

Best Practices for Property Prediction from Molecular Simulations

Bryce C. Manubay,^{1,*} John D. Chodera,^{2,†} and Michael R. Shirts^{1,‡}

¹University of Colorado

²Computational Biology Program, Sloan Kettering Institute,
Memorial Sloan Kettering Cancer Center, New York, NY 10065, United States

(Dated: July 4, 2016)

This document describes a collected set of best practices for computing various physical properties from molecular simulations of liquid mixtures.

Keywords: best practices; molecular dynamics simulation; physical property computation

TODO LIST

[JDC1]: This is a TODO item example.

I. PRELIMINARIES

Definitions

- V : Volume
- U : Total energy (including potential and kinetic, excluding external energy such as due to gravity, etc)
- S : Entropy
- N : Number of particles
- T : Temperature
- P : Pressure
- k_B : Boltzmann constant
- $\beta: (k_B T)^{-1}$
- M : Molar mass
- ρ : Density (M/V)
- H : Enthalpy
- G : Gibbs Free Energy (free enthalpy)
- A : Helmholtz Free Energy
- μ : Chemical potential
- D : Total dipole moment
- u : reduced energy
- f : reduced free energy

Macroscopically, the quantities V , U , N are constants (assuming the system is not perturbed in any way), as we assume that the fluctuations are essentially zero, and any uncertainty comes from our inability to measure that constant precisely. For a mole of compound (about 18 mL for water), the relative uncertainty in any of these quantities is about 10^{-12} , far lower than any thermodynamics experiment.

[JDC1]: This is a TODO item example.

However, in a molecular simulation, these quantities are not necessarily constant. For example, in a NVT simulation, U is allowed to vary. For a long enough simulation (assuming ergodicity, which can pretty much always be assumed with correct simulations and simple fluids), then the ensemble average value of $U = \langle U \rangle$ will be constant, and in the limit of large simulations/long time will converge to the macroscopic value U ; at least, the macroscopic value of that given model, though perhaps not the U for the real system. In an NVT simulation, clearly V is constant. In a NPT simulation, however, V is a variable, and we must estimate what the macroscopic value would be with the ensemble estimate $\langle V \rangle$.

The quantities T , P , and μ are *always* constants in both simulation and experiment. There are a number of quantities that can be used to ESTIMATE these constants. For example, $\langle \frac{1}{3Nk_B} \sum_i m_i |v_i|^2 \rangle$, where m is the mass of each particle and $|v_i|$ is the magnitude of the velocity of each particle, is an estimator of T , and its average will be equal to the temperature. But it is not the temperature. This quantity fluctuates, but the temperature remains constant; otherwise the simulation could not be at constant temperature.

Ensemble averages of some quantity X ($\langle X \rangle$) are assumed to be averages over the appropriate Boltzmann weighting, i.e. in the NVT ensemble with classical statistical mechanics, they would be $\int X(\vec{x}, \vec{p}) e^{-\beta U(\vec{x}, \vec{p})} d\vec{x} d\vec{p}$. We note that in the limit of very large systems, $\langle X \rangle_{NPT} = \langle X \rangle_{NVT} = \langle X \rangle_{\mu VT}$.

Ensemble averages can be computed by one of two ways. First, they can be computed directly, by running a simulation that produces samples with the desired Boltzmann distribution. In that case ensemble averages can be computed as simple averages, $\langle V \rangle = \frac{1}{N} \sum_i V_i$, where the sum is over all observations. Uncertainties can be estimated in a number of different ways, but usually require estimating the number of uncorrelated samples. Secondly, they can be calculated as reweighted estimates from several different sim-

* email@email.com

† john.chodera@choderalab.org

‡ Corresponding author; michael.shirts@virginia.edu

ulations, as $\langle V \rangle = \frac{1}{\sum_i w_i} V_i w_i$ where w_i is a reweighting factor that can be derived from importance sampling theory.

To simplify our discussion of reweighting, we use some additional notation. We define the reduced potential $u = \beta U(\vec{x})$ in the canonical (NVT) ensemble, $u = \beta U + \beta PV$ in the isobaric-isothermal (NPT) ensemble, and $u = \beta U - \beta N\mu$ in the grand canonical ensemble (similar potentials can be defined in other ensembles). We then define $f = \int e^{-u} dx$, where the integral is over all of the DOF of the system (x for NVT, x, V for NPT, and x, N for μ VT). For NPT, we then have $f = \beta G$, and for NVT we have $f = \beta A$, while for μVT we have $f = -\beta \langle P \rangle V$.

To calculate expectations at one set of parameters generated with parameters that give rise to a different set of probability distributions, we start with the definition of an ensemble average given a probability distribution $p_i(x)$.

$$\langle X \rangle_i = \int X(x) p_i(x) dx \quad (1)$$

We then multiply and divide by $p_j(x)$, to get

$$\langle X \rangle_i = \int X(x) p_i(x) \frac{p_j(x)}{p_j(x)} dx = \int X(x) p_j(x) \frac{p_i(x)}{p_j(x)} dx \quad (2)$$

We then note that this last integral can be estimated by the Monte Carlo estimate

$$\langle X \rangle_i = \int X(x) p_j(x) \frac{p_i(x)}{p_j(x)} dx = \frac{1}{N} \sum_{n=1}^N X(x_n) \frac{p_i(x_n)}{p_j(x_n)} \quad (3)$$

Where the x_k are sampled from probability distribution $p_j(x)$

We now define the mixture distribution of K other distributions as: $p_m(x) = \frac{1}{N} \sum_{i=1}^N N_k p_k(x)$, where $N = \sum_k N_k$. We can construct a sample from the mixture distribution by simply pooling all the samples from k individual simulations. The formula for calculating ensemble averages in a distribution $p_i(x)$ from samples from the mixture distribution is:

$$\langle X \rangle_i = \sum_{n=1}^N X(x_n) \frac{p_i(x_n)}{\sum_{k=1}^{N_k} p_k(x_n)} \quad (4)$$

In the case of Boltzmann averages, then $p_i(x) = e^{f_i - u_i(x)}$, where the reduced free energy f is unknown. Reweighting from the mixture distribution becomes.

$$\langle X \rangle_i = \sum_{n=1}^N X(x_n) \frac{e^{f_i - u_i(x_n)}}{\sum_{k=1}^{N_k} e^{f_k - u_k(x_n)}} \quad (5)$$

which can be seen to be the same formula as the MBAR formula for expectations. The free energies can be obtained by setting $X=1$, and looking at the K equations obtained by reweighting to the K different distributions.

Finite differences at different temperatures and pressures can be calculated by including states with different reduced

potentials. For example, $u_j(x) = \beta_i U(x) + \beta_i (P_i + \Delta P) V$, or $u_j = \frac{1}{k_B(T_i + \Delta T)} U(x) + \frac{1}{k_B(T_i + \Delta T)} P_i V$. However, the relationship between f and G can be problematic when looking at differences in free energy with respect to temperature, because $G_2 - G_1 = \beta_2 f_2 - \beta_1 f_1$. [MRS: needs to find notes on how this was dealt with last time]

Since with MBAR, one can make the differences as small as one would like (you don't have to actually carry out a simulation at those points), we can use the simplest formulas: central difference for first derivatives:

$$\frac{dA}{dx} \approx \frac{1}{2\Delta x} (A(x + \Delta x) - A(x - \Delta x))$$

And for 2nd derivatives:

$$\frac{d^2 A}{dx^2} \approx \frac{1}{\Delta x^2} (A(x + \Delta x) - 2A(x) + A(x - \Delta x))$$

Thus, only properties at two additional points need to be evaluated to calculate both first and 2nd derivatives.

Note that if the finite differences are reevaluated using reweighting approaches, it is important that the simulation used generates the correct Boltzmann distribution. If not, reweighted observables will be incorrect, and the results of the finite difference approach will have significant error.

II. PURE SOLVENT PROPERTIES

A. Density

1. Direct calculation

Starting with the equation used to calculate the density experimentally,

$$\rho = \frac{M}{V} \quad (6)$$

We replace the average with the ensemble estimate (calculated either directly, or with reweighting) to obtain:

$$\rho = \frac{M}{\langle V \rangle} \quad (7)$$

2. Derivative Estimate

From the differential definition of the Gibbs free energy $dG = VdP - SdT + \sum_i \mu_i dN_i$ that V can be calculated from the Gibbs free energy as:

$$V = \left(\frac{\partial G}{\partial P} \right)_{T,N} \quad (8)$$

The density can therefore be estimated from the Gibbs free energy.

$$\rho = \frac{M}{\left(\frac{\partial G}{\partial P}\right)_{T,N}} \quad (9)$$

The derivative can be estimated using a central difference numerical method utilizing Gibbs free energies reweighted to different pressures.

$$\left(\frac{\partial G}{\partial P}\right)_{T,N} \approx \frac{G_{P+\Delta P} - G_{P-\Delta P}}{2\Delta P} \quad (10)$$

The density can then finally be estimated.

$$\rho \approx \frac{M}{\frac{G_{P+\Delta P} - G_{P-\Delta P}}{2\Delta P}} \quad (11)$$

This can be calculated from the reduced free energy f if desired by simply substituting:

$$\rho \approx \frac{\beta M}{\frac{f_{P+\Delta P} - f_{P-\Delta P}}{2\Delta P}} \quad (12)$$

B. Molar Enthalpy

Section on relation of enthalpy to Gibbs free energy (should we need it). This is not an experimental quantity, but will be helpful in calculating related properties of interest. The enthalpy, H , can be found from the Gibbs free energy, G , by the Gibbs-Helmholtz relation:

$$H = -T^2 \left(\frac{\partial \left(\frac{G}{T} \right)}{\partial T} \right)_{P,N} \quad (13)$$

Transforming the derivative in the Gibbs-Helmholtz relation to be in terms of β instead of T yields:

$$H = -T^2 \frac{\beta^2}{\beta^2} \left(\frac{\partial \left(\frac{G}{T} \right)}{\partial T} \frac{\partial T}{\partial \beta} \frac{\partial \beta}{\partial T} \right)_{P,N} \quad (14)$$

Recall that $\beta = \frac{1}{k_B T}$, therefore $\frac{\partial \beta}{\partial T} = -\frac{1}{k_B T^2}$. Substituting these values into the enthalpy equation gives:

$$\begin{aligned} H &= \frac{1}{k_B^3 T^2 \beta^2} \left(\frac{\partial \left(\frac{G}{T} \right)}{\partial \beta} \right)_{P,N} = \frac{1}{k_B} \left(\frac{\partial \left(\frac{G}{T} \right)}{\partial \beta} \right)_{P,N} \\ &= \frac{\partial f}{\partial \beta}_{P,N} \end{aligned} \quad (15)$$

C. Heat Capacity

The definition of the isobaric heat capacity is:

$$C_P = \left(\frac{\partial H}{\partial T} \right)_{P,N} \quad (16)$$

$$= \frac{\partial \left(\frac{\partial f}{\partial \beta} \right)}{\partial T}_{P,N} \quad (17)$$

$$= k_B \beta^2 \frac{\partial^2 f}{\partial \beta^2} \quad (18)$$

This could be computed by finite differences approach or analytical derivation using MBAR.

The enthalpy fluctuation formula can also be used to calculate C_P [1],

$$C_P = \frac{\langle H^2 \rangle - \langle H \rangle^2}{N k_B \langle T \rangle^2}. \quad (19)$$

This form is equivalent for isochoric heat capacity, but with derivatives at constant volume rather than pressure.

Horn et al.[1] suggest a number of vibrational corrections be applied to the calculation of C_P due to a number of approximations made during the simulation of the liquid [1]. The following terms were added as a correction:

$$\begin{aligned} \left(\frac{\partial E_{vib,l}}{\partial T} \right)_P &= \left(\frac{\partial E_{vib,l,intra}^{QM}}{\partial T} \right)_P + \left(\frac{\partial E_{vib,l,inter}^{QM}}{\partial T} \right)_P \\ &\quad - \left(\frac{\partial E_{vib,l,inter}^{CM}}{\partial T} \right)_P \end{aligned} \quad (20)$$

where

$$\left(\frac{\partial E_{vib}^{CM}}{\partial T} \right)_P = k_B n_{vib} \quad (21)$$

and

$$\left(\frac{\partial E_{vib}^{QM}}{\partial T} \right)_P = \sum_{i=1}^{n_{vib}} \left(\frac{h^2 v_i^2 e^{\frac{h v_i}{k_B T}}}{k_B T^2 \left(e^{\frac{h v_i}{k_B T}} - 1 \right)^2} \right) \quad (22)$$

Above, n_{vib} is the number of vibrational modes, h is Planck's constant and v_i is the vibrational frequency of mode i .

D. Isothermal Compressibility

The definition of isothermal compressibility is:

$$\kappa_T = -\frac{1}{V} \left(\frac{\partial V}{\partial P} \right)_T \quad (23)$$

1. First Derivative

Thus, it can be estimated by the finite difference of $\langle V \rangle$

$$\kappa_T = -\frac{1}{2V(T, P)^2} (\langle V(P + \Delta P, T) \rangle - \langle V(P - \Delta P, T) \rangle) \quad (24)$$

Or by the finite differences evaluation of:

$$\kappa_T = -\frac{\left(\frac{\partial^2 G}{\partial P^2}\right)_{T,N}}{\left(\frac{\partial G}{\partial P}\right)_{T,N}} \quad (25)$$

κ_T can also be estimated from the ensemble average and fluctuation of volume (in the NPT ensemble) or particle number (in the μVT ensemble)[2]:

$$\kappa_T = \beta \frac{\langle \Delta V^2 \rangle_{NTP}}{\langle V \rangle_{NTP}} = V\beta \frac{\langle \Delta N^2 \rangle_{VT}}{\langle N \rangle_{VT}} \quad (26)$$

$$\left(\frac{\partial S}{\partial V}\right)_P = \left(\frac{\partial S}{\partial T}\right)_P \left(\frac{\partial T}{\partial V}\right)_P = \frac{C_P}{T} \left(\frac{\partial T}{\partial V}\right)_P = \frac{C_P}{TV\alpha} \quad (31)$$

Where $\alpha = \frac{1}{V} \left(\frac{\partial V}{\partial T}\right)_P = \left(\frac{\partial \ln V}{\partial T}\right)_P$ is the coefficient of thermal expansion. The second term is our triple product rule $\left(\frac{\partial S}{\partial P}\right)_V$ can be expressed as follows.

$$\left(\frac{\partial S}{\partial P}\right)_V = \left(\frac{\partial S}{\partial T}\right)_V \left(\frac{\partial T}{\partial P}\right)_V = \frac{C_V}{T} \left(\frac{\partial T}{\partial P}\right)_V = \frac{C_V}{T\gamma_V} \quad (32)$$

Thus our derivation yields:

$$\left(\frac{\partial P}{\partial V}\right)_S = \frac{C_P\gamma_V}{C_VV\alpha} \quad (33)$$

Horn et al set out several ways for calculating α [1].

E. Speed of Sound

The definition of the speed of sound is[3]:

$$c^2 = \left(\frac{\partial P}{\partial \rho}\right)_S = -\frac{V^2}{M} \left(\frac{\partial P}{\partial V}\right)_S \quad (27)$$

$$c^2 = \frac{V^2}{\beta M} \left[\left(\frac{\gamma_V}{k_B}\right)^2 + \frac{\beta}{V\kappa_T} \right] \quad (28)$$

Where:

$$\gamma_V = \left(\frac{\partial P}{\partial T}\right)_V \quad (29)$$

γ_V is known as the isochoric pressure coefficient. κ_T is the same isothermal compressibility from section A.1.3

An alternate derivation, applying the triple product rule to $\left(\frac{\partial P}{\partial V}\right)_S$ yields the following.

$$\left(\frac{\partial P}{\partial V}\right)_S = \frac{\left(\frac{\partial S}{\partial V}\right)_P}{\left(\frac{\partial S}{\partial P}\right)_V} \quad (30)$$

1. Analytical derivative of density with respect to temperature

$$\alpha = -\frac{d \ln \langle \rho \rangle}{dT} \quad (34)$$

2. Numerical derivative of density over range of T of interest

Same finite differences approach on ρ can be used as was shown for isothermal compressibility.

3. Using the enthalpy-volume fluctuation formula

$$\alpha = \frac{\langle VH \rangle - \langle V \rangle \langle H \rangle}{k_B \langle T \rangle^2 \langle V \rangle} \quad (35)$$

Finite differences approximations and/or analytical derivation can also be used to calculate γ_V .

F. Enthalpy of Vaporization

The definition of the enthalpy of vaporization is[4]:

$$\Delta H_{vap} = H_{gas} - H_{liq} = E_{gas} - E_{liq} + P(V_{gas} - V_{liq}) \quad (36)$$

223
224

225 If we assume that $V_{gas} \gg V_{liq}$ and that the gas is ideal
226 (and therefore kinetic energy terms cancel):

$$\Delta H_{vap} = E_{gas,potential} - E_{liq,potential} + RT \quad (37)$$

227
228

229 An alternate, but similar, method is recommended by
230 Horn et al [1].

$$\Delta H_{vap} = -\frac{E_{liq,potential}}{N} + RT - PV_{liq} + C \quad (38)$$

231

232 In the above equation C is a correction factor for vibra-
233 tional energies, polarizability, non-ideality of the gas and
234 pressure. It can be calculated as follows.
235

$$C_{vib} = C_{vib,intra} + C_{vib,inter} = (E_{vib,QM,gas,intra} - E_{vib,QM,liq,intra}) + (E_{vib,QM,liq,inter} - E_{vib,CM,liq,inter}) \quad (39)$$

236
237

238 The QM and CM subscripts stand for quantum and clas-
239 sical mechanics, respectively.

$$C_{pol} = \frac{N}{2} \frac{(d_{gas} - d_{liq})^2}{\alpha_{p,gas}} \quad (40)$$

240
241

242 Where d_i is the dipole moment of a molecule in phase i
243 and $\alpha_{p,gas}$ is the mean polarizability of a molecule in the gas
244 phase.

$$C_{ni} = P_{vap} \left(B - T \frac{dB}{dT} \right) \quad (41)$$

245
246

247 Where B is the second virial coefficient.

$$C_x = \int_{P_{ext}}^{P_{vap}} [V(P_{ext}) [1 - (P - P_{ext}) \kappa_T] - TV\alpha] dP \quad (42)$$

248
249

250 Where P_{ext} is the external pressure and $V(P_{ext})$ is the
251 volume at P_{ext} .

252 This is frequently done as a single simulation calculation
253 by assuming the average intramolecular energies remains
254 constant during the phase change, which is rigorously cor-
255 rect for something like a rigid water molecule (intramolecu-
256 lar energies are zero), but less true for something with struc-
257 tural rearrangement between gas and liquid phases.

258 As discussed by myself and MRS, we have decided to not
259 initially begin the parametrization process using enthalpy of
260 vaporization data. While force field parametrization is com-
261 monly done using said property we have ample reason to
262 not follow classical practice. First of all, the enthalpy data is
263 usually not collected at standard temperature and pressure,
264 but at the saturation conditions of the liquid being vapor-

265 ized [5]. This would require corrections to be made to get the
266 property at STP (the process will be explained below) using
267 fitted equations for heat capacity. Not only is this inconve-
268 nient, but it adds an unknown complexity to correcting un-
269 certainties in the experimental data. Often times the uncer-
270 tainties of these "experimental" enthalpies are unrecorded
271 because they are estimated from fitted Antoine equation co-
272 efficients [5].

273 An additional issue is the necessity of having to use gas
274 phase simulation data in order to validate a parametriza-
275 tion process meant for small organic liquids and their mix-
276 tures. Following an example of Wang et al. [6] we plan to in-
277 stead use enthalpy of vaporization calculations as an unbi-
278 ased means of testing the success of the parametrization. If
279 the parametrization procedure is expanded to use enthalpy
280 of vaporization, corrections can be made to the experimen-
281 tal heat of vaporization in order to get a value at STP using
282 the following equation.

$$\Delta H_{vap}(T) = \Delta H_{vap}^{ref} + \int_{T_{ref}}^T (C_{P,gas} - C_{P,liq}) dT \quad (43)$$

283
284

1. Dielectric Constant

285 This equation was provided by a literature reference au-
286 thored by CJ Fennell[7]. Below, $\epsilon(0)$ is the zero frequency
287 dielectric constant, V is the system volume and D is the to-
288 tal system dipole moment.
289

$$\epsilon(0) = 1 + \frac{4\pi}{3k_B T \langle V \rangle} (\langle D^2 \rangle - \langle D \rangle^2) \quad (44)$$

III. BINARY MIXTURE PROPERTIES

A. Mass Density, Speed of Sound and Dielectric Constant

The methods for these calculations are the same for a multicomponent system.

1. Activity Coefficient

The definition of chemical potential in a pure substance is:

$$\mu(T, P) = \left(\frac{\partial G}{\partial N} \right)_{T, P} \quad (45)$$

which is a function of only temperature and pressure.

Then the definition of the chemical potential μ_i of component i in a mixture is:

$$\mu_i(T, P, \vec{N}) = \left(\frac{\partial G}{\partial N_i} \right)_{T, P, N_{j \neq i}} \quad (46)$$

N_i refers to a molecule of component i and $N_{j \neq i}$ refers to all molecules other than component i , with \vec{N} the vector of all component numbers. Since μ_i is intensive, this is equivalently a function of the vector of mole fractions \vec{x}_i instead of simply of N_i .

For an ideal solution, the chemical potential μ_i can be related to the pure chemical potential by

$$\mu_i(T, P, \vec{x}_i) = \mu(T, P) + k_B T \ln(\gamma_i) \quad (47)$$

By analogy to this form, we can

$$\mu_i(T, P, \vec{x}_i) = \mu(T, P) + k_B T \ln(x_i \gamma_i) \quad (48)$$

Where γ_i is the activity coefficient of component i , and is a function of T, P , and \vec{x}_i . Rearrangement of the previous equation yields:

$$\gamma_i = \frac{e^{\left(\frac{\mu_i(T, P, \vec{x}_i) - \mu(T, P)}{k_B T} \right)}}{x_i} \quad (49)$$

Although chemical potentials cannot be directly calculated from simulation, chemical potential differences can.

We can calculate the difference $\mu_i(T, P, \vec{x}_i) - \mu(T, P)$ by calculating $\Delta\mu(T, P)_{liquid} - \Delta\mu(T, P)_{gas}$ using a standard alchemical simulation of the pure substance, followed by the calculation of $\mu_i(T, P, \vec{x}_i)_{liquid} - \Delta\mu(T, P, \vec{x}_i)_{gas}$, and assuming that $\Delta\mu(T, P, \vec{x}_i)_{gas} = \Delta\mu(T, P)_{gas}$ (note: there are a few subtleties here relating to the $\ln x_i$ factor, but it appears that with alchemical simulations with a only one particle that is allowed to change, this will cancel out (need to follow up)).

2. Excess Molar Properties

The general definition of an excess molar property can be stated as follows:

$$y^E = y^M - \sum_i x_i y_i \quad (50)$$

Where y^E is the excess molar quantity, y^M is the mixture quantity, x_i is the mole fraction of component i in the mixture and y_i is the pure solvent quantity. In general, the simplest methods for calculating excess molar properties for binary mixtures will require three simulations. One simulation is run for each pure component and a third will be run for the specific mixture of interest. We note that only one set of pure simulations are needed to calculate excess properties at all compositions.

3. Excess Molar Heat Capacity and Volume

Excess molar heat capacities and volume will be calculated using the methods for the pure quantities in section A.1 in combination with the general method for excess property calculation above.

4. Excess Molar Enthalpy

Excess molar enthalpy can be calculated using the general relation of molar enthalpy as it relates to Gibbs Free Energy from section A.1 and the generalized method of excess molar property calculation above or by the following[8]:

$$H^E = \langle E^M \rangle + PV^E - \sum_i x_i \langle E_i \rangle \quad (51)$$

Where $\langle \rangle$ denotes an ensemble average and V^E is calculated using the general method of excess molar properties.

[1] H. et al., Journal of Chemical Physics **120**, 9665 (2004).

[2] D. et al., The Journal of Physical Chemistry B **105**, 715 (2000).

- 362 [3] R. Lustig, *Molecular Simulation* **37**, 457 (2011).
363 [4] W. et al., *Journal of Chemical Theory and Computation* **7**, 2151
364 (2011).
365 [5] C. et al., *Journal of Physical and Chemical Reference Data* **32**,
366 519 (2003).
367 [6] W. et al., *Journal of Physical Chemistry B* **116**, 7088 (2012).
368 [7] C. Fennell, *The Journal of Physical Chemistry B* **116**, 6936
369 (2012).
370 [8] D. et al., *Fluid Phase Equilibria* **289**, 156 (2010).