitute the cofactors with the elements of the inverse matrix, we get:

$$R = \frac{\sum_{j=1}^N d_{ij}(r^{\text{new}}) C_{ij}(r^{\text{old}})}{\sum_{j=1}^N d_{ij}(r^{\text{old}}) C_{ij}(r^{\text{old}})} = \frac{\sum_{j=1}^N d_{ij}(r^{\text{new}}) d_{ji}^{-1}(r^{\text{old}})}{\sum_{j=1}^N d_{ij}(r^{\text{old}}) d_{ji}^{-1}(r^{\text{old}})} \quad (5)$$

Efficient calculation of Slater determinants

Now by eq. (2) the denominator of the rightmost expression must be unity, so that we finally arrive at:

$$R = \sum_{j=1}^N d_{ij}(\mathbf{r}^{\text{new}}) d_{ji}^{-1}(\mathbf{r}^{\text{old}}) = \sum_{j=1}^N \phi_j(\mathbf{r}_i^{\text{new}}) d_{ji}^{-1}(\mathbf{r}^{\text{old}}) \quad (6)$$

What this means is that in order to get the ratio when only the i -th particle has been moved, we only need to calculate the dot product of the vector $(\phi_1(\mathbf{r}_i^{\text{new}}), \dots, \phi_N(\mathbf{r}_i^{\text{new}}))$ of single particle wave functions evaluated at this new position with the i -th column of the inverse matrix \hat{D}^{-1} evaluated at the original position. Such an operation has a time scaling of $O(N)$. The only extra thing we need to do is to maintain the inverse matrix $\hat{D}^{-1}(\mathbf{x}^{\text{old}})$.

Efficient calculation of Slater determinants

If the new position \mathbf{r}^{new} is accepted, then the inverse matrix can be suitably updated by an algorithm having a time scaling of $O(N^2)$. This algorithm goes as follows. First we update all but the i -th column of \hat{D}^{-1} . For each column $j \neq i$, we first calculate the quantity:

$$S_j = (\hat{D}(\mathbf{r}^{\text{new}}) \times \hat{D}^{-1}(\mathbf{r}^{\text{old}}))_{ij} = \sum_{l=1}^N d_{il}(\mathbf{r}^{\text{new}}) d_{lj}^{-1}(\mathbf{r}^{\text{old}}) \quad (7)$$

Efficient calculation of Slater determinants

The new elements of the j -th column of \hat{D}^{-1} are then given by:

$$d_{kj}^{-1}(\mathbf{r}^{\text{new}}) = d_{kj}^{-1}(\mathbf{r}^{\text{old}}) - \frac{S_j}{R} d_{ki}^{-1}(\mathbf{r}^{\text{old}}) \quad \begin{array}{l} \forall \ k \in \{1, \dots, N\} \\ j \neq i \end{array} \quad (8)$$

Efficient calculation of Slater determinants

Finally the elements of the i -th column of \hat{D}^{-1} are updated simply as follows:

$$d_{ki}^{-1}(r^{\text{new}}) = \frac{1}{R} d_{ki}^{-1}(r^{\text{old}}) \quad \forall \quad k \in \{1, \dots, N\} \quad (9)$$

We see from these formulas that the time scaling of an update of \hat{D}^{-1} after changing one row of \hat{D} is $O(N^2)$.

The scheme is also applicable for the calculation of the ratios involving derivatives. It turns out that differentiating the Slater determinant with respect to the coordinates of a single particle r_i changes only the i -th row of the corresponding Slater matrix.

The gradient and the Laplacian

The gradient and the Laplacian can therefore be calculated as follows:

$$\frac{\vec{\nabla}_i |\hat{D}(\mathbf{r})|}{|\hat{D}(\mathbf{r})|} = \sum_{j=1}^N \vec{\nabla}_i d_{ij}(\mathbf{r}) d_{ji}^{-1}(\mathbf{r}) = \sum_{j=1}^N \vec{\nabla}_i \phi_j(\mathbf{r}_i) d_{ji}^{-1}(\mathbf{r})$$

and

$$\frac{\nabla_i^2 |\hat{D}(\mathbf{r})|}{|\hat{D}(\mathbf{r})|} = \sum_{j=1}^N \nabla_i^2 d_{ij}(\mathbf{r}) d_{ji}^{-1}(\mathbf{r}) = \sum_{j=1}^N \nabla_i^2 \phi_j(\mathbf{r}_i) d_{ji}^{-1}(\mathbf{r})$$

How to compute the derivatives of the Slater determinant

Thus, to calculate all the derivatives of the Slater determinant, we only need the derivatives of the single particle wave functions ($\vec{\nabla}_i \phi_j(r_i)$ and $\nabla_i^2 \phi_j(r_i)$) and the elements of the corresponding inverse Slater matrix ($\hat{D}^{-1}(r_i)$). A calculation of a single derivative is by the above result an $O(N)$ operation. Since there are $d \cdot N$ derivatives, the time scaling of the total evaluation becomes $O(d \cdot N^2)$. With an $O(N^2)$ updating algorithm for the inverse matrix, the total scaling is no worse, which is far better than the brute force approach yielding $O(d \cdot N^4)$.

Important note: In most cases you end with closed form expressions for the single-particle wave functions. It is then useful to calculate the various derivatives and make separate functions for them.

The Slater determinant

The Slater determinant takes the form

$$\Phi(r_1, r_2, r_3, r_4, \alpha, \beta, \gamma, \delta) = \frac{1}{\sqrt{4!}} \begin{vmatrix} \psi_{100\uparrow}(r_1) & \psi_{100\uparrow}(r_2) & \psi_{100\uparrow}(r_3) & \psi_{100\uparrow}(r_4) \\ \psi_{100\downarrow}(r_1) & \psi_{100\downarrow}(r_2) & \psi_{100\downarrow}(r_3) & \psi_{100\downarrow}(r_4) \\ \psi_{200\uparrow}(r_1) & \psi_{200\uparrow}(r_2) & \psi_{200\uparrow}(r_3) & \psi_{200\uparrow}(r_4) \\ \psi_{200\downarrow}(r_1) & \psi_{200\downarrow}(r_2) & \psi_{200\downarrow}(r_3) & \psi_{200\downarrow}(r_4) \end{vmatrix}$$

The Slater determinant as written is zero since the spatial wave functions for the spin up and spin down states are equal. But we can rewrite it as the product of two Slater determinants, one for spin up and one for spin down.

Rewriting the Slater determinant

We can rewrite it as

$$\begin{aligned}\Phi(r_1, r_2, r_3, r_4, \alpha, \beta, \gamma, \delta) = & \det \uparrow (1, 2) \det \downarrow (3, 4) - \det \uparrow (1, 3) \det \downarrow (2, 4) \\ & - \det \uparrow (1, 4) \det \downarrow (3, 2) + \det \uparrow (2, 3) \det \downarrow (1, 4) - \det \uparrow (2, 4) \det \downarrow (1, 3) \\ & + \det \uparrow (3, 4) \det \downarrow (1, 2),\end{aligned}$$

where we have defined

$$\det \uparrow (1, 2) = \frac{1}{\sqrt{2}} \begin{vmatrix} \psi_{100\uparrow}(r_1) & \psi_{100\uparrow}(r_2) \\ \psi_{200\uparrow}(r_1) & \psi_{200\uparrow}(r_2) \end{vmatrix},$$

and

$$\det \downarrow (3, 4) = \frac{1}{\sqrt{2}} \begin{vmatrix} \psi_{100\downarrow}(r_3) & \psi_{100\downarrow}(r_4) \\ \psi_{200\downarrow}(r_3) & \psi_{200\downarrow}(r_4) \end{vmatrix}.$$

The total determinant is still zero!

Splitting the Slater determinant

We want to avoid to sum over spin variables, in particular when the interaction does not depend on spin.

It can be shown, see for example Moskowitz and Kalos, [Int. J. Quantum Chem. 20 1107 \(1981\)](#), that for the variational energy we can approximate the Slater determinant as

$$\Phi(r_1, r_2, , r_3, r_4, \alpha, \beta, \gamma, \delta) \propto \det \uparrow (1, 2) \det \downarrow (3, 4),$$

or more generally as

$$\Phi(r_1, r_2, \dots r_N) \propto \det \uparrow \det \downarrow,$$

where we have the Slater determinant as the product of a spin up part involving the number of electrons with spin up only (2 for beryllium and 5 for neon) and a spin down part involving the electrons with spin down.

This ansatz is not antisymmetric under the exchange of electrons with opposite spins but it can be shown (show this) that it gives the same expectation value for the energy as the full Slater

Spin up and spin down parts

We will thus factorize the full determinant $|\hat{D}|$ into two smaller ones, where each can be identified with \uparrow and \downarrow respectively:

$$|\hat{D}| = |\hat{D}|_{\uparrow} \cdot |\hat{D}|_{\downarrow}$$

Factorization

The combined dimensionality of the two smaller determinants equals the dimensionality of the full determinant. Such a factorization is advantageous in that it makes it possible to perform the calculation of the ratio R and the updating of the inverse matrix separately for $|\hat{D}|_{\uparrow}$ and $|\hat{D}|_{\downarrow}$:

$$\frac{|\hat{D}|_{\text{new}}}{|\hat{D}|_{\text{old}}} = \frac{|\hat{D}|_{\uparrow}^{\text{new}}}{|\hat{D}|_{\uparrow}^{\text{old}}} \cdot \frac{|\hat{D}|_{\downarrow}^{\text{new}}}{|\hat{D}|_{\downarrow}^{\text{old}}}$$

This reduces the calculation time by a constant factor. The maximal time reduction happens in a system of equal numbers of \uparrow and \downarrow particles, so that the two factorized determinants are half the size of the original one.

Number of operations

Consider the case of moving only one particle at a time which originally had the following time scaling for one transition:

$$O_R(N) + O_{\text{inverse}}(N^2)$$

For the factorized determinants one of the two determinants is obviously unaffected by the change so that it cancels from the ratio R .

Counting the number of FLOPS

Therefore, only one determinant of size $N/2$ is involved in each calculation of R and update of the inverse matrix. The scaling of each transition then becomes:

$$O_R(N/2) + O_{\text{inverse}}(N^2/4)$$

and the time scaling when the transitions for all N particles are put together:

$$O_R(N^2/2) + O_{\text{inverse}}(N^3/4)$$

which gives the same reduction as in the case of moving all particles at once.

Computation of ratios

Computing the ratios discussed above requires that we maintain the inverse of the Slater matrix evaluated at the current position. Each time a trial position is accepted, the row number i of the Slater matrix changes and updating its inverse has to be carried out.

Getting the inverse of an $N \times N$ matrix by Gaussian elimination has a complexity of order of $\mathcal{O}(N^3)$ operations, a luxury that we cannot afford for each time a particle move is accepted. We will use the expression

$$d_{kj}^{-1}(\mathbf{x}^{\text{new}}) = \begin{cases} d_{kj}^{-1}(\mathbf{x}^{\text{old}}) - \frac{d_{ki}^{-1}(\mathbf{x}^{\text{old}})}{R} \sum_{l=1}^N d_{il}(\mathbf{x}^{\text{new}}) d_{lj}^{-1}(\mathbf{x}^{\text{old}}) & \text{if } j \neq i \\ \frac{d_{ki}^{-1}(\mathbf{x}^{\text{old}})}{R} \sum_{l=1}^N d_{il}(\mathbf{x}^{\text{old}}) d_{lj}^{-1}(\mathbf{x}^{\text{old}}) & \text{if } j = i \end{cases}$$

Scaling properties

This equation scales as $O(N^2)$. The evaluation of the determinant of an $N \times N$ matrix by standard Gaussian elimination requires $O(N^3)$ calculations. As there are Nd independent coordinates we need to evaluate Nd Slater determinants for the gradient (quantum force) and Nd for the Laplacian (kinetic energy). With the updating algorithm we need only to invert the Slater determinant matrix once. This can be done by standard LU decomposition methods.

How to get the determinant

Determining a determinant of an $N \times N$ matrix by standard Gaussian elimination is of the order of $O(N^3)$ calculations. As there are $N \cdot d$ independent coordinates we need to evaluate Nd Slater determinants for the gradient (quantum force) and $N \cdot d$ for the Laplacian (kinetic energy)

With the updating algorithm we need only to invert the Slater determinant matrix once. This is done by calling standard LU decomposition methods.

Expectation value of the kinetic energy

The expectation value of the kinetic energy expressed in atomic units for electron i is

$$\langle \hat{K}_i \rangle = -\frac{1}{2} \frac{\langle \Psi | \nabla_i^2 | \Psi \rangle}{\langle \Psi | \Psi \rangle},$$

$$K_i = -\frac{1}{2} \frac{\nabla_i^2 \Psi}{\Psi}. \quad (10)$$

$$\begin{aligned} \frac{\nabla^2 \Psi}{\Psi} &= \frac{\nabla^2(\Psi_D \Psi_C)}{\Psi_D \Psi_C} = \frac{\nabla \cdot [\nabla(\Psi_D \Psi_C)]}{\Psi_D \Psi_C} = \frac{\nabla \cdot [\Psi_C \nabla \Psi_D + \Psi_D \nabla \Psi_C]}{\Psi_D \Psi_C} \\ &= \frac{\nabla \Psi_C \cdot \nabla \Psi_D + \Psi_C \nabla^2 \Psi_D + \nabla \Psi_D \cdot \nabla \Psi_C + \Psi_D \nabla^2 \Psi_C}{\Psi_D \Psi_C} \end{aligned} \quad (11)$$

$$\frac{\nabla^2 \Psi}{\Psi} = \frac{\nabla^2 \Psi_D}{\Psi_D} + \frac{\nabla^2 \Psi_C}{\Psi_C} + 2 \frac{\nabla \Psi_D}{\Psi_D} \cdot \frac{\nabla \Psi_C}{\Psi_C} \quad (12)$$

Second derivative of the Jastrow factor

The second derivative of the Jastrow factor divided by the Jastrow factor (the way it enters the kinetic energy) is

$$\left[\frac{\nabla^2 \psi_C}{\psi_C} \right]_x = 2 \sum_{k=1}^N \sum_{i=1}^{k-1} \frac{\partial^2 g_{ik}}{\partial x_k^2} + \sum_{k=1}^N \left(\sum_{i=1}^{k-1} \frac{\partial g_{ik}}{\partial x_k} - \sum_{i=k+1}^N \frac{\partial g_{ki}}{\partial x_i} \right)^2$$

Functional form

But we have a simple form for the function, namely

$$\psi_C = \prod_{i < j} \exp f(r_{ij}) = \exp \left\{ \sum_{i < j} \frac{ar_{ij}}{1 + \beta r_{ij}} \right\},$$

and it is easy to see that for particle k we have

$$\frac{\nabla_k^2 \psi_C}{\psi_C} = \sum_{ij \neq k} \frac{(r_k - r_i)(r_k - r_j)}{r_{ki}r_{kj}} f'(r_{ki})f'(r_{kj}) + \sum_{j \neq k} \left(f''(r_{kj}) + \frac{2}{r_{kj}} f'(r_{kj}) \right)$$

Second derivative of the Jastrow factor

Using

$$f(r_{ij}) = \frac{ar_{ij}}{1 + \beta r_{ij}},$$

and $g'(r_{kj}) = dg(r_{kj})/dr_{kj}$ and $g''(r_{kj}) = d^2g(r_{kj})/dr_{kj}^2$ we find that for particle k we have

$$\frac{\nabla_k^2 \psi_C}{\psi_C} = \sum_{ij \neq k} \frac{(r_k - r_i)(r_k - r_j)}{r_{ki}r_{kj}} \frac{a}{(1 + \beta r_{ki})^2} \frac{a}{(1 + \beta r_{kj})^2} + \sum_{j \neq k} \left(\frac{2a}{r_{kj}(1 + \beta r_{kj})} \right)$$

Gradient and Laplacian

The gradient and Laplacian can be calculated as follows:

$$\frac{\nabla_i |\hat{D}(\mathbf{r})|}{|\hat{D}(\mathbf{r})|} = \sum_{j=1}^N \vec{\nabla}_i d_{ij}(\mathbf{r}) d_{ji}^{-1}(\mathbf{r}) = \sum_{j=1}^N \vec{\nabla}_i \phi_j(\mathbf{r}_i) d_{ji}^{-1}(\mathbf{r})$$

and

$$\frac{\nabla_i^2 |\hat{D}(\mathbf{r})|}{|\hat{D}(\mathbf{r})|} = \sum_{j=1}^N \nabla_i^2 d_{ij}(\mathbf{r}) d_{ji}^{-1}(\mathbf{r}) = \sum_{j=1}^N \nabla_i^2 \phi_j(\mathbf{r}_i) d_{ji}^{-1}(\mathbf{r})$$

The gradient for the determinant

The gradient for the determinant is

$$\frac{\nabla_i |\hat{D}(\mathbf{r})|}{|\hat{D}(\mathbf{r})|} = \sum_{j=1}^N \nabla_i d_{ij}(\mathbf{r}) d_{ji}^{-1}(\mathbf{r}) = \sum_{j=1}^N \nabla_i \phi_j(\mathbf{r}_i) d_{ji}^{-1}(\mathbf{r}).$$

Jastrow gradient in quantum force

We have

$$\psi_C = \prod_{i < j} g(r_{ij}) = \exp \left\{ \sum_{i < j} \frac{ar_{ij}}{1 + \beta r_{ij}} \right\},$$

the gradient needed for the quantum force and local energy is easy to compute. We get for particle k

$$\frac{\nabla_k \psi_C}{\psi_C} = \sum_{j \neq k} \frac{r_{kj}}{r_{kj}} \frac{a}{(1 + \beta r_{kj})^2},$$

which is rather easy to code. Remember to sum over all particles when you compute the local energy.

Metropolis Hastings part

We need to compute the ratio between wave functions, in particular for the Slater determinants.

$$R = \sum_{j=1}^N d_{ij}(\mathbf{r}^{\text{new}}) d_{ji}^{-1}(\mathbf{r}^{\text{old}}) = \sum_{j=1}^N \phi_j(\mathbf{r}_i^{\text{new}}) d_{ji}^{-1}(\mathbf{r}^{\text{old}})$$

What this means is that in order to get the ratio when only the i -th particle has been moved, we only need to calculate the dot product of the vector $(\phi_1(\mathbf{r}_i^{\text{new}}), \dots, \phi_N(\mathbf{r}_i^{\text{new}}))$ of single particle wave functions evaluated at this new position with the i -th column of the inverse matrix \hat{D}^{-1} evaluated at the original position. Such an operation has a time scaling of $O(N)$. The only extra thing we need to do is to maintain the inverse matrix $\hat{D}^{-1}(\mathbf{x}^{\text{old}})$.

Proof for updating algorithm for Slater determinant

As a starting point we may consider that each time a new position is suggested in the Metropolis algorithm, a row of the current Slater matrix experiences some kind of perturbation. Hence, the Slater matrix with its orbitals evaluated at the new position equals the old Slater matrix plus a perturbation matrix,

$$d_{jk}(x^{\text{new}}) = d_{jk}(x^{\text{old}}) + \Delta_{jk}, \quad (13)$$

where

$$\Delta_{jk} = \delta_{ik} [\phi_j(x_i^{\text{new}}) - \phi_j(x_i^{\text{old}})] = \delta_{ik} (\Delta\phi)_j. \quad (14)$$

Proof for updating algorithm for Slater determinant

Computing the inverse of the transposed matrix we arrive at

$$d_{kj}(\mathbf{x}^{\text{new}})^{-1} = [d_{kj}(\mathbf{x}^{\text{old}}) + \Delta_{kj}]^{-1}. \quad (15)$$

Proof for updating algorithm for Slater determinant

The evaluation of the right hand side (rhs) term above is carried out by applying the identity $(A + B)^{-1} = A^{-1} - (A + B)^{-1}BA^{-1}$. In compact notation it yields

$$\begin{aligned} [D^T(x^{\text{new}})]^{-1} &= [D^T(x^{\text{old}}) + \Delta^T]^{-1} \\ &= [D^T(x^{\text{old}})]^{-1} - [D^T(x^{\text{old}}) + \Delta^T]^{-1} \Delta^T [D^T(x^{\text{old}})]^{-1} \\ &= [D^T(x^{\text{old}})]^{-1} - \underbrace{[D^T(x^{\text{new}})]^{-1}}_{\text{By Eq.15}} \Delta^T [D^T(x^{\text{old}})]^{-1}. \end{aligned}$$

Proof for updating algorithm for Slater determinant

Using index notation, the last result may be expanded by

$$\begin{aligned}d_{kj}^{-1}(x^{\text{new}}) &= d_{kj}^{-1}(x^{\text{old}}) - \sum_l \sum_m d_{km}^{-1}(x^{\text{new}}) \Delta_{ml}^T d_{lj}^{-1}(x^{\text{old}}) \\&= d_{kj}^{-1}(x^{\text{old}}) - \sum_l \sum_m d_{km}^{-1}(x^{\text{new}}) \Delta_{lm} d_{lj}^{-1}(x^{\text{cur}}) \\&= d_{kj}^{-1}(x^{\text{old}}) - \sum_l \sum_m d_{km}^{-1}(x^{\text{new}}) \delta_{im}(\Delta\phi)_l d_{lj}^{-1}(x^{\text{old}}) \\&= d_{kj}^{-1}(x^{\text{old}}) - d_{ki}^{-1}(x^{\text{new}}) \sum_{l=1}^N (\Delta\phi)_l d_{lj}^{-1}(x^{\text{old}}) \\&= d_{kj}^{-1}(x^{\text{old}}) - d_{ki}^{-1}(x^{\text{new}}) \sum_{l=1}^N [\phi_l(r_i^{\text{new}}) - \phi_l(r_i^{\text{old}})] D_{lj}^{-1}(x^{\text{old}}).\end{aligned}$$

Proof for updating algorithm for Slater determinant

Using

$$D^{-1}(x^{\text{old}}) = \frac{\text{adj}D}{|D(x^{\text{old}})|} \quad \text{and} \quad D^{-1}(x^{\text{new}}) = \frac{\text{adj}D}{|D(x^{\text{new}})|},$$

and dividing these two equations we get

$$\frac{D^{-1}(x^{\text{old}})}{D^{-1}(x^{\text{new}})} = \frac{|D(x^{\text{new}})|}{|D(x^{\text{old}})|} = R \Rightarrow d_{ki}^{-1}(x^{\text{new}}) = \frac{d_{ki}^{-1}(x^{\text{old}})}{R}.$$

Proof for updating algorithm for Slater determinant

We have

$$d_{kj}^{-1}(x^{\text{new}}) = d_{kj}^{-1}(x^{\text{old}}) - \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N [\phi_l(r_i^{\text{new}}) - \phi_l(r_i^{\text{old}})] d_{lj}^{-1}(x^{\text{old}}),$$

or

$$\begin{aligned} d_{kj}^{-1}(x^{\text{new}}) &= d_{kj}^{-1}(x^{\text{old}}) - \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N \phi_l(r_i^{\text{new}}) d_{lj}^{-1}(x^{\text{old}}) \\ &\quad + \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N \phi_l(r_i^{\text{old}}) d_{lj}^{-1}(x^{\text{old}}) \\ &= d_{kj}^{-1}(x^{\text{old}}) - \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N d_{il}(x^{\text{new}}) d_{lj}^{-1}(x^{\text{old}}) \\ &\quad + \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N d_{il}(x^{\text{old}}) d_{lj}^{-1}(x^{\text{old}}). \end{aligned}$$

Proof for updating algorithm for Slater determinant

In this equation, the first line becomes zero for $j = i$ and the second for $j \neq i$. Therefore, the update of the inverse for the new Slater matrix is given by

$$d_{kj}^{-1}(x^{\text{new}}) = \begin{cases} d_{kj}^{-1}(x^{\text{old}}) - \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N d_{il}(x^{\text{new}}) d_{lj}^{-1}(x^{\text{old}}) & \text{if } j \neq i \\ \frac{d_{ki}^{-1}(x^{\text{old}})}{R} \sum_{l=1}^N d_{il}(x^{\text{old}}) d_{lj}^{-1}(x^{\text{old}}) & \text{if } j = i \end{cases}$$