

A simple method for automated equilibration detection in molecular simulations

John D. Chodera^{1,*}

¹Computational Biology Program, Sloan Kettering Institute,
Memorial Sloan Kettering Cancer Center, New York, NY 10065

(Dated: August 14, 2015)

Molecular simulations intended to compute equilibrium properties are often initiated from configurations that are highly atypical of equilibrium samples, a practice which can generate a distinct initial transient in mechanical observables computed from the simulation trajectory. Traditional practice in simulation data analysis recommends this initial portion be discarded to *equilibration*, but no simple, general, and automated procedure for this process exists. Here, we suggest a conceptually simple automated procedure that does not make strict assumptions about the distribution of the observable of interest, in which the equilibration time is chosen to maximize the number of effectively uncorrelated samples in the production timespan used to compute equilibrium averages. We present a simple Python reference implementation of this procedure, and demonstrate its utility on typical molecular simulation data.

Keywords: molecular dynamics (MD); Metropolis-Hastings; Monte Carlo (MC); Markov chain Monte Carlo (MCMC); equilibration; burn-in; timeseries analysis; statistical inefficiency; integrated autocorrelation time

INTRODUCTION

Molecular simulations use Markov chain Monte Carlo (MCMC) techniques [1] to sample configurations x from an equilibrium distribution $\pi(x)$, either exactly (using Monte Carlo methods such as Metropolis-Hastings) or approximately (using molecular dynamics integrators without Metropolization) [2].

Due to the sensitivity of the equilibrium probability density $\pi(x)$ to small perturbations in configuration x and the difficulty of producing sufficiently good guesses of typical equilibrium configurations $x \sim \pi(x)$, these molecular simulations are often started from highly atypical initial conditions. For example, simulations of biopolymers might be initiated from a fully extended conformation unrepresentative of behavior in solution, or a geometry derived from a fit to diffraction data collected from a cryocooled crystal; solvated systems may be prepared by periodically replicating a small solvent box equilibrated under different conditions, yielding atypical densities and solvent structure; liquid mixtures or lipid bilayers may be constructed by using methods that fulfill spatial constraints (e.g. PackMol [3]) but create locally atypical geometries, requiring long simulation times to relax to typical configurations.

As a result, traditional practice in molecular simulation has recommended some initial portion of the trajectory be discarded to *equilibration* (also called *burn-in*¹ in the MCMC literature [4]). While the process of discarding initial samples is strictly unnecessary for the time-average of quantities of interest to eventually converge to the desired expectations [5], this nevertheless often allows the practitioner to avoid what may be impractically long run times to eliminate the bias in computed properties in finite-length simulations

induced by atypical initial starting conditions. It is worth noting that a similar procedure is not a practice universally recommended by statisticians when sampling from posterior distributions in statistical inference [4]; the differences in complexity of probability densities typically encountered in statistics and molecular simulation may explain the difference in historical practice.

As a motivating example, consider the computation of the average density of liquid argon under a given set of reduced temperature and pressure conditions shown in Figure 1. To initiate the simulation, an initial dense liquid geometry at reduced density $\rho^* \equiv \rho\sigma^3 = 0.960$ was prepared and subjected to local energy minimization. The upper panel of Figure 1 depicts the average relaxation behavior of simulations initiated from the same configuration with different random initial velocities and integrator random number seeds (see *Simulation Details*). The average (black line) and 95% confidence interval (shaded grey) of 500 realizations of this process show a characteristic relaxation behavior away from the initial density toward the equilibrium density. The expectation of the running average of the density over many realizations of this procedure (Figure 1, lower panel) significantly deviates from the true expectation (dashed line), leading to significantly biased estimates of the expectation unless simulations are sufficiently long to eliminate this starting point dependent bias—a surprisingly long 30 ns in this case. Note that this bias is present even in the average of many realizations because the *same* atypical starting condition is used for every realization of this simulation process.

To develop an automatic approach to eliminating this bias, we take motivation from the concept of *reverse cumulative averaging* from Yang et al. [6], in which the trajectory statistics over the production region of the trajectory are examined for different choices of the end of the discarded equilibration region to determine the optimal production region to use for computing expectations and other statistical properties. We begin by first formalizing our objectives mathematically.

* Corresponding author; john.chodera@choderalab.org

¹ The term *burn-in* comes from the field of electronics, in which a short “burn-in” period is used to ensure that a device is free of faulty components—which often fail quickly—and is operating normally [4].

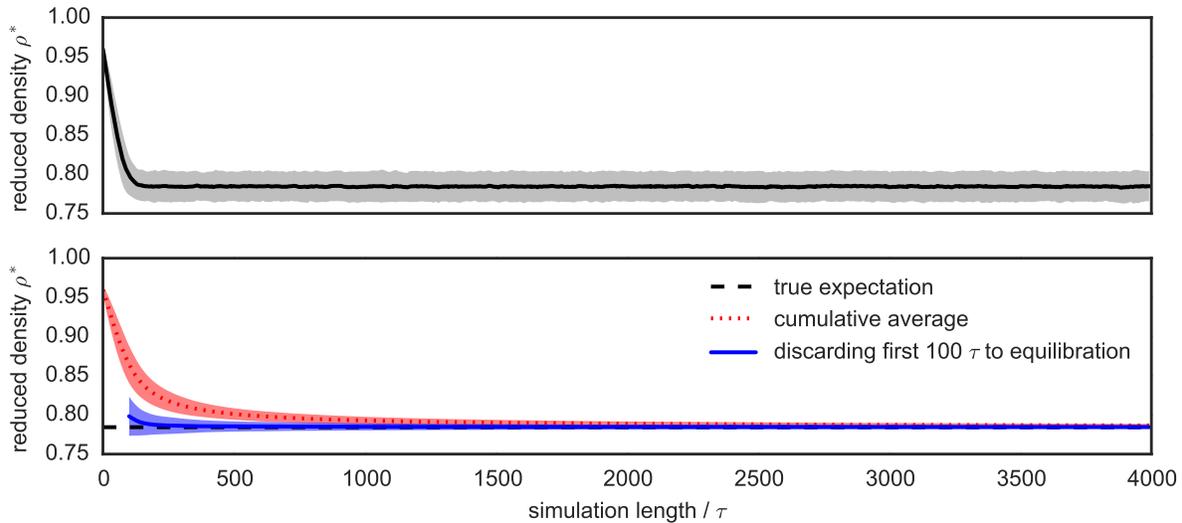


FIG. 1. Illustration of the motivation for discarding data to equilibration. To illustrate the bias in expectations induced by relaxation away from initial conditions, 500 replicates of a simulation of liquid argon were initiated from the same energy-minimized initial configuration constructed with initial reduced density $\rho^* \equiv \rho\sigma^3 = 0.960$ but different random number seeds for stochastic integration. **Top:** The average of the reduced density (black line) over the replicates relaxes to the region of typical equilibrium densities over the first $\sim 90 \tau$ of simulation time, where τ is a natural time unit (see *Simulation Details*). **Bottom:** If the average density is estimated by a cumulative average from the beginning of the simulation (red dotted line), the estimate will be heavily biased by the atypical starting density even beyond 1000τ . Discarding even a small amount of initial data—in this case 500 initial samples—results in a cumulative average estimate that converges to the true average (black dashed line) much more rapidly. Shaded regions denote 95% confidence intervals.

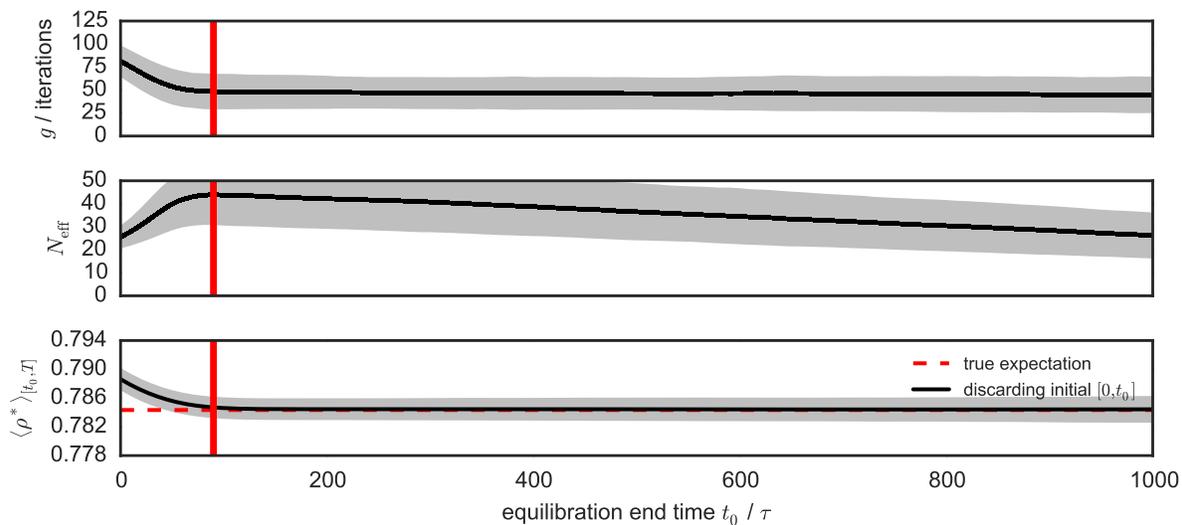


FIG. 2. Statistical inefficiency, number of uncorrelated samples, and bias for different equilibration times. Trajectories of length $T = 2000 \tau$ for the argon system described in Figure 1 were analyzed as a function of equilibration time choice t_0 . Averages over all 500 replicate simulations (all starting from the same initial conditions) are shown as dark lines, with shaded lines showing standard deviation of estimates among replicates. **Top:** The statistical inefficiency g as a function of equilibration time choice t_0 is initially very large, but diminishes rapidly after the system has relaxed to equilibrium. **Middle:** The number of effectively uncorrelated samples $N_{\text{eff}} = (T - t_0 + 1)/g$ shows a maximum at $t_0 \sim 90 \tau$ (red vertical lines), suggesting the system has equilibrated by this time. **Bottom:** The cumulative average density $\langle \rho^* \rangle_{[t_0, T]}$ computed over the span $[t_0, T]$ shows that the bias (deviation from the true estimate, shown as red dashed lines) is minimized for choices of $t_0 \geq 90 \tau$. The standard deviation among replicates (shaded region) grows with t_0 because fewer data are included in the estimate. The choice of optimal t_0 that maximizes N_{eff} (red vertical line) strikes a good balance between bias and variance. The true estimate (red dashed lines) is computed from averaging over the range $[5\,000, 10\,000] \tau$ over all 500 replicates.

STATEMENT OF THE PROBLEM

BIAS-VARIANCE TRADEOFF

78 Consider T successively sampled configurations x_t from
79 a molecular simulation, with $t = 1, \dots, T$, initiated from x_0 .
80 We presume we are interested in computing the expectation

$$\langle A \rangle \equiv \int dx A(x) \pi(x) \quad (1)$$

81 of a mechanical property $A(x)$. For convenience, we will refer
82 to the timeseries $a_t \equiv A(x_t)$, with $t \in [1, T]$. The estimator
83 $\hat{A} \approx \langle A \rangle$ constructed from the entire dataset is given
84 by

$$\hat{A}_{[1,T]} \equiv \frac{1}{T} \sum_{t=1}^T a_t. \quad (2)$$

85 While $\lim_{T \rightarrow \infty} \hat{A}_{[1,T]} = \langle A \rangle$ for an infinitely long simulation²,
86 the bias in $\hat{A}_{[1,T]}$ may be significant in a simulation of
87 finite length T .

88 By discarding samples $t < t_0$ to equilibration, we hope to
89 exclude the initial transient from our sample average, and
90 provide a less biased estimate of $\langle A \rangle$,

$$\hat{A}_{[t_0,T]} \equiv \frac{1}{T - t_0 + 1} \sum_{t=t_0}^T a_t. \quad (3)$$

91 We can quantify the overall error in an estimator $\hat{A}_{[t_0,T]}$
92 in a sample average that starts at x_0 and excludes samples
93 where $t < t_0$ by the expected error $\delta^2 \hat{A}_{[t_0,T]}$,

$$\delta^2 \hat{A}_{[t_0,T]} \equiv E_{x_0} \left[\left(\hat{A}_{[t_0,T]} - \langle A \rangle \right)^2 \right] \quad (4)$$

94 where $E_{x_0}[\cdot]$ denotes the expectation over independent realizations
95 of the specific simulation process initiated from configuration x_0 ,
96 but with different velocities and random number seeds.

98 We can rewrite the expected error $\delta^2 \hat{A}$ by separating it
99 into two components:

$$\delta^2 \hat{A}_{[t_0,T]} = E_{x_0} \left[\left(\hat{A}_{[t_0,T]} - E_{x_0}[\hat{A}_{[t_0,T]}] \right)^2 \right] + \left(E_{x_0}[\hat{A}_{[t_0,T]}] - \langle A \rangle \right)^2 \quad (5)$$

100 The first term denotes the variance in the estimator \hat{A} ,

$$\text{var}_{x_0}(\hat{A}_{[t_0,T]}) \equiv E_{x_0} \left[\hat{A}_{[t_0,T]} - E_{x_0}[\hat{A}_{[t_0,T]}] \right]^2 \quad (6)$$

101 while the second term denotes the contribution from the
102 squared bias,

$$\text{bias}_{x_0}^2(\hat{A}_{[t_0,T]}) \equiv \left(E_{x_0}[\hat{A}_{[t_0,T]}] - \langle A \rangle \right)^2 \quad (7)$$

² We note that this equality only holds for simulation schemes that sample from the true equilibrium density $\pi(x)$, such as Metropolis-Hastings Monte Carlo or Metropolized dynamical integration schemes such as hybrid Monte Carlo (HMC). Molecular dynamics simulations utilizing finite timestep integration without Metropolization will produce averages that may deviate from the true expectation $\langle A \rangle$ [2].

104 With increasing equilibration time t_0 , bias is reduced, but
105 the variance—the contribution to error due to random variation
106 from having a finite number of uncorrelated samples—
107 will increase because less data is included in the estimate.
108 This can be seen in the bottom panel of Figure 2, where
109 the shaded region (95% confidence interval of the mean) increases
110 in width with increasing equilibration time t_0 .

111 To examine the tradeoff between bias and variance explicitly,
112 Figure 3 plots the bias and variance (here, shown as standard error)
113 contributions against each other as a function of t_0 (denoted by color)
114 as computed from statistics over all 500 replicates. At $t_0 = 0$, the bias is large
115 but variance is minimized. With increasing t_0 , bias is eventually
116 eliminated but then variance rapidly grows as fewer uncorrelated
117 samples are included in the estimate. There is a clear optimal choice
118 at $t_0 \sim 90 \tau$ that minimizes variance while also effectively
119 eliminating bias (where τ is a natural time unit—see *Simulation Details*).
120
121

SELECTING THE EQUILIBRATION TIME

123 Is there a simple approach to choosing an optimal equilibration
124 time t_0 that provides a significantly improved estimate
125 $\hat{A}_{[t_0,T]}$, even when we do not have access to multiple
126 realizations? At worst, we hope that such a procedure would
127 at least give some improvement over the naive estimate,
128 such that $\delta^2 \hat{A}_{[t_0,T]} < \delta^2 \hat{A}_{[0,T]}$; at best, we hope that we can
129 achieve a reasonable bias-variance tradeoff close to the optimal
130 point identified in Figure 3 that minimizes bias without greatly
131 increasing variance. We remark that, for cases in which the
132 simulation is not long enough to reach equilibrium, no choice of
133 t_0 will eliminate bias completely; the best we can hope for is
134 to minimize this bias.

135 While automated methods for selecting the equilibration time
136 t_0 have been proposed, these approaches have shortcomings that
137 have greatly limited their use. The reverse cumulative averaging
138 (RCA) method proposed by Yang et al. [6], for example, uses a
139 statistical test for normality to determine the point before which
140 the observable time-series deviates from normality when examining
141 the time-series in reverse. While this concept may be reasonable
142 for experimental data, where measurements often represent the
143 sum of many random variables such that the central limit theorem's
144 guarantee of asymptotic normality ensures the distribution of the
145 observable will be approximately normal, there is no such guarantee
146 that instantaneous measurements of a simulation property of interest
147 will be normally distributed. In fact, many properties will be
148 decidedly *non-normal*. For a biomolecule such as a protein, for
149 example, the radius of gyration, end-to-end distance, and torsion
150 angles sampled during a simulation will all be highly non-normal.
151 Instead, we require a method that makes no assumptions about the
152 nature of the distribution of the property under study.
153
154
155

AUTOCORRELATION ANALYSIS

The set of successively sampled configurations $\{x_t\}$ and their corresponding observables $\{a_t\}$ compose a correlated timeseries of observations. To estimate the statistical error or uncertainty in a stationary timeseries free of bias, we must be able to quantify the *effective number of uncorrelated samples* present in the dataset. This is usually accomplished through computation of the *statistical inefficiency* g , which quantifies the number of correlated timeseries samples needed to produce a single effectively uncorrelated sample of the observable of interest. While these concepts are well-established for the analysis of both Monte Carlo and molecular dynamics simulations [7–10], we review them here for the sake of clarity.

For a given equilibration time choice t_0 , the statistical uncertainty in our estimator $\hat{A}_{[t_0, T]}$ can be written as,

$$\begin{aligned} \delta^2 \hat{A}_{[t_0, T]} &\equiv E_{x_0} \left[\left(\hat{A}_{[t_0, T]} - \langle \hat{A} \rangle \right)^2 \right] \\ &= E_{x_0} \left[\hat{A}_{[t_0, T]}^2 \right] - E_{x_0} \left[\hat{A}_{[t_0, T]} \right]^2 \\ &= \frac{1}{T_{t_0}^2} \sum_{t, t'=t_0}^T \{ E_{x_0} [a_t a_{t'}] - E_{x_0} [a_t] E_{x_0} [a_{t'}] \} \\ &= \frac{1}{T_{t_0}^2} \sum_{t=t_0}^T \{ E_{x_0} [x_t^2] - E_{x_0} [x_t]^2 \} \\ &\quad + \frac{1}{T_{t_0}^2} \sum_{t \neq t'=t_0}^T \{ E_{x_0} [a_t a_{t'}] - E_{x_0} [a_t] E_{x_0} [a_{t'}] \}, \end{aligned} \quad (8)$$

where $T_{t_0} \equiv T - t_0 + 1$, the number of correlated samples in the timeseries $\{a_t\}_{t_0}^T$. In the last step, we have split the double-sum into two separate sums—a term capturing the variance in the observations a_t , and a remaining term capturing the correlation between observations.

If t_0 is sufficiently large for the initial bias to be eliminated, the remaining timeseries $\{a_t\}_{t_0}^T$ will obey the properties of both *stationarity* and *time-reversibility*, allowing us to write,

$$\begin{aligned} \delta^2 \hat{A}_{[t_0, T]}^{\text{equil}} &= \frac{1}{T_{t_0}} [\langle a_t^2 \rangle - \langle a_t \rangle^2] \\ &\quad + \frac{2}{T_{t_0}} \sum_{n=1}^{T-t_0} \left(\frac{T-t_0-n}{T_{t_0}} \right) [\langle a_t a_{t+n} \rangle - \langle a_t \rangle \langle a_{t+n} \rangle] \\ &\equiv \frac{\sigma_{t_0}^2}{T_{t_0}} (1 + 2\tau_{t_0}) = \frac{\sigma_{t_0}^2}{T_{t_0}/g_{t_0}}, \end{aligned} \quad (9)$$

where the variance σ^2 , statistical inefficiency g , and integrated autocorrelation time τ (in units of the sampling interval) are given by

$$\sigma^2 \equiv \langle a_t^2 \rangle - \langle a_t \rangle^2, \quad (10)$$

$$\tau \equiv \sum_{t=1}^{T-1} \left(1 - \frac{t}{T} \right) C_t, \quad (11)$$

$$g \equiv 1 + 2\tau, \quad (12)$$

with the discrete-time normalized fluctuation autocorrelation function C_t defined as

$$C_t \equiv \frac{\langle a_n a_{n+t} \rangle - \langle a_n \rangle^2}{\langle a_n^2 \rangle - \langle a_n \rangle^2}. \quad (13)$$

In practice, it is difficult to estimate C_t for $t \sim T$, due to growth in the statistical error, so common estimators of g make use of several additional properties of C_t to provide useful estimates (see *Practical Computation of Statistical Inefficiencies*).

The t_0 subscript for the variance σ^2 , the integrated autocorrelation time τ , and the statistical inefficiency t_0 mean that these quantities are only estimated over the production portion of the timeseries, $\{a_t\}_{t=t_0}^T$. Since we assumed that the bias was eliminated by judicious choice of the equilibration time t_0 , this estimate of the statistical error will be poor for choices of t_0 that are too small.

THE ESSENTIAL IDEA

Suppose we choose some arbitrary time t_0 and discard all samples $t \in [0, t_0)$ to equilibration, keeping $[t_0, T]$ as the dataset to analyze. How much data remains? We can determine this by computing the statistical inefficiency g_{t_0} for the interval $[t_0, T]$, and computing the effective number of uncorrelated samples $N_{\text{eff}}(t_0) \equiv (T - t_0 + 1)/g_{t_0}$. If we start at $t_0 \equiv T$ and move t_0 to earlier and earlier points in time, we expect that the effective number of uncorrelated samples $N_{\text{eff}}(t_0)$ will continue to grow until we start to include the highly atypical initial data. At that point, the integrated autocorrelation time τ (and hence the statistical inefficiency g) will greatly increase (a phenomenon observed earlier, e.g. Figure 2 of [6]). As a result, the effective number of samples N_{eff} will start to plummet.

Figure 2 demonstrates this behavior for the liquid argon system described above, using averages of the statistical inefficiency g_{t_0} and $N_{\text{eff}}(t_0)$ computed over 500 independent replicate trajectories. At short t_0 , the average statistical inefficiency g (Figure 2, top panel) is large due to the contribution from slow relaxation from atypical initial conditions, while at long t_0 the statistical inefficiency estimate is much shorter and nearly constant of a large span of time origins. As a result, the average effective number of uncorrelated samples N_{eff} (Figure 2, middle panel) has a peak at $t_0 \sim 90 \tau$ (Figure 2, vertical red lines). The effect on bias in the estimated average reduced density $\langle \rho^* \rangle$ (Figure 2, bottom panel) is striking—the bias is essentially eliminated for the choice of equilibration time t_0 that maximizes the number of uncorrelated samples N_{eff} .

This suggests an alluringly simple algorithm for identifying the optimal equilibration time—pick the t_0 which maximizes the number of uncorrelated samples N_{eff} . In mathematical terms,

$$\begin{aligned} t_0^{\text{opt}} &= \underset{t_0}{\operatorname{argmax}} N_{\text{eff}}(t_0) \\ &= \underset{t_0}{\operatorname{argmax}} \frac{T - t_0 + 1}{g_{t_0}} \end{aligned} \quad (14)$$

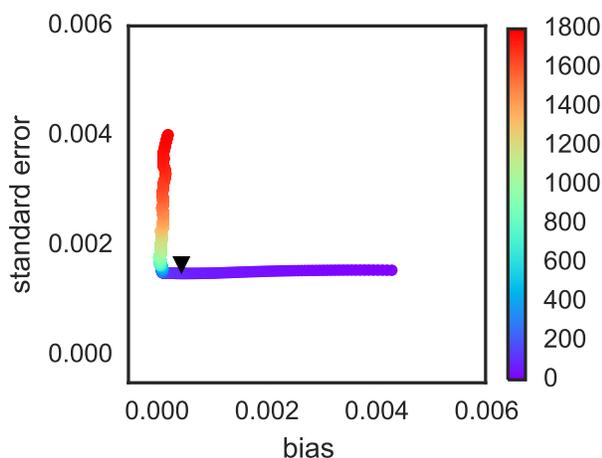


FIG. 3. Bias-variance tradeoff for fixed equilibration time versus automatic equilibration time selection. Trajectories of length $T = 2000\tau$ for the argon system described in Figure 1 were analyzed as a function of equilibration time choice t_0 , with colors denoting the value of t_0 (in units of τ) corresponding to each plotted point. Using 500 replicate simulations, the average bias (average deviation from true expectation) and standard deviation (random variation from replicate to replicate) were computed as a function of a prespecified fixed equilibration time t_0 , with colors running from violet (0τ) to red (1800τ). As is readily discerned, the bias for small t_0 is initially large, but minimized for larger t_0 . By contrast, the standard error (a measure of variance, estimated here by standard deviation among replicates) grows as t_0 grows above a certain critical time (here, $\sim 90\tau$). If the t_0 that maximizes N_{eff} is instead chosen *individually* for each trajectory based on that trajectory's estimates of statistical inefficiency $g_{[t_0, T]}$, the resulting bias-variance tradeoff (black triangle) does an excellent job minimizing bias and variance simultaneously, comparable to what is possible for a choice of equilibration time t_0 based on knowledge of the true bias and variance among many replicate estimates.

Bias-variance tradeoff. How will the simple strategy of selecting the equilibration time t_0 using Eq 14 work for cases where we do not know the statistical inefficiency g as a function of the equilibration time t_0 precisely? When all that is available is a single simulation, our best estimate of g_{t_0} is derived from that simulation alone over the span $[t_0, T]$ —will this affect the quality of our estimate of equilibration time? Empirically, this does not appear to be the case—the black triangle in Figure 3 shows the bias and variance contributions to the error for estimates computed over the 500 replicates where t_0 is individually determined from each simulation using this simple scheme based on selecting t_0 to maximize N_{eff} for each individual realization. Despite not having knowledge about multiple realizations, this strategy effectively achieves a near-optimal balance between minimizing bias without increasing variance.

Overall RMS error. How well does this strategy perform in terms of decreasing the overall error $\delta\hat{A}_{[t_0, T]}$ compared to $\delta\hat{A}_{[0, T]}$? Figure 4 compares the expected standard error (denoted $\delta\hat{A}$) as a function of a fixed initial equilibration time t_0 (black line with shaded region denoting 95% confidence interval) with the strategy of selecting t_0 to maximize N_{eff} for each realization (red line with shaded region denoting 95% confidence interval). While the minimum error for the fixed- t_0 strategy (0.00154 ± 0.00005) is achieved at 90τ —a fact that could only be determined from knowledge of multiple realizations—the simple strategy of selecting t_0 using Eq. 14 achieves a minimum error of 0.00171 ± 0.00006 , only 11% worse (compared to errors of 0.00456 ± 0.00007 , or 296% worse, should no data have been discarded).

DISCUSSION

The scheme described here—in which the equilibration time t_0 is computed using Eq. 14 as the choice that maximizes the number of uncorrelated samples in the production region $[t_0, T]$ —is both conceptually and computationally straightforward. It provides an approach to determining the optimal amount of initial data to discard to equilibration in order to minimize variance while also minimizing initial bias, and does this without employing statistical tests that require generally unsatisfiable assumptions of normality of the observable of interest. As we have seen, this scheme empirically appears to select a practical compromise between bias and variance even when the statistical inefficiency g is estimated directly from the trajectory using Eq. 12.

A word of caution is necessary. One can certainly envision pathological scenarios where this algorithm for selecting an optimal equilibration time will break down. In cases where the simulation is not long enough to reach equilibrium—let alone collect many uncorrelated samples from it—no choice of equilibration time will bestow upon the experimenter the ability to produce an unbiased estimate of the true expectation. Similarly, in cases where insufficient data is available for the statistical inefficiency to be estimated well, this algorithm is expected to perform poorly. However, in these cases, the data itself should be suspect if the trajectory is not at least an order of magnitude longer than the minimum estimated autocorrelation time.

SIMULATION DETAILS

All molecular dynamics simulations described here were performed with OpenMM 6.2 [11] (available at openmm.org) using the Python API. All scripts used to retrieve the software versions used here, run the simulations, analyze data, and generate plots—along with the simulation data itself and scripts for generating figures—are available on GitHub³.

³ All Python scripts necessary to reproduce this work—along with data plotted in the published version—are available at:

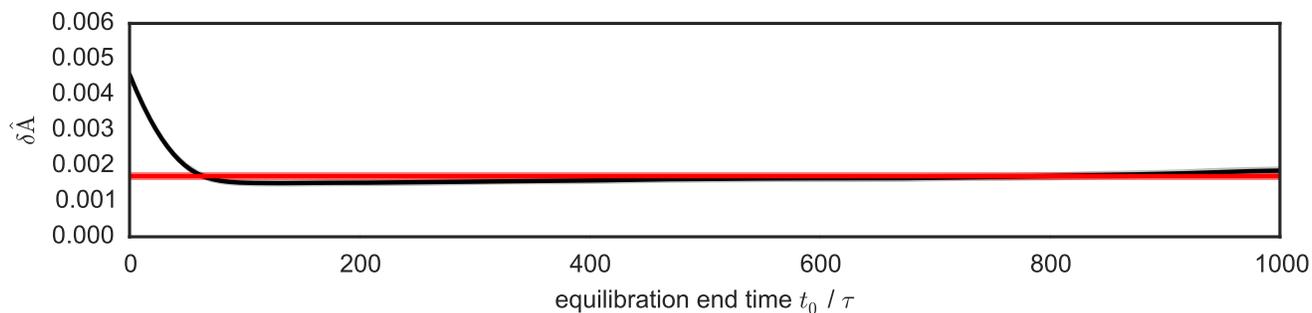


FIG. 4. RMS error for fixed equilibration time versus automatic equilibration time selection. Trajectories of length $T = 2000\tau$ for the argon system described in Figure 1 were analyzed as a function of fixed equilibration time choice t_0 . Using 500 replicate simulations, the root-mean-squared (RMS) error (Eq. 4) was computed (black line) along with 95% confidence interval (gray shading). The RMS error is minimized for fixed equilibration time choices in the range 90–200 τ . If the t_0 that maximizes N_{eff} is instead chosen *individually* for each trajectory based on that trajectory’s estimated statistical inefficiency $g_{[t_0, T]}$ using Eq. 14, the resulting RMS error (red line, 95% confidence interval shown as red shading) is quite close to the minimum RMS error achieved from any particular *fixed* choice of equilibration time t_0 , suggesting that this simple automated approach to selecting t_0 achieves close to optimal performance.

295 To model liquid argon, the LennardJonesFluid model
 296 system in the `openmmtools` package⁴ was used with param-
 297 eters appropriate for liquid argon ($\sigma = 3.4 \text{ \AA}$, $\epsilon = 0.238$
 298 kcal/mol). All results are reported in reduced (dimension-
 299 less) units. A cubic switching function was employed, with
 300 the potential gently switched to zero over $r \in [\sigma, 3\sigma]$, and
 301 a long-range isotropic dispersion correction accounting for
 302 this switching behavior used to include neglected contribu-
 303 tions. Simulations were performed using a periodic box of
 304 $N = 500$ atoms at reduced temperature $T^* \equiv k_B T / \epsilon =$
 305 0.850 and reduced pressure $p^* \equiv p\sigma^3 / \epsilon = 1.266$ using a
 306 Langevin integrator [12] with timestep $\Delta t = 0.01\tau$ and col-
 307 lision rate $\nu = \tau^{-1}$, with characteristic oscillation timescale
 308 $\tau = \sqrt{mr_0^2 / 72\epsilon}$ and $r_0 = 2^{1/6}\sigma$ [13]. All times are reported
 309 in multiples of the characteristic timescale τ . A molecu-
 310 lar scaling Metropolis Monte Carlo barostat with Gaussian
 311 simulation volume change proposal moves attempted every
 312 τ (100 timesteps), using an adaptive algorithm that ad-
 313 justs the proposal width during the initial part of the simu-
 314 lation [11]. Densities were recorded every τ (100 timesteps).
 315 The true expectation $\langle \rho^* \rangle$ was estimated from the sample
 316 average over all 500 realizations over [5000, 10000] τ .

317 The automated equilibration detection scheme is also
 318 available in the `timeseries` module of the `pybar` pack-
 319 age as `detectEquilibration()`, and can be accessed us-
 320 ing the following code:

```

from pybar.timeseries import detectEquilibration
# determine equilibrated region
[t0, g, Neff_max] = detectEquilibration(A_t)
# discard initial samples to equilibration
A_t = A_t[t0:]
    
```

321 PRACTICAL COMPUTATION OF STATISTICAL INEFFICIENCIES

322 The robust computation of the statistical inefficiency g
 323 (defined by Eq. 12) for a finite timeseries $a_t, t = 0, \dots, T$
 324 deserves some comment. There are, in fact, a variety of
 325 schemes for estimating g described in the literature, and
 326 their behaviors for finite datasets may differ, leading to dif-
 327 ferent estimates of the equilibration time t_0 using the algo-
 328 rithm of Eq. 14.

329 The main issue is that a straightforward approach to es-
 330 timating the statistical inefficiency using Eqs. 11–13 in which
 331 the expectations are simply replaced with sample estimates
 332 causes the statistical error in the estimated correlation func-
 333 tion C_t to grow with t in a manner that allows this error to
 334 quickly overwhelm the sum of Eq. 11. As a result, a number of
 335 alternative schemes—generally based on controlling the er-
 336 ror in the estimated C_t or truncating the sum of Eq. 11 when
 337 the error grows too large—have been proposed.

338 For stationary, irreducible, reversible Markov chains,
 339 Geyer observed that a function $\Gamma_k \equiv \gamma_{2k} + \gamma_{2k+1}$ of the
 340 unnormalized fluctuation autocorrelation function $\gamma_t \equiv$
 341 $\langle a_i a_{i+t} \rangle - \langle a_i \rangle^2$ has a number of pleasant properties (The-
 342 orem 3.1 of [14]): It is strictly positive, strictly decreasing,
 343 and strictly convex. Some or all of these properties can be
 344 exploited to define a family of estimators called *initial se-*
 345 *quence methods* (see Section 3.3 of [14] and Section 1.10.2
 346 of [4]), of which the *initial convex sequence* (ICS) estimator is
 347 generally agreed to be optimal, if somewhat more complex
 348 to implement.⁵

349 All computations in this manuscript used the fast mul-
 350 tiscule method described in Section 5.2 of [10], which we
 351 found performed equivalently well to the Geyer estimators
 352 (data not shown). This method is related to a multiscale

⁴ <http://github.com/choderalab/automatic-equilibration-detection>
⁵ available at <http://github.com/choderalab/openmmtools> Implementations of these methods are provided with the code dis-
 tributed with this manuscript.

353 variant of the *initial positive sequence* (IPS) method of Geyer
354 [15], where contributions are accumulated at increasingly
355 longer lag times and the sum of Eq. 11 is truncated when the
356 terms become negative. We have found this method to be
357 both fast and to provide useful estimates of the statistical
358 inefficiency, but it may not perform well for all problems.

ACKNOWLEDGMENTS

360 We are grateful to William C. Swope (IBM Almaden Re-
361 search Center) for his illuminating introduction to the use
362 of autocorrelation analysis for the characterization of sta-

363 tistical error, as well as Michael R. Shirts (University of Vir-
364 ginia), David L. Mobley (University of California, Irvine),
365 Michael K. Gilson (University of California, San Diego), Kyle
366 A. Beauchamp (MSKCC), and Robert C. McGibbon (Stan-
367 ford University) for valuable discussions on this topic, and
368 Joshua L. Adelman (University of Pittsburgh) for helpful
369 feedback and encouragement. We are grateful to Michael
370 K. Gilson (University of California, San Diego) and Wei
371 Yang (Florida State University) for critical feedback on the
372 manuscript itself. JDC acknowledges a Louis V. Gerstner
373 Young Investigator Award, NIH core grant P30-CA008748,
374 and the Sloan Kettering Institute for funding during the
375 course of this work.

-
- 376 [1] J. S. Liu, *Monte Carlo strategies in scientific computing*, 2nd ed.
377 ed. (Springer-Verlag, New York, 2002).
378 [2] D. Sivak, J. Chodera, and G. Crooks, *Physical Review X* **3**,
379 011007 (2013), bibtex: Sivak:2013:Phys.Rev.X.
380 [3] L. Martínez, R. Andrade, E. G. Birgin, and J. M. Martínez, *J.*
381 *Chem. Theor. Comput.* **30**, 2157 (2009).
382 [4] S. Brooks, A. Gelman, G. L. Jones, and X.-L. Meng, in *Hand-*
383 *book of Markov chain Monte Carlo*, Chapman & Hall/CRC *Hand-*
384 *books of Modern Statistical Methods* (CRC Press, ADDRESS,
385 2011), Chap. Introduction to Markov chain Monte Carlo.
386 [5] C. Geyer, Burn-in is unnecessary., [http://users.stat.umn.](http://users.stat.umn.edu/~geyer/mcmc/burn.html)
387 [edu/~geyer/mcmc/burn.html](http://users.stat.umn.edu/~geyer/mcmc/burn.html).
388 [6] W. Yang, R. Bittetti-Putzer, and M. Karplus, *J. Chem. Phys.* **120**,
389 2618 (2004).
390 [7] H. Müller-Krumbhaar and K. Binder, *J. Stat. Phys.* **8**, 1 (1973).
391 [8] W. C. Swope, H. C. Andersen, P. H. Berens, and K. R. Wilson, *J.*
392 *Chem. Phys.* **76**, 637 (1982).
393 [9] W. Janke, in *Quantum Simulations of Complex Many-Body Sys-*
394 *tems: From Theory to Algorithms*, edited by J. Grotendorst, D.
395 Marx, and A. Murmatsu (John von Neumann Institute for Com-
396 puting, ADDRESS, 2002), Vol. 10, pp. 423–445.
397 [10] J. D. Chodera, W. C. Swope, J. W. Pitera, C. Seok, and K. A. Dill,
398 *J. Chem. Theor. Comput.* **3**, 26 (2007).
399 [11] P. Eastman, M. Friedrichs, J. D. Chodera, R. Radmer, C. Bruns,
400 J. Ku, K. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, T. Tye,
401 M. Houston, T. Stitch, and C. Klein, *J. Chem. Theor. Comput.* **9**,
402 461 (2012).
403 [12] D. A. Sivak, J. D. Chodera, and G. E. Crooks, *J. Phys. Chem. B*
404 **118**, 6466 (2014).
405 [13] B. Veytsman and M. Kotelyanskii, Lennard-Jones poten-
406 tial revisited., [http://borisv.lk.net/matsc597c-1997/](http://borisv.lk.net/matsc597c-1997/simulations/Lecture5/node3.html)
407 [simulations/Lecture5/node3.html](http://borisv.lk.net/matsc597c-1997/simulations/Lecture5/node3.html).
408 [14] C. J. Geyer, *Stat. Sci.* **76**, 473 (1992).
409 [15] C. J. Geyer and E. A. Thompson, *J. Royal Stat. Soc. B* **54**, 657
410 (1992).