1 Algorithmic details

1.1 Notation

- x_n : Three-dimensional positions of the fine-scale atoms (typically taken from a PDB file). An entire configuration of N fine-scale atoms is stored in the $N \times 3$ matrix $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)'$.
- X_k : Positions of CG beads (pseudo-atoms) where k = 1, ..., K. Typically, we have $N \gg K$. A CG configuration is $X = (X_1, ..., X_K)'$.
- P_k : Conjugate momentum variables for the CG positions that will be used in Hamiltonian Monte Carlo.
- Z_{nk} : Binary assignment variable indicating which of the CG particles represents the n-th fine-scale atom.
- Z: All assignments collected in a $N \times K$ binary matrix.
- s: Error / resolution of the CG representation.
- $E(X; \lambda)$: Force field used to regularize CG configurations.
- λ : Vector of L force field parameters: $\lambda = (\lambda_1, \dots, \lambda_L)'$.
- f(X): Vector of features $f_l(X)$ such that the force field is a linear combination of the features: $E(X; \lambda) = \sum_l \lambda_l f_l(X) = \langle \lambda, f(X) \rangle$.

1.2 Gaussian mixture model

The Gaussian mixture model can be interpreted as a marginal distribution obtained by summing out the assignment variables:

$$\Pr(\boldsymbol{x}_n \mid \boldsymbol{X}, s) = \sum_{Z_{n1} \in \{0, 1\}} \cdots \sum_{Z_{nK} \in \{0, 1\}} \delta\left(\sum_k Z_{nk} - 1\right) \prod_{k=1}^K \left[\frac{1}{K} \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{X}_k, s^2)\right]^{Z_{nk}} . (1)$$

If we use this relation for all atoms, we can write the full likelihood as a marginal distribution over all $N \times K$ assignment variables:

$$\Pr(\boldsymbol{x} \mid \boldsymbol{X}, s) = \sum_{\boldsymbol{Z}} \prod_{n=1}^{N} \delta\left(\sum_{k'} Z_{nk'} - 1\right) \prod_{k=1}^{K} \left[\frac{1}{K} \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{X}_k, s^2)\right]^{Z_{nk}}$$
(2)

where the sum $\sum_{\boldsymbol{Z}}$ is a short-hand for

$$\sum_{\mathbf{Z}} \equiv \sum_{Z_{11} \in \{0,1\}} \cdots \sum_{Z_{1K} \in \{0,1\}} \cdots \sum_{Z_{N1} \in \{0,1\}} \cdots \sum_{Z_{NK} \in \{0,1\}}.$$

The likelihood (2) can be decomposed into a new likelihood which assumes that we know the CG mapping:

$$\Pr(\boldsymbol{x} \mid \boldsymbol{X}, \boldsymbol{Z}, s) = \prod_{n=1}^{N} \prod_{k=1}^{K} \left[\mathcal{N}(\boldsymbol{x}_n; \boldsymbol{X}_k, s^2) \right]^{Z_{nk}}$$
(3)

and a prior distribution over the binary assignment variables:

$$\Pr(\mathbf{Z}) = \prod_{n} \delta\left(\sum_{k} Z_{nk} - 1\right) \prod_{k=1}^{K} K^{-Z_{nk}} = \prod_{n} \mathcal{M}(\mathbf{Z}_{n}; 1, 1/K)$$
(4)

where \mathbf{Z}_n is the *n*-th row of the assignment matrix (a *K*-dimensional binary vector) and $\mathbf{1}$ is a *K*-dimensional vector whose elements are 1. Moreover, we introduced the multinomial distribution in Eq. (4):

$$\mathcal{M}(\boldsymbol{m}; M, \boldsymbol{p}) = \frac{M!}{m_1! \cdots m_K!} \,\delta(M - \sum_k m_k) \prod_k p_k^{m_k} \tag{5}$$

for counts $m_k \in \{0, ..., M\}$ that sum to a total of $M = \sum_k m_k$ and probabilities $p_k \in [0, 1], \sum_k p_k = 1$ represented by a probability vector \boldsymbol{p} .

With these choices we have:

$$Pr(\boldsymbol{x} \mid \boldsymbol{X}, s) = \sum_{\boldsymbol{Z}} Pr(\boldsymbol{x} \mid \boldsymbol{X}, \boldsymbol{Z}, s) \times Pr(\boldsymbol{Z}).$$
 (6)

Instead of summing over all assignment variables analytically, we keep the assignment variables and sample them together with the CG positions using a Gibbs sampler. To see that this is useful, we rewrite the augmented likelihood Pr(x | X, Z, s) (Eq. 3):

$$\Pr(\boldsymbol{x} \mid \boldsymbol{X}, \boldsymbol{Z}, s) = (2\pi s^2)^{-3N/2} \exp\left\{-\frac{1}{2s^2} \sum_{k} N_k \left[\|\boldsymbol{X}_k - \boldsymbol{\mu}_k\|^2 + s_k^2 \right] \right\}$$
(7)

where we introduced the summary statistics

$$N_k = \sum_n Z_{nk}, \quad \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_n Z_{nk} \boldsymbol{x}_n, \quad s_k^2 = \frac{1}{N_k} \sum_n Z_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|^2.$$
 (8)

 N_k is the number of atoms that have been assigned to the k-th bead, μ_k is the center of mass of the assigned fine-scale atoms and s_k the spatial extent of the cluster. The likelihood will pull the beads towards the cluster centers with a harmonic potential whose force constant is N_k/s^2 .

1.3 Inference

We use Markov chain Monte Carlo (MCMC) [1] to infer the model parameters X, Z, s, and λ . The general MCMC strategy is a Gibbs sampler [2], which updates groups of

parameters successively by drawing from the conditional posteriors:

$$X|Z, s, \lambda \sim \Pr(X|Z, s, \lambda, x) \propto \exp\left\{-\frac{1}{2s^2} \sum_k N_k ||X_k - \mu_k||^2 - \langle \lambda, f(X) \rangle\right\}$$
 (9)

$$\boldsymbol{Z}|\boldsymbol{X},s \sim \Pr\{\boldsymbol{Z}|\boldsymbol{X},s,\boldsymbol{x}\} \propto \prod_{n} \mathcal{M}(\boldsymbol{Z}_{n};1,\boldsymbol{p}_{n})$$
 (10)

$$s^{-2}|\boldsymbol{X}, \boldsymbol{Z} \sim \Pr\{s^{-2}|\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{x}\} = \mathcal{G}(s^{-2}; 3N/2, \frac{1}{2} \sum_{k} N_{k}[\|\boldsymbol{\mu}_{k} - \boldsymbol{X}_{k}\|^{2} + s_{k}^{2}])$$
 (11)

$$\lambda | X \sim \Pr{\{\lambda | X\}} \propto \frac{1}{Z(\lambda)} \exp{\{-\langle \lambda, f(X) \rangle\}}$$
 (12)

where we introduced the Gamma distribution

$$\mathcal{G}(t;\alpha,\beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} t^{\alpha-1} e^{-\beta t}, \quad t,\alpha,\beta \ge 0$$
 (13)

and the assignment probabilities

$$p_{nk} = \frac{\exp\{-\frac{1}{2s^2} \|\boldsymbol{x}_n - \boldsymbol{X}_k\|^2\}}{\sum_{k'} \exp\{-\frac{1}{2s^2} \|\boldsymbol{x}_n - \boldsymbol{X}_{k'}\|^2\}}.$$
 (14)

Steps (10) and (11) are straightforward: We simply have to use standard random number generators of the multinomial distribution and of the Gamma distribution to update the CG mapping Z_{nk} and the precision s^{-2} . Sampling X and λ is more challenging.

1.4 Hamiltonian Monte Carlo

The idea of Hamiltonian Monte Carlo (HMC) is to first introduce auxiliary parameters, the momenta P, and simulate the phase-space distribution:

$$\Pr(\boldsymbol{X}, \boldsymbol{P}) \propto \exp\left\{-\frac{1}{2}\|\boldsymbol{P}\|^2 - U(\boldsymbol{X})\right\} = \exp\left\{-H(\boldsymbol{X}, \boldsymbol{P})\right\}$$
(15)

where we introduced the Hamiltonian

$$H(X, P) = \frac{1}{2} ||P||^2 + U(X).$$
 (16)

Because

$$\int \Pr(\boldsymbol{X}, \boldsymbol{P}) d\boldsymbol{P} \propto \Pr(\boldsymbol{X}),$$

we can sample the configurations from Pr(X) by simulating the phase space ensemble $\exp\{-H(X, P)\}$. This is achieved by combining MD simulation with a Metropolis step.

The idea is to start from a current configuration X(0) and generate initial momenta by drawing from a K three-dimensional standard Normal distributions:

$$\mathbf{P}_k(0) \sim \mathcal{N}(\mathbf{P}_k(0); \mathbf{0}, 1). \tag{17}$$

The initial positions $X_k(0)$ could be the cluster centers μ_k defined in Eq. (8) or the previous configuration obtained during Gibbs sampling. We generate a proposal state in phase space $(X(\tau), P(\tau))$ by integrating Hamilton's equations of motion

$$\dot{X} = P, \quad \dot{P} = -\nabla U(X) \tag{18}$$

until an integration time τ . We accept $(X(\tau), P(\tau))$ as a new state with probability

$$\Pr(\boldsymbol{X}(\tau), \boldsymbol{P}(\tau) \mid \boldsymbol{X}(0), \boldsymbol{P}(0)) = \min \left\{ 1, \exp\{H(\boldsymbol{X}(0), \boldsymbol{P}(0)) - H(\boldsymbol{X}(\tau), \boldsymbol{P}(\tau))\} \right\}.$$

If we could integrate Hamilton's equations of motion (18) exactly, we would always accept because the Hamiltonian is conserved under Hamiltonian dynamics. However, for realistic systems we have to use numerical methods to integrate Hamilton's equations of motion

A popular integrator for Hamilton's equations of motion is the *leapfrog* algorithm. It uses a stepsize ϵ and consists of alternating momentum and position update:

$$egin{array}{lll} m{P}(t+\epsilon/2) &=& m{P}(t)-\epsilon/2m{
abla}U(m{X}(t)) \ m{X}(t+\epsilon) &=& m{X}(t)+\epsilonm{P}(t+\epsilon/2) \ m{P}(t+\epsilon) &=& m{P}(t+\epsilon/2)-\epsilon/2m{
abla}U(m{X}(t+\epsilon)) \end{array}$$

1.5 Modeling cryo-EM maps with bead models

It is straight forward to extend our CG graining algorithm to density maps obtained with cryo-EM. Let us assume that the structure for which we would like to find a CG representation is given in the form of a density map with real-valued intensities ρ_n at grid points x_n . Our model for the density at position x is again a mixture of Gaussians centered at the unknown CG positions:

$$\rho(\boldsymbol{x}; \boldsymbol{X}, s) = \sum_{k} \frac{1}{(2\pi s^2)^{3/2}} \exp\left\{-\frac{1}{2s^2} \|\boldsymbol{x} - \boldsymbol{X}_k\|^2\right\}$$

We use the Poisson distribution as the likelihood of observing the experimental map:

$$\Pr(\rho \mid \boldsymbol{X}, s, \alpha) = \prod_{n} [\alpha \rho(\boldsymbol{x}_n; \boldsymbol{X}, s)]^{\rho_n} e^{-\alpha \rho(\boldsymbol{x}_n; \boldsymbol{X}, s)} / (\rho_n!)$$

where $\alpha > 0$ is an unknown scaling parameter. The resulting log likelihood is:

$$\log \Pr(\rho \mid \boldsymbol{X}, s, \alpha) = \sum_{n} \rho_{n} \log[\alpha \rho(\boldsymbol{x}_{n}; \boldsymbol{X}, s)] - \alpha \rho(\boldsymbol{x}_{n}; \boldsymbol{X}, s)$$

$$\approx \sum_{n} \rho_{n} \log \rho(\boldsymbol{x}_{n}; \boldsymbol{X}, s) + \log \alpha \sum_{n} \rho_{n} - K\alpha$$
(19)

where we assumed $\sum_{n} \rho(\mathbf{x}_n; \mathbf{X}, s) \approx \int \rho(\mathbf{x}; \mathbf{X}, s) d\mathbf{x} = K$. Let us simplify the \mathbf{X} dependent part of the log likelihood:

$$\sum_{n} \rho_{n} \log \rho(\boldsymbol{x}_{n}; \boldsymbol{X}, s) = \sum_{n} \rho_{n} \log \sum_{k} \frac{1}{(2\pi s^{2})^{3/2}} \exp \left\{ -\frac{1}{2} \|\boldsymbol{x}_{n} - \boldsymbol{X}_{k}\|^{2} / s^{2} \right\}$$

$$\geq -\frac{1}{2} \sum_{n} \sum_{k} \rho_{n} Z_{nk} \left[\|\boldsymbol{x}_{n} - \boldsymbol{X}_{k}\|^{2} / s^{2} + 3 \log(2\pi s^{2}) \right]$$

where we introduced the assignment variables:

$$Z_{nk} = \frac{\exp\{-\frac{1}{2s^2} \|\boldsymbol{x}_n - \boldsymbol{X}_k\|^2\}}{\sum_{k'} \exp\{-\frac{1}{2s^2} \|\boldsymbol{x}_n - \boldsymbol{X}_{k'}\|^2\}}$$

and used Jensen's inequality to derive a lower bound of the log likelihood. Maximization of the lower bound results in an estimate of the CG positions:

$$\boldsymbol{\mu}_k = \frac{1}{\sum_n \rho_n Z_{nk}} \sum_n \rho_n Z_{nk} \, \boldsymbol{x}_n \tag{20}$$

and of the precision of the representation:

$$s^{2} = \frac{1}{3\sum_{n} \rho_{n}} \sum_{n,k} \rho_{n} Z_{nk} \|\boldsymbol{x}_{n} - \boldsymbol{X}_{k}\|^{2}.$$
 (21)

The CG energy has the same functional form as before (Eq. ??):

$$\frac{1}{2s^2} \sum_{k} N_k \|\boldsymbol{\mu}_k - \boldsymbol{X}_k\|^2, \quad N_k = \sum_{n} \rho_n Z_{nk}$$

where the only difference is that the fine-grained positions x_n are weighted by their density values ρ_n .

References

- [1] J. S. Liu. Monte Carlo strategies in scientific computing. Springer, 2001.
- [2] S. Geman and D. Geman. Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. *IEEE Trans. PAMI*, 6(6):721–741, 1984.