

## 1 On the average of rotation matrices sharing a common (deterministic) axis

Let  $\mathcal{Q} = \{Q_k\}_{k=1}^N \subset \text{SO}(3)$  an ensemble of rotation matrices sharing a common axis of rotation denoted by  $\mathbf{u}$ . Show that  $\|\langle Q \rangle\| = 1$ , where  $\langle Q \rangle := \frac{1}{N} \sum_{k=1}^N Q_k$  and  $\|A\| = \sup_{\|x\|=1} \|Ax\|$ . Moreover show that if at least one rotation matrix in  $\mathcal{Q}$  has a different rotation axis, then  $\|\langle Q \rangle\| < 1$ .

## 2 Monte Carlo simulation with the cgDNA model II

Download here, [http://lcvwww.epfl.ch/teaching/modelling\\_dna/public\\_files/lambda\\_phage.mat](http://lcvwww.epfl.ch/teaching/modelling_dna/public_files/lambda_phage.mat), the sequences you will work on in this exercise. The five sequences, called `lambda_1...5` have 300 base-pairs and are all part of the  $\lambda$  phage genome (see the details in Sanger et al. 1982).

### 2.1 cgDNA reconstruction of the five lambda phage segments

Using the cgDNA MATLAB package reconstruct the ground state of the five sequences using `cgDNAparamset2`, then compare the 3D reconstruction of all the ground states.

### 2.2 Monte Carlo simulation of the cgDNA model with MATLAB

See Ex. 1.2 Session 6 for details about Monte Carlo simulation with MATLAB. Use only, for brevity, the `cgDNAparamset2`.

Sample 250 configurations for each of the five sequences, and plot each sampled configuration plus the ground state on one figure. Here you can just plot the xyz position of all the base-pairs for each configuration. Are the five plots similar? Compare them to the plot obtained with the six poly-dinucleotides (Ex. 1.1 Session 6).

### 2.3 Cloud of points of last base-pair positions

Use `cgDNAmc` to sample  $10^5$  configurations for each sequence and save for each sample the last base-pair position (see Ex. 1.3 Session 6). Plot using MATLAB the cloud of points, what can you say? Also compare the obtained clouds with the poly-dinucleotide ones from Ex. 1.3 Session 6.

### 2.4 Compute the persistence length using the cgDNAmc code

#### The tangent–tangent correlation

For the following Monte Carlo simulations use  $10^5$  number of samples.

- i) Plot in a single figure all the tangent–tangent correlation (ttc) values for each sequence. Explain the different behaviour and compare the obtained plot with the plot of the ttc values of the six poly-dinucleotides.
- ii) Compute the persistence length for each sequence.

- iii) Define by  $\hat{\mathbf{t}}_1 \cdot \hat{\mathbf{t}}_i$  the intrinsic tangent-tangent for the base pair  $i$ , i.e, the tangent-tangent computed at the  $i$ -th base pair of the ground state. For all the six lambda sequence plot the following

$$\ln\langle \mathbf{t}_1 \cdot \mathbf{t}_i \rangle - \ln \hat{\mathbf{t}}_1 \cdot \hat{\mathbf{t}}_i \text{ vs base pair } i \quad (1)$$

Equation (1) is the so called factorized tangent tangent correlation..

### The Flory vector

For the following Monte Carlo simulations use  $10^5$  number of samples. Plot in a single figure ( use `plot3` and plot only each base-pair position ) the ground state and the Flory persistence vector (F), moreover on the Flory persistence vector plot a cross each 25 base-pairs. What can you say about the position of the crosses? Is the Flory persistence vector converged? ( the convergence here is on the number of base-pairs and not on the number of samples).

## 3 Effect of the sparsity on the Monte Carlo simulation efficiency

Nowadays Monte Carlo simulation for univariate normal distribution can be done in an extremely efficient way. Hence the univariate case is also used to sample from a multivariate normal distribution, but for taking advantage of the univariate case one as to diagonalise the stiffness matrix. Here you have two methods to diagonalise the cgDNA stiffness matrix.

- Use the spectral decomposition to diagonalise the cgDNA stiffness matrix, i.e, write  $K = M\Lambda M^T$ .
- Use the Cholesky factorisation to decompose the cgDNA stiffness matrix, i.e, write  $K = L^T L$ .

1. Using the MATLAB function `mvnrnd` sample 500 configurations using both methods of diagonalisation on  $\text{poly}(AT)_N$  with different  $N$ .
2. The computation of `ttc` and `F` depends only on the value of the inter variables. Thus, formally the expectation of these two quantities is with respect to the marginal distribution of the inter variables. More precisely, let  $x = (y, z) \in \mathbb{R}^{12n-6}$  be a cgDNA configuration where  $y \in \mathbb{R}^{6n}$  are all the intra variables and  $z \in \mathbb{R}^{6(n-1)}$  are all the inter variables. Denote by  $\phi(z)$  a non linear function of the inters. Thus the expected value of  $\phi$  with respect to the cgDNA distribution  $\rho(y, z)$  is

$$\langle \phi(z) \rangle = \int_{\mathbb{R}^{6n}} \int_{\mathbb{R}^{6(n-1)}} \phi(z) \rho(y, z) dz dy = \int_{\mathbb{R}^{6(n-1)}} \phi(z) \rho(z)_y dz = \langle \phi(z) \rangle_z, \quad (2)$$

where  $\langle \cdot \rangle_z$  is the expectation with respect to the inter's marginal distribution  $\rho(z)_y$ . We want now to compare the efficiency of sampling from the whole cgDNA distribution against sampling from the marginal with respect to the intra variables, using Monte Carlo simulation. For doing that consider different repetitions of `poly(AT)` and, as method of decomposition, use the Cholesky factorisation. For this exercise you can fix the number of samples to 500.

**Remark:** The marginal distribution, of a subset of variables, of a Gaussian distribution can be obtained by selecting, in the covariance and in the mean, the corresponding blocks. For example, in the case of a cgDNA distribution the marginal distribution of the inters is obtained by extracting all the inter-inter blocks of the covariance matrix, and the inter part of the mean.