# TUTORIAL: MONTE CARLO METHODS

## I. Simulation of two repulsive particles in a harmonic potential

We consider two particles 1 and 2 of equal mass $m$ in a one-dimensional harmonic potential of stiffness $\kappa$ at temperature $T$. The two particles also interact repulsively with a typical interaction strength $\gamma$. The Hamiltonian of the system reads

$$\mathcal{H} = \frac{p_1^2}{2m} + \frac{p_2^2}{2m} + \frac{1}{2}\kappa\left(x_1^2 + x_2^2\right) - \gamma|x_1 - x_2|. \tag{1}$$

The first two terms stand for the kinetic energy (with $p_a$ the momentum of particle $a$), the third term corresponds to the quadratic potential energy of each particle, the fourth term is the interaction potential between the two particles which depends on their relative distance.

Integrating over the momentum degrees of freedom, the Boltzmann distribution gives the joint probability distribution density of the positions $(x_1, x_2)$ of the two particles:

$$p(x_1, x_2) = \frac{1}{\mathcal{Z}}\exp\left[-\frac{\kappa}{2k_{\mathrm{B}}T}\left(x_1^2 + x_2^2\right) + \frac{\gamma}{k_{\mathrm{B}}T}|x_1 - x_2|\right] \propto \exp\left[-\frac{V(x_1, x_2)}{k_{\mathrm{B}}T}\right], \tag{2}$$

with $\mathcal{Z}$ the normalization constant. After some algebra (try it at home, this is a good statistical mechanics exercise!), we can prove that $\langle x_1 \rangle = \langle x_2 \rangle = 0$ while

$$\langle x_1^2 \rangle = \langle x_2^2 \rangle = \frac{k_{\mathrm{B}}T}{\kappa}\left[1 + \frac{\gamma}{\sqrt{\pi\kappa k_{\mathrm{B}}T}}\frac{\exp\left(-\frac{\gamma^2}{\kappa k_{\mathrm{B}}T}\right)}{1 + \mathrm{erf}\left(\frac{\gamma}{\sqrt{\kappa k_{\mathrm{B}}T}}\right)} + \frac{\gamma^2}{\kappa k_{\mathrm{B}}T}\right]. \tag{3}$$

In particular, when $\gamma = 0$, we recover the equipartition theorem ($\langle x_1^2 \rangle = \langle x_2^2 \rangle = k_{\mathrm{B}}T/\kappa$).

We want to use the Metropolis algorithm to sample this distribution. We propose the following Markov chain:

1. Pick at random one of the two particles.

2. Propose a new trial position $y_a$ for the chosen particle $a$ by randomly shifting its current position $x_a$ by a random number between $-\epsilon$ and $\epsilon$.

3. Accept or reject this trial position following Metropolis procedure.

For the simulations, we set $m = 1$, $\kappa = 1$ and $k_{\mathrm{B}} = 1$ (rescaled units).

**Question 1:** Implement the above Markov chain. You can choose $(x_1, x_2) = (0, 0)$ to initiate the chain.

**Question 2:** We start with $\gamma = 0$.

a. For several values of $T \in [0.2, 2]$ and for several values of $\epsilon \in [1, 10]$, perform $N = 10000$ steps and compute the average acceptance rate of trial moves. For each temperature, plot the acceptance rate as a function of $\epsilon$ and display all curves on the same graph. How should you choose $\epsilon$ such that the acceptance rate is roughly between $20\,\%$ and $50\,\%$ for all temperatures?

b. For several values of $T \in [0.2, 2]$, perform $N = 100000$ steps and plot the energy $V(x_1, x_2)$ as a function of the number of steps. How should you choose the number of steps $N_{\mathrm{b-i}}$ in the burn-in phase?

c. For the same simulations, compute the correlation function of $x_1$ in the equilibrium phase and plot the curve as a function of the number of steps. Estimate the number of steps $N_{\mathrm{corr}}$ over which the samples are correlated. How should you choose the number of steps $N_{\mathrm{eq}}$ in the equilibrium phase (with respect to $N_{\mathrm{corr}}$)?

  **d.** For several values of $T \in [0.2, 2]$, compute the variance of $x_1$ in the equilibrium phase and plot $\langle x_1^2 \rangle$ as a function of $T$. Is the equipartition theorem verified?

**Question 3:** We now study the case $\gamma = 1.0$ (in rescaled units).

  **a.** For $T = 0.2$, perform $N = 10000$ steps and plot $x_1$ and $x_2$ as a function of the number of steps on the same graph. What do you observe? Can you rationalize this behavior? How should you choose the number of steps $N_{\mathrm{b-i}}$ in the burn-in phase?

  **b.** For the same temperature, perform $N = 100000$ steps, compute the correlation function of $x_1$ in the equilibrium phase and plot it as a function of the number of steps. Estimate the number of steps $N_{\mathrm{corr}}$ over which the samples are correlated. How should you choose the number of steps $N_{\mathrm{eq}}$ in the equilibrium phase (with respect to $N_{\mathrm{corr}}$)? Compare with the case $\gamma = 0$.

  **c.** For several values of $T \in [0.2, 2]$, compute the variance of $x_1$ in the equilibrium phase and plot $\langle x_1^2 \rangle$ as a function of $T$. Check that Eq. (3) is verified.

## II. Monte Carlo integration

**Question 1:** Compute the volume of the unit ball (the set of vectors of norm smaller or equal than 1) in $D$ dimensions using a Monte Carlo integration. Compare with the exact results $4\pi/3$ and $\pi^5/120$ for $D = 3$ and $D = 10$ respectively.

**Question 2:** We now want to compute numerically the integral

$$I = \int_0^1 \left( \frac{1}{x^{1/3}} + \frac{x}{10} \right) \mathrm{d}x. \tag{4}$$

Its exact value is $31/20$.

  **a.** Compute the integral using a Monte Carlo integration with the uniform distribution for $N = 100$, $N = 1000$ and $N = 10000$ samples. For each value of $N$, you can do the calculation 20 times and average the result. Plot $|I - 31/20|$ as a function of $N$ in a loglog plot. Comment on the accuracy of the method as a function of $N$ (recall that the accuracy of Monte Carlo integration is of order $N^{-1/2}$).

  **b.** We want to perform another Monte Carlo integration using the probability distribution density $p(x) = 2/(3x^{1/3})$ of support $I = [0, 1[$. Plot the functions $p(x)$ and $f(x) = 1/x^{1/3} + x/10$ on the same graph. Justify the choice of $p(x)$ (importance sampling).

  **c.** Generate random numbers following the distribution $p(x)$ using the inverse transform sampling.

  **d.** From the set of $N$ random numbers obtained in the previous question, compute an estimate of $I$ for $N = 100$, $N = 1000$ and $N = 10000$ samples (you can do the calculation 20 times and again average the result). Plot $|I - 31/20|$ as a function of $N$ in a loglog plot. Comment on the accuracy of the method as a function of $N$ and compare with the uniform distribution.

## III. The best concert tour

A band wants to make a concert tour in France and has selected 13 cities among the biggest in Metropolitan France. The group wants to start from Paris, perform only once in each city and come back to Paris eventually. Because of ecological concerns, the band wants to travel the least number of kilometers. The geographical coordinates of the 13 French cities where they want to perform are listed in the array below. We recall that on Earth, given two locations of coordinates $(\lambda_1, \phi_1)$ and $(\lambda_2, \phi_2)$ (with $\lambda_a$ the longitudes measuring the W/E deviation from the Greenwich Meridian and $\phi_a$ the latitudes measuring the N/S deviation from the Equator), their relative distance reads

$$d_{12} = R \arccos \left[ \cos \phi_1 \cos \phi_2 \cos(\lambda_1 - \lambda_2) + \sin \phi_1 \sin \phi_2 \right], \tag{5}$$

| City | Longitude (O/E) $\lambda$ | Latitude (N/S) $\phi$ |
|------|---------------------------|------------------------|
| Paris | 2° 21′ 07″ (E) | 48° 51′ 24″ (N) |
| Lyon | 4° 49′ 56″ (E) | 45° 45′ 28″ (N) |
| Toulouse | 1° 26′ 38″ (E) | 43° 36′ 16″ (N) |
| Nice | 7° 16′ 17″ (E) | 43° 41′ 45″ (N) |
| Nantes | 1° 33′ 10″ (O) | 47° 13′ 05″ (N) |
| Montpellier | 3° 52′ 38″ (E) | 43° 36′ 43″ (N) |
| Strasbourg | 7° 45′ 08″ (E) | 48° 34′ 24″ (N) |
| Bordeaux | 0° 34′ 46″ (O) | 44° 50′ 16″ (N) |
| Lille | 3° 03′ 48″ (E) | 50° 38′ 14″ (N) |
| Le Havre | 0° 06′ 00″ (E) | 49° 29′ 24″ (N) |
| Clermont-Ferrand | 3° 04′ 56″ (E) | 45° 46′ 33″ (N) |
| Limoges | 1° 15′ 00″ (E) | 45° 51′ 00″ (N) |
| Orléans | 1° 54′ 32″ (E) | 47° 54′ 09″ (N) |

Table 1: **Geographical coordinates of the 13 French cities chosen by a band for their French concert tour.** The longitudes $\lambda$ and latitudes $\phi$ are given in degrees (°), minutes (′) and seconds (″). We recall that $1° = 60′ = 3600″$.

with $R = 6371\,\mathrm{km}$ the average radius of the Earth.

A direct brute-force test of the $12! = 479001600$ possible paths takes about 30 minutes to converge and shows that the smallest tour is: Paris → Orléans → Limoges → Clermont-Ferrand → Lyon → Nice → Montpellier → Toulouse → Bordeaux → Nantes → Le Havre → Lille → Strasbourg → Paris. The objective of this exercise is to find circuits which minimize the total length travelled using a simulated annealing algorithm.

**Question 1:** Start from the Python code `band_circuit.py` and construct the $n \times n$ symmetric matrix $d_{ij}$ of all distances between cities $i$ and $j$ (with $n$ the number of cities).

**Question 2:** Compute the length $L$ of the smallest tour described above.

**Question 3:** We propose the following simulated annealing scheme to find an approximation of the optimal tour. We start from a random circuit connecting all cities once, starting from Paris and coming back to Paris. At each step, we propose the following procedure:

1. Try to swap two cities in the path except Paris and compute the new length of the circuit.

2. Accept and reject the new trial path following the Metropolis criterion assuming that the probability of a path of total length $d$ is $\propto e^{-d/T}$ (with $T$ the temperature).

3. Update the temperature $T$ by decreasing it by a factor $\Delta$ at each step.

The simulated annealing scheme starts with an initial temperature $T_{\mathrm{i}}$ and ends with a final temperature $T_{\mathrm{f}}$. Implement the above procedure.

**Question 4:** The performance of the algorithm depends on the annealing rate $\Delta$.

    **a.** For $T_{\mathrm{i}} = 100$ and $T_{\mathrm{f}} = 10^{-5}$, run the algorithm 20 times for $\Delta = 10^{-1}$, $10^{-2}$, $10^{-3}$, $10^{-4}$, $10^{-5}$ and compute the average length $\tilde{L}(\Delta)$ of the optimal path.

    **b.** Plot $\tilde{L}(\Delta) - L$ as a function of $\Delta$ in a loglog plot. How does the difference scales with $\Delta$?

**Question 4:** For $\Delta = 10^{-4}$, $T_{\mathrm{i}} = 100$ and $T_{\mathrm{f}} = 10^{-5}$, run the simulated annealing algorithm and plot the entire path at the end of the simulated annealing procedure. Is the simulated annealing algorithm able to find the minimal path?

# IV. Parameter inference from experimental data

We are interested in biological data on blood types, see the array below. The observed blood types A, B, AB or O (phenotype) depend on the pair of two alleles of the genotype and on their dominance properties. We want to estimate the probabilities $p_A$, $p_B$ and $p_O$ that an allele is represented in the population (where $p_A + p_B + p_O = 1$).

| Genotype | Probability | Phenotype | Number of observed phenotypes |
|:---:|:---:|:---:|:---:|
| AA | $p_A^2$ | A | $N_A = 186$ |
| AO | $2p_A p_O$ | | |
| BB | $p_B^2$ | B | $N_B = 38$ |
| BO | $2p_B p_O$ | | |
| AB | $2p_A p_B$ | AB | $N_{AB} = 13$ |
| OO | $p_O^2$ | O | $N_O = 284$ |

Table 2: **Results from a biology experiment.** The blood types of $N = 521$ people and their corresponding genotype are reported.

Given the allele probabilities $p_A$, $p_B$, $p_O$, the probability to observe such values of $N_A$, $N_B$, $N_{AB}$ and $N_O$ are given by a multinomial distribution:

$$P(N_A, N_B, N_{AB}, N_O) = \frac{N!}{N_A! N_B! N_{AB}! N_O!} (p_A^2 + 2p_A p_O)^{N_A} (p_B^2 + 2p_B p_O)^{N_B} (2p_A p_B)^{N_{AB}} p_O^{2N_O}, \qquad (6)$$

with $N = N_A + N_B + N_{AB} + N_O$. To estimate the allele probabilities, we propose to maximise the above probability with respect to $p_A$, $p_B$ and $p_O$, which amounts to maximizing the following likelihood

$$L(p_A, p_B, p_O) = N_A \ln(p_A^2 + 2p_A p_O) + N_B \ln(p_B^2 + 2p_B p_O) + N_{AB} \ln(p_A p_B) + 2N_O \ln(p_O) \qquad (7)$$

with the constraint $p_A + p_B + p_O = 1$.
A direct calculation using a Lagrange multiplier to enforce the above constraint leads to tedious algebra and multiple solutions. We instead propose an algorithm to find the maximum. We introduce the two hidden variables $Z_A$ and $Z_B$ corresponding to the number of allele pairs (A,A) and (B,B) respectively. Given $Z_A$ and $Z_B$, the new likelihood to maximize reads

$$\begin{aligned} L'(p_A, p_B, p_O; Z_A, Z_B) = {}& 2Z_A \ln(p_A) + (N_A - Z_A) \ln(2p_A p_O) + 2Z_B \ln(p_B) + (N_B - Z_B) \ln(2p_B p_O) \\ & + N_{AB} \ln(p_A p_B) + 2N_O \ln(p_O). \end{aligned} \qquad (8)$$

By construction $Z_A$ is a random variable which follows a binomial distribution of parameters $N_A$ (number of trials) and $p_A^2/(p_A^2 + 2p_A p_O)$ [relative probability to observe a pair of alleles (A,A)]. Similarly, $Z_B$ follows a binomial distribution of parameters $N_B$ and $p_B^2/(p_B^2 + 2p_B p_O)$.

We propose the following Expectation–Maximization algorithm to find the allele probabilities, starting from random values of $p_A$, $p_B$ and $p_O$.

1. At step $n$, for the estimated values $p_A^{(n)}$, $p_B^{(n)}$ and $p_O^{(n)}$ of the allele probabilities, compute the average of $L'$ with respect to $Z_A$ and $Z_B$. We denote this average $L_{avg}(p_A, p_B, p_O)$.

2. Update the allele probabilities by setting their values at step $n + 1$ to the location of the maximum of the average of $L_{avg}$ computed before:

$$(p_A^{(n+1)}, p_B^{(n+1)}, p_O^{(n+1)}) = \mathrm{argmax}_{(p_A, p_B, p_O)} L_{avg}(p_A, p_B, p_O). \qquad (9)$$

3. Stop the algorithm when you reach fixed points for the allele probabilities, namely, when the relative change of the probabilities between two steps is smaller than a given threshold $\epsilon$. If such a criterion cannot be met, stop the algorithm after a given number of steps.

**Question 1:** Implement the function $L_{\mathrm{avg}}(p_{\mathrm{A}}, p_{\mathrm{B}}, p_{\mathrm{O}})$. You can simplify its expression using pen and paper before coding it, or you can simply generate samples of $Z_{\mathrm{A}}$ and $Z_{\mathrm{B}}$ following binomial distributions.

**Question 2:** Implement the part of the algorithm which finds the location of the maximum of $L_{\mathrm{avg}}$. You can first use pen and paper to find its expression before coding it, or you can directly use optimization routines implemented in Python.

**Question 3:** Find the best estimates of $p_{\mathrm{A}}$, $p_{\mathrm{B}}$ and $p_{\mathrm{O}}$ using the above algorithm.