This article was downloaded by:

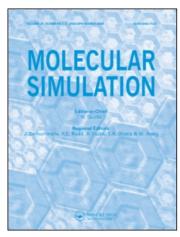
On: 14 January 2011

Access details: Access Details: Free Access

Publisher Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-

41 Mortimer Street, London W1T 3JH, UK



Molecular Simulation

Publication details, including instructions for authors and subscription information: http://www.informaworld.com/smpp/title~content=t713644482

Basic statistics and variational concepts behind the reverse Monte Carlo technique

F. L. B. da Silva^{ab}; W. Olivares-Rivas^c; P. J. Colmenares^c

^a Departamento de Física e Química, Faculdade de Ciências Farmacêuticas de Ribeirão Preto, Universidade de São Paulo, Ribeirão Preto, SP, Brazil ^b Theoretical Chemistry, Lund University, Lund, Sweden ^c Grupo de Química Teórica: Quimicofísica de Fluidos y Fenómenos Interfaciales (QUIFFIS), Departamento de Química, Universidad de Los Andes, Mérida, Venezuela

To cite this Article da Silva, F. L. B., Olivares-Rivas, W. and Colmenares, P. J.(2007) 'Basic statistics and variational concepts behind the reverse Monte Carlo technique', Molecular Simulation, 33: 8, 639-647

To link to this Article: DOI: 10.1080/08927020701361884 URL: http://dx.doi.org/10.1080/08927020701361884

PLEASE SCROLL DOWN FOR ARTICLE

Full terms and conditions of use: http://www.informaworld.com/terms-and-conditions-of-access.pdf

This article may be used for research, teaching and private study purposes. Any substantial or systematic reproduction, re-distribution, re-selling, loan or sub-licensing, systematic supply or distribution in any form to anyone is expressly forbidden.

The publisher does not give any warranty express or implied or make any representation that the contents will be complete or accurate or up to date. The accuracy of any instructions, formulae and drug doses should be independently verified with primary sources. The publisher shall not be liable for any loss, actions, claims, proceedings, demand or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of this material.



Basic statistics and variational concepts behind the reverse Monte Carlo technique

F. L. B. DA SILVA†‡, W. OLIVARES-RIVAS¶* and P. J. COLMENARES¶

†Departamento de Física e Química, Faculdade de Ciências Farmacêuticas de Ribeirão Preto, Universidade de São Paulo, Avenue do café, s/no., 14040-903 Ribeirão Preto, SP, Brazil

‡Theoretical Chemistry, Lund University, POB 124, S-221 00 Lund, Sweden

¶Grupo de Química Teórica: Quimicofísica de Fluidos y Fenómenos Interfaciales (QUIFFIS), Departamento de Química, Universidad de Los Andes, Mérida 5101, Venezuela

(Received November 2006; in final form March 2007)

We revise the statistical foundations of the reverse Monte Carlo (RMC) technique by constructing the associated functional of a variational principle which incorporates, without any *ad hoc* assumptions, the inherent errors accompanying the simulation and the experimental data. We propose a Bayesian criteria for acceptance/rejection of configurations, in terms of the variations of the functional. The loss function and variational functional minimization approaches are compared.

Keywords: Statistical mechanics; Reverse Monte Carlo; Functional analysis; Loss function

1. Introduction

Even though reverse Monte Carlo (RMC) method is a successful molecular modeling technique [1-4], used now-a-days as routinely as traditional techniques (e.g. Metropolis Monte Carlo (MMC) [5,6] and molecular dynamics (MD) [7,8]), there are aspects not yet rigorously established [9,10]. Basically, it works like the MMC method, except that in this case, an input structural property, the radial distribution function (rdf) or the structure factor (sf), is used instead of the MMC energy criterion. This does not mean that RMC is the reversed form of MMC. RMC is closer to a sophisticated fitting of equations that provides configurations containing on average the same structural information as described by the input function. That is, one distributes points in the space that, on the average, reproduce the input rdf or sf. This idea goes back to the work of Kaplow and collaborators [11]. Later, it has been modified and improved by different groups [1,3,12-14]. There has also appeared closely related techniques such as Soper's empirical potential Monte Carlo (EPMC) [15,16]. We recognize that the acronym RMC, originally adopted by McGreevy school [1,9,12], is used in the literature to denote, in general, the inverse problem of a MC

simulation; so, whenever we refer in this work to the original RMC work of McGreevy and Putztai [1], we shall explicitly denote it as the MP-RMC method.

The advantage of such techniques is that no interatomic (pair) potential is needed, only structural information. Typically, this information may be the experimental rdf or the sf. Therefore, the technique may be employed to extract more detailed information from these functions, e.g. orientational or high-order correlations. In fact, it has already been claimed that this kind of procedure will have a major effect on our understanding of solution's structure, and will open up the field of "liquid state crystallography" [17]. This will, of course, have an impact on the understanding of general aqueous solutions properties, even those containing bio-molecules that are far more complicated. Furthermore, although the method was designed primarily as a tool to analyze and interpret experimental data, there is also a possibility that the method can help in the development of effective pair potentials [18,19].

The major problem with the RMC is the lack of a solid theoretical background for the technique. There are only arguments based on Henderson's uniqueness theorem for correlation functions [20], which assures that there is a unique relation between the pair potential and the rdf, g(r)

^{*}Corresponding author. Email: wilmer@ula.ve

(see also Ref. [21]). It also guarantees that g(r) determines implicitly all higher-order correlations functions, if the interaction potential is assumed to be pairwise additive. Thus, for a given rdf, at a specific density and temperature, RMC may be able to extract this "hidden information".

In a previous work by one of us [3], referred to as paper I, a statistical based RMC (SRMC) method was suggested and applied it with great success to liquid water [4]. This new algorithm does, in contrast to several other earlier algorithms [1,14], meet the fundamental requirements. That is, it reproduces the input structural information and relevant thermodynamic properties with good accuracy both for monatomic and polyatomic liquids [3,4,22-24]. The SRMC method differs from the original MP-RMC method of McGreevy and Putztai [1], essentially on the way the statistical mechanical rdf function g(r) is evaluated, and on the criteria for accepting configurations. The adequacy of the SRMC approach was shown numerically in paper I [3]. A small number of particles, typically less than 100, is enough to give proper convergence and there is no need for additional constraints. No problems within the hard-core range were found and thermodynamic properties, as the configurational average energy or the excess chemical potential, calculated using an ad-hoc model, were well reproduced. In paper I, it was also shown that SRMC is able to extract three-body correlations from the input two-body rdf of an accurately simulated MMC LJ fluid [3]. Despite some initial criticism [22,23], in a second work [4], referred to as paper II, the SRMC method was applied successfully to the case of a polyatomic fluid data, like water SPC models and accurate experimental liquid water data. For instance, SRMC was able to show how the hydrogen bonds were distributed and oriented. Nevertheless, no statistical mechanical argument was given to support it.

The purpose of this work is to give a variational and statistical view for the foundation of the RMC technique in general and, of the SRMC in particular. A second issue here, is an analysis on how the SRMC inverse simulation samples the phase space, in comparison to the corresponding MC simulation that uses the Metropolis acceptance criteria.

2. A variational approach to reverse Monte Carlo

The RMC simulation requires only the experimental rdf, $g_{\rm exp}(r)=g^e(r)$, and the system density as input. The fundamental requirement is that the RMC generated configurations are, on the average, consistent with such input function, i.e. the RMC rdf corresponding to these configurations, $g_{\rm RMC}(r)=g(r)$, should satisfy, for every point in space

$$g(r) = g^{e}(r). \tag{1}$$

It is well established in statistical mechanics that the rdf is a functional of the intermolecular potential [20,25].

In a classical system, the available energy states correspond to the possible particle configurations in phase space. Therefore, by sampling configurations, one samples phase space for states with an associated trial rdf, consistent with the previous relationship equation (1). To establish this identity we assume in this section that both, the experimental rdf function $g_{\rm exp}(r)$ and the statistical rdf $g_{\rm RMC}(r)$, are very accurately measured or evaluated. The more realistic case, where there are implicit errors in the measurement of $g_{\rm exp}(r)$ and in the evaluation of $g_{\rm RMC}(r)$, will be carefully analyzed in next section. So here, we sample for configurations that let the following composed function, $\mathcal{B}[g]$, vanish:

$$\mathcal{B}[g(r)] = f(r)[g(r) - g^e(r)] \to 0, \tag{2}$$

where f(r) is any arbitrary well behaved function. That is, one has to construct a functional $\mathcal{F}[g(r)]$ consistent with equation (2), in the sense that the functional derivative gives the vanishing function,

$$\frac{\delta \mathcal{F}[g(r)]}{\delta g(r)} = \mathcal{B}[g(r)]. \tag{3}$$

Following the prescriptions of Olivares-Rivas [26] based on the variational method of Arthurs [27], it is straightforward to show that the required functional is given as

$$\mathcal{F}[g(r)] = \mathcal{F}^* + \int d\vec{\mathbf{r}} \,\mathcal{H}[g(r)], \tag{4}$$

where

$$\mathcal{H}[g(r)] = \int_0^1 \mathrm{d}t f(r)g(r)B[tg(r)],\tag{5}$$

and \mathcal{F}^* is a quantity independent of g(r). Equation (5) can be readily integrated over t to give

$$\mathcal{H}[g(r)] = \frac{1}{2}g^{2}(r) - g(r)g^{e}(r). \tag{6}$$

In order to complete the square, we chose the arbitrary constant in equation (4) as

$$\mathcal{F}^* = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) (g^e(r))^2, \tag{7}$$

to get

$$\mathcal{F}[g(r)] = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) [g(r) - g^{e}(r)]^{2}.$$
 (8)

By a simple application of the definition of the functional derivative,

$$\frac{\delta \mathcal{F}[g(r)]}{\delta g(r)} = \int d\vec{\mathbf{r}} \frac{d\mathcal{H}[g(r)]}{dg(r)} \delta g(r), \tag{9}$$

we recover equation (2) from equation (8). So, the RMC variational prescription requires the mimization of the funtional $\mathcal{F}[g(r)]$ of equation (8).

Now, since the function f(r) is arbitrary, we have several choices. The most trivial one, f(r) = 1, gives the least

square integral over the volume. Another option is $f(r) = (2/4\pi r^2)$, which gives

$$\mathcal{F}[g(r)] = \int_{0}^{\infty} [g(r) - g^{e}(r)]^{2} dr.$$
 (10)

Discretizing over a grid of n_p points and a step dr, corresponding to the bin width [3], we can simply take the functional as

$$\mathcal{F}[g(r)] = \sum_{i}^{n_p} [g(r_i) - g^e(r_i)]^2.$$
 (11)

In practice, the minimization of the functional $\mathcal{F}[g(r)]$ with respect to the variations $\delta g(r)$ of g(r) is obtained by randomly generating several different configurations, as in a regular MMC simulation, but only accepting those configurations that lead to a lower value of this functional [3]. This form of the functional $\mathcal{F}[g(r)]$ is similar to the standard definition of the χ^2 parameter, used in the MP–RMC literature [1,9].

$$\chi^2 = \frac{1}{\sigma^2} \sum_{i=1}^{n_p} \left[g_{\text{RMC}}(r_i) - g_{\text{EXP}}(r_i) \right]^2, \tag{12}$$

where n_p is the number of points in the experimental rdf data and, $g_{\rm RMC}(r_i)$ is the calculated rdf in the RMC run. The parameter σ^2 does not appear in the variational functional. It is arbitrarily introduced in the standard MP–RMC prescription and is conventionally assumed to be the experimental error in $g_{\rm EXP}(r_i)$ without further physical arguments. Typically, a value around 0.1 is adopted. The fact of the matter is, that there is no rigorous foundation for this parameter, which is commonly related to the acceptance of the configurations, and is taken by most authors as a weighting parameter. Some authors have already used σ^2 in this context, applying some sort of "annealing" in order to improve the simulation convergence [28]. Other authors have even suggested the minimization of the functional [10]

$$\chi^{2} = \frac{\sum_{i=1}^{n_{p}} [g_{\text{RMC}}(r_{i}) - g_{\text{EXP}}(r_{i})]^{2}}{\sum_{i=1}^{n_{p}} [g_{\text{EXP}}(r_{i})]^{2}},$$
 (13)

where the normalizing factor in the denominator corresponds to an integral in r space. In our simple variational formulation, however, there is no justification or need to invoke the standard deviation error parameter and, more importantly, there is no justification for the Metropolis type of acceptance criterium involving the exponential of $\Delta\chi^2$, as largely used in earlier algorithms [1,14].

It has been argued that the introduction of a σ^2 parameter is, perhaps, justified on a curve fitting of experimental data with large errors, where the target function $g^e(r)$ is not well known [9,14].

In order to address to this point, in the next section, we shall combine the variational construction of a suitable functional, with some basic Bayesian statistical concepts, to introduce, in a natural manner, the experimental error in the inverse problem, posed by the reverse MC methods.

3. Statistical RMC and experimental error

The goal of a fitting procedure, like RMC, is to minimize the loss of accuracy of the estimation of a given experimental data by a proposed model. Statistically, this is commonly accomplished by numerical minimization of the so called loss function [29]. The most popular loss function is the likelihood function L, which measures how likely is the observed distribution, given some model. A common statistical test is based on assuming that the conditional probability is maximized for every point in space, or, for every bin, in a discrete space. So, the goodness of fit is conjectured to correspond to maximizing a likelihood estimator directly given by the conditional probability of the event, namely,

$$L = P[g_1, g_2, g_3, \dots | g_1^e, g_2^e, g_3^e, \dots] = \prod_i P[g_i | g_i^e],$$
(14)

or, more conveniently, to minimize the loss function $F = -2 \ln L$

$$F = -2 \ln L = -2 \sum_{i} \ln P[g_i | g_i^e], \tag{15}$$

where $P[g_i|g_i^e]$ is the Bayesian conditional probability of getting g_i given g_i^e , in the *i*th bin. When $P[g_i|g_i^e]$ is a Poisson distribution this is called the Pearson χ^2 fit

$$F \approx \chi^2 = \sum_i \frac{(g_i - g_i^e)^2}{g_i^e}.$$
 (16)

Since, for a Poisson distribution, the squared standard deviation is equal to the mean value, it is common to replace in the previous expression $\sigma_i^2 = \bar{g}_i$ instead of g_i^e . On the other hand, if the probability distribution is gaussian, this is equivalent to the so called least square fit. The *a priori* assumption that the conditional probability $P[g_i|g_i^e]$ is Gaussian distributed around the value g_i^e with a standard deviation σ_i

$$P[g_i|g_i^e] = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{(g_i - g_i^e)^2}{2\sigma_i^2}\right],$$
 (17)

is then equivalent to, arbitrarily, choose the associated loss function as

$$F = \sum_{i} \frac{(g_i - g_i^e)^2}{\sigma_i^2}.$$
 (18)

Since $P[g_i|g_i^e]$ is a conditional probability, the meaning of σ_i^2 in equations (17) and (18) is not as transparent as an experimental error. The difficulty of an inverse MC simulation is that the set $g(r) = \{g_i\}$ itself, is formally a well defined, statistical mechanical distribution function. How this statistical quantity is evaluated or approximated clearly depends on the method used. In the original MP-RMC method, g(r) is approximated by the average obtained over a single configuration. In the SRMC method of paper I, the average is obtained over an ensemble of

collected configurations. We shall revise here the statistical problem behind any RMC technique, to further clarify this point and the meaning of the terms in equation (18).

The problem in a MMC simulation is to sample phase space by generating molecular configurations based on mechanical information and using an importance acceptance criteria, that ensures the microscopic reversibility of the Markov chain, when moving from one configuration to another. Each accepted configuration gives a rdf, g(r), and the ensemble average $\langle g(r) \rangle$ is the output. In an inverse problem, like RMC, phase space is also sampled, but no mechanical information is available.

We have argued before [3,4,23] that the output of SRMC should be a collection of randomly generated configurations that gives an ensemble average $\langle g(r) \rangle$, which mimics the experimentally measured rdf or sf. The principle behind the SRMC of da Silva *et al.* [3,4] is the minimization of the differences between the input rdf and the calculated one, by randomly generating particle configurations, and accumulating them, to generate a new trial rdf. The fundamental basis is the equivalence between particles and fields. For pairwise potentials such an equivalence follows from the fact that the rdf is a unique functional of the intermolecular potential [30].

One expects that an accurately measured rdf or sf, does contain implicitly a great deal of information about both, two-body and higher-order correlations. However, the information content decreases with the loss of accuracy and range of the data. The inherent error accompanying the experimental measurement can be represented by a gaussian probability density function

$$P[g_i^e] = \frac{1}{\sqrt{2\pi}(\sigma_e)_i} \exp\left[-\frac{(g_i^e - \overline{g_e}_i)^2}{2(\sigma_e^2)_i}\right], \quad (19)$$

where $\overline{g^e}_i$ and $(\sigma_e^2)_i$ are the experimental average rdf and the standard deviation of a set of experimentally measured rdf or sf in the *i*th bin, respectively.

In the same manner, the collected trial rdf can also be assumed to be distributed about the answer trial rdf, $g_a^i = \langle g^i \rangle$, averaged over all the accepted trials, with a simulation error or standard deviation σ_i

$$P[g_i] = \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{(g_i - \langle g_i \rangle)^2}{2\sigma_i^2}\right]. \tag{20}$$

The experimental and simulation standard deviations, σ_{ei} and σ_i , are in principle different from each other, for every r space bin. From now on, we will denote the experimental mean, $\overline{g^e(r)}$, and the SRMC simulation mean $\langle g(r) \rangle$ differently, to stress their different nature. In fact, the first is the result of $n_{\rm exp}$ independent and identical experimental measurements and the latter corresponds to an average over the accumulated $n_{\rm obs}$ observations in all the RMC simulation cycles.

If the events, described by these probability densities functions, are uncorrelated (zero covariance), the joint

probability $P[g^e(r), g(r)]$ is given by their product

$$P[g(r)|g^{e}(r)]P[g^{e}(r)] = P[g^{e}(r),g(r)] = P[g(r)]P[g^{e}(r)],$$
(21)

for every point in r space.

So, when the experimental data set is noisy, one is not interested in reproducing such experimental or training set exactly, but to make null the average difference between the model predicted rdf and the experimentally measured rdf, as pointed out by Webb [31]. Thus, following the variational prescription described above, the vanishing composed function should be, instead

$$\langle g(r) - g^e(r) \rangle_{(g,g^e)} \to 0,$$
 (22)

where the subscript on the brackets means that the average has to be taken on the joint probability distribution of the RMC rdf and the experimental value at that point. This is the essential conjecture of this work which, as we will see later, justifies on theoretically and statistical grounds the SRMC method of da Silva *et al.* [3,4].

The variational functional $\mathcal{F}[g(r)]$ associated with this problem, then have to include such averaging, i.e. as before, we require that

$$\frac{\delta \mathcal{F}[g(r)]}{\delta g(r)} = f(r) \langle g(r) - g^e(r) \rangle_{(g,g^e)} \to 0, \quad (23)$$

where f(r) is any arbitrary well behaved function. Following the prescription given in the previous section

$$\mathcal{F}[g(r)] = \mathcal{F}^* + \int d\vec{\mathbf{r}} \langle \mathcal{H}[g(r)] \rangle_{(g,g^e)}, \tag{24}$$

where \mathcal{F}^* is a quantity independent of g(r) and $\langle \mathcal{H}[g(r)] \rangle_{(g,g^e)}$ is given by

$$\langle \mathcal{H}[g(r)] \rangle_{(g,g^e)} = f(r) \left\langle \frac{1}{2} g^2(r) - g(r) g^e(r) \right\rangle_{(g,g^e)}.$$
 (25)

Therefore, by choosing

$$\mathcal{F}^* = \frac{1}{2} \left\langle \int d\vec{\mathbf{r}} f(r) [g^e(r)]^2 \right\rangle_{(q,q^e)},\tag{26}$$

we find that

$$\mathcal{F}[g(r)] = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) \langle [g(r) - g^e(r)]^2 \rangle_{(g,g^e)}. \tag{27}$$

If we assume that g(r) and $g^{e}(r)$ are statistical independent and use the definition of the standard deviation of the experimental rdf, σ_{e}

$$\sigma_e^2 = \langle [g^e(r)]^2 \rangle_{g^e} - \langle g^e(r) \rangle_{g^e}^2, \tag{28}$$

the integrand of equation (27) can be written out in the form

$$\mathcal{F}[g(r)] = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) \left[\left\langle \left[g(r) - \overline{g^e(r)} \right]^2 \right\rangle_g + \sigma_e^2 \right], \quad (29)$$

where $\overline{g^e(r)} = \langle g^e(r) \rangle_{g_e}$ and the average in equation (29) is taken over the distribution of g(r). This form of the

functional reminds the Taguchi loss functional [32], extensively used in quality control and robust design in manufacturing processes. It is very suitable for computation on a simulation, since it only requires the mean of the experimental measurement at each point of space. In the Appendix A, we show that the functional derivative of this functional recovers the desired relationship, equation (23).

We can use the definition of the SRMC simulation standard deviation σ_{σ}^2

$$\sigma_g^2 = \langle g(r)^2 \rangle_g - \langle g(r) \rangle_g^2, \tag{30}$$

to further rewrite equation (29) as:

$$\mathcal{F}[g(r)] = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) \left\{ \left[\langle g(r) \rangle - \overline{g^e(r)} \right]^2 + \sigma_e^2(r) + \sigma_g^2(r) \right\}. \tag{31}$$

Again, using $f(r) = (2/4\pi r^2)$ and representing the r space in n_p bins, we finally get:

$$\mathcal{F}[g(r)] = \sum_{i}^{n_p} \left\{ \left[\langle g_i \rangle - \overline{g_i^e} \right]^2 + (\sigma_e^2)_i + (\sigma_g^2)_i \right\}. \tag{32}$$

The standard deviation $(\sigma_e^2)_i$ is a fixed number, determined experimentally, which indicates the error of the g e measurement and the quality of the information it contains. Therefore, it is $(\sigma_{\sigma}^2)_i$ the parameter that monitors the dispersion of the simulated rdf answer g(r). As mentioned above, MP-RMC g_i is usually evaluated over a single trial configuration [1], so MP-RMC σ_g depends strongly on the number N of central particles used in the simulation. In SRMC (paper I) [3], g_i is evaluated also over the N central particles in each configuration, but then the n_c collected configurations in the simulation are taken into account; for this reason, the number of particles does not need to be as large as in MP-RMC. For the same number of particles N, a better accuracy is obtained in the SRMC method as the number of configurations increases. We must stress that even in the MP-RMC method, with large N, the error σ_g is expected to be much smaller than the experimental error σ_e . One can evaluate it during the simulation or chose it as a fixed parameter to measure the desired tolerance of the RMC simulation.

In the second case, the σ terms can be omitted and the functional in equation (32) reduces to the expression in the previous section, equation (11), with the understanding that g_i^e and g_i are average values, as originally proposed by da Silva *et al.* [3,4].

4. A Bayesian metropolis importance sampling

The MMC importance sampling criteria is often justified invoking microscopic reversibility in moving in phase space, from one configuration with an energy state U_n to another with energy state U_m . So the transition probabilities π_{mn} and π_{nm} are related to the canonical

probabilities of the given states n and m, P_n and P_m , as

$$\frac{\pi_{mn}}{\pi_{nm}} = \frac{P_n}{P_m} = \exp\{-(U_m - U_n)/k_B T\},\tag{33}$$

where $k_{\rm B}$ is Boltzmann constant and, T the temperature. Thus, the new configurations are accepted with probability

$$P_{\rm acc} = \min[1, \exp\{-(U_{\rm new} - U_{\rm old})/k_{\rm B}T\}].$$
 (34)

In the case of RMC, we can use an analogous argument, using Bayes theorem. We assume that, in the production cycles the trial rdf's fluctuate about the target function $\overline{g^e}$ with a gaussian normal probability density function. Then, a given configuration n giving rise to a rdf $\langle g_n(r_i) \rangle$, for the ith bin, has a sampling distribution of the mean given by the probability distribution

$$P[\langle g_n(r_i)\rangle] \frac{1}{\sqrt{2\pi}\sigma_{\langle g\rangle}} \exp\left[-\frac{\left(\langle g_n(r_i)\rangle - \overline{g^e}\right)^2}{2\sigma_{\langle g\rangle}^2}\right], \quad (35)$$

where $\sigma_{\langle g \rangle}^2$ is the variance of the sample distribution of the mean, a measure of the standard error of the mean $\langle g \rangle$. For simplicity, we assumed the same standard error $\sigma_{\langle g \rangle}$ for every bin. The probability of observing the rdf $\langle g_n \rangle$ for all the bins is the joint probability $P[g_n] = \prod_i^{n_p} P[\langle g_n(r_i) \rangle]$. The log loss function associated to this probability is, within a constant, given by the functional $\mathcal{F}[g_n]$ as given by equation (32).

When one moves from a configuration n to a new configuration m, the *transition* probability is the conditional probability, of observing the configuration m given that the configuration n was accepted, $P[g_m|g_n]$. Bayes theorem then gives

$$P[g_m|g_n]P[g_n] = P[g_n|g_m]P[g_m]. (36)$$

Therefore, we get an expression analogous to the Metropolis condition, equation (33),

$$\frac{P[g_m|g_n]}{P[g_n|g_m]} = \frac{P[g_m]}{P[g_n]} = \exp\left\{-(\mathcal{F}[g_m] - \mathcal{F}[g_n])/2\sigma_{\langle g \rangle}^2\right\}. \tag{37}$$

This justifies the use of the acceptance probability condition

$$P_{\text{acc}} = \min \left[1, \exp\{ -(\mathcal{F}_{\text{new}} - \mathcal{F}_{\text{old}}) / 2\sigma_{(g)}^2 \} \right], \quad (38)$$

as originally proposed by McGreevy [1]. However, the quantity $\sigma_{(g)}^2$ here is well defined, as the standard deviation of the fluctuations of the average rdf trial functions, about the target mean experimental rdf, during the RMC simulation. It is not to be confused with σ_e , the experimental error in the input rdf. In fact, one could have a very noisy experimental data with a large σ_e and get an excellent RMC fit with a low σ_g . Clearly, in that case, the obtained molecular configurations would have scarce physical meaning. For input data with low σ_e , as in tests we carried out with accurately calculated MMC simulations for Lennard–Jones (LJ) or water models in

papers I and II and in the next section, better results were reported for values of σ_g as low as 10^{-15} . The use of the average over the ensemble of configurations used in our SRMC simulations ensured good three particles correlation functions in the case of the LJ fluid and predicted nicely the hydrogen bond and and angular distributions in the case of the water models. In all studied cases, the statistical sampling tests carried out, like the linearity of mean squared displacements with the number of simulation cycles (i.e. the RMC *time*), were excellent. For such input data no additional constraints were needed.

The two criteria given by equations (34) and (38) should sample phase space differently. In the next section we study the MMC LJ model to analyze this point.

5. A comparison of the sampling in MMC and SRMC

In this section we would like to analyze the performance of the SRMC method, based on the variational prescription discussed above, in predicting physically valid configurations. For this purpose, we shall compare configurations obtained from a regular MMC simulation with those obtained with the corresponding SRMC, for a Lennard–Jones (LJ) fluid. As mentioned above, this kind of analysis, as well as the model system and simulation details, have been completely described elsewhere [3,4]. Briefly, accurate rdfs are generated by a MMC simulation, using the LJ potential; then the SRMC algorithm is used to generate trial configurations, using the MMC rdf as the experimental target.

For each trial configuration, accepted or rejected by the SRMC criterion, we also "tested" the Metropolis MC criterion, to check if that configuration would have been accepted or rejected by the MMC. In the later step, we calculated the variation of the total energy (ΔU) for new and old configurations and the corresponding Boltzmann weight. We also recorded the number of cases that both techniques would reject (or accept) this particular configuration. We also counted the cases, where one technique accepts the configuration and the other, does not. All these measurements were made during the production phase cycles and did not affect the SRMC configurational acceptance criterion. After a SRMC run, different probabilities were then calculated: (a) both SRMC and MMC reject the configuration $(P_{R,R})$, (b) SRMC rejects and MMC would accept $(P_{R,A})$, (c) SRMC accepts and MMC would reject $(P_{A,R})$ and (d) both SRMC and MMC accept the trial configuration $(P_{A,A})$. Several different runs were used to estimate the uncertainty in the probability calculations.

In SRMC, the rdfs were calculated from histograms obtained from *all* generated configurations. That is, at a particular configuration k, the histogram for the intermolecular distances, $H_k(r_i)$, contains $n_{\rm obs} = N_{\rm cycles} \times N(N-1)/2$ observations. At the beginning of the reverse simulation, the histograms have no physical significance, since no statistical averaging of the fluctuations have been

produced. In fact, the first histogram $H_0(r_i)$ contains only N(N-1)/2 observations corresponding to the number of pair interactions for N molecules, at their initial coordinates. Large fluctuations are therefore, expected in the value of the functional F[g(r)] at the beginning. Hence, there is a need to accumulate a large number of observations in the histogram used to compute the trial functions, $g_{\text{new}}(r)$ and $g_{\text{old}}(r)$. This is what it was called "memory effect" in SRMC, which gives the rdf its statistical meaning.

For the kth move, generating a new configuration k, an auxiliary incremental histogram h_k is generated, over *just* the N-1 distance counts (or observations). This histogram records just what one sees sitting only in the particle that one is trying to move. After k-1 such large number of observations, the ensemble histogram $H_{k-1} = \sum_{l=0}^{k-1} h_l$ contains the statistical memory effect. A new configuration is generated, and we can define

$$\left\langle g_k^{\text{new}}(r_i) \right\rangle = \frac{1}{n_k a_i} [H_{k-1} + h_k],$$
 (39)

where $a_i = 4\pi r_i^2 dr$ and $n_k = N(N-1)/2 + k(N-1)$.

In order to apply the functional variational principle numerically, equation (23), the functional differential of g(r), namely δg , must be clearly defined, for each trial move. From the definition of the incremental histogram h_k , it follows that δ_g , must be

$$\delta g = \frac{1}{n_k a_i} [h_k - h_{k-1}]. \tag{40}$$

Now, letting $\delta g = g^{\text{new}}(r) - g^{\text{old}}(r)$, is equivalent to define $g^{\text{old}}(r)$ as [3]

$$\langle g_k^{\text{old}}(r_i) \rangle = \frac{1}{n_k a_i} [H_{k-1} + h_{k-1}],$$
 (41)

Observe that, therefore, $[H_{k-1} + h_k]$ and $[H_{k-1} + h_{k-1}]$ correspond to $H_{\text{new}}(r)$ and $H_{\text{old}}(r)$, respectively, and that both have the same number of observations.

The functionals $\mathcal{F}[g(r)]^{\text{new}}$ and $\mathcal{F}[g(r)]^{\text{old}}$ are constructed with these functions and the kth configuration is accepted if

$$\Delta \mathcal{F}[g(r)] = \mathcal{F}[g(r)]^{\text{new}} - \mathcal{F}[g(r)]^{\text{old}} \le 0.$$
 (42)

Since the target *experimental* rdf is an accurate MMC value, the σ_g^2 was taken vanishingly small in equation (38).

In table 1, we present some results for specific conditions of a LJ fluid. MMC and SRMC simulations were performed at reduced densities $\rho^* = \rho \sigma^3$ ranging from 0.1 to 0.8, where ρ is the number density and σ is the LJ particle size. Both, MMC and SRMC, simulations were carried out with a small number of particles (N = 256). Tests with a larger N showed virtually the same outcomes [23].

The displacement parameter was chosen to give 50% of acceptance for MMC experiments, although for SRMC we found better convergence with a $\sim 40\%$ of acceptance. As in previous results [3,4], SRMC was more sensitive to the choice of the displacement parameter.

78.77(2)

8

System $-U_{MMC}/NkT$ $P_{A,A}$ $-U_{SRMC}/NkT$ $P_{R,R}$ $P_{R,A}$ $P_{A,R}$ 0.6925(1)94.12(2) 0.1 0.700(2)12.58(1)87.21(1) 5.88(2)1.3406(1) 1.356(8) 21.31(3) 78.70(2) 9.03(1)90.91(4) 0.2 3 0.3 1.8850(1)1.908(7) 27.70(2)72.30(2)11.13(2) 88.86(1) 0.4 2.3560(2)2.386(9)32.73(3) 67.27(3)12.70(4)87.30(4) 5 37.98(3) 0.5 2.8403(2) 2.879(3)62.02(3) 14.31(2) 85.69(2) 6 0.6 3.3667(2) 44.44(2) 55.57(2) 16.32(2) 83.68(2) 3.418(2)0.7 3.8920(2) 3.963(1) 51.81(3) 48.19(2) 18.69(3) 81.32(2)

Table 1. Analysis of the configurations in the production runs of a SRMC simulation for several densities of a LJ fluid.

The columns $-U_{SRMC}/NkT$ and $-U_{MMC}/NkT$ are the same as those in paper I [3]. The average SRMC reduced energies obtained assuming a LJ potential, are compared with the corresponding values obtained in the regular MMC simulation. The percentages of relative acceptance $P_{R,R}$, $P_{R,A}$, $P_{A,R}$ and $P_{A,A}$ are described in the text. The standard deviation of these numbers are indicated in parenthesis. They were obtained averaging five different simulation runs.

59.88(3)

4.468(2)

A total of 10^4 equilibration and production runs, without any hard-core constraints, were performed both for MMC and SRMC. A cut-off at half (L/2) of the cubic simulation box ($V=L^3$) was assumed, and a tail correction term was applied to correct the configurational energy U [33]. The reduced temperature was fixed at $T^* = k_{\rm B}T/\epsilon = 1.2$, which is below the liquid–gas critical temperature. Configurations collected during the runs were subjected of energy analysis, which is one initial way to assure that we obtained an appropriate sampling of configurations. As mentioned before, the same LJ potential, used in the MMC calculation, was assumed in order to calculate the averaged configurational energies ($U_{\rm SRMC}$) of the SRMC configurations generated in the production runs.

4.3677(3)

The agreement between SRMC and MMC calculations is very good with small relative error of about 2.0%. Here, we used bin sizes equal to 0.05 (in units of the LJ distance parameter). In paper I [3], it was pointed out that this difference could be reduced by decreasing the bin size. For example, for a density of 0.6, a bin size of 0.01 gives $U_{\text{SRMC}}/NkT0 - 3.416(2)$, in excellent agreement with the MMC value.

From the percentages of relative acceptance given in table 1, one can see that, in general, SRMC and MMC accept almost the same configurations. However, for low density systems, SRMC rejects many configurations that would be accepted by MMC. Just a few configurations are accepted by SRMC that would be rejected by MMC. This number is not more than 0.22, in the worst case (system label 8), for the very high density. Although there are similar trends when comparing $P_{A,R}$ and $P_{A,A}$, the behavior found for $P_{R,R}$ and $P_{R,A}$ indicates that each technique has its own way to search phase space. This might be taken as an indication that SRMC is not the reversed form of MMC. Nevertheless, both simulations provide configurations that, on the average, correspond to the same system.

6. Conclusions

We have revised the RMC technique in terms of a functional minimization, constructing an adequate

variational principle, which is in accordance with the basic ideas of goodness of fit and Bayes statistics. We have given a rigorous foundation to the acceptance/rejection criteria, establishing an analogy of the Metropolis importance sampling with Bayes theorem, in terms of the functional $\mathcal{F}[g(r)]$.

21.22(3)

40.10(2)

This formulation gives a formal basis for the newer version of RMC successfully used by da Silva and collaborators in the analysis of monatomic and polyatomic fluids [3,4]. We have named that version statistical RMC (SRMC) in order to point out the convenience of using a collection of randomly generated configurations, which gives an ensemble average $\langle g(r) \rangle$, as the trial functions that mimic the experimentally measured average $g^e(r)$. That is, the answer in SRMC is not a single configuration, but a large set of configurations which, on average renders the experimental target rdf.

We believe that our analysis unifies the existing approaches and dissipates possible controversies. For instance, in the case of experimental data with large errors or in the case of more complex systems, other order parameters can be incorporated to the variational formalism. Recently, Pikunic et al. [10] successfully introduced such weighting parameters, chosen with a trial-and-error procedure and simulated annealing, to minimize a linear combination of cost functions. So, in general, σ_{q}^{2} is much smaller than a typical experimental σ_e^2 . A Bayesian Metropolis-like criterium, as equation (38), could be used with a still smaller $\sigma_{\langle g \rangle}$. Alternatively, the standard deviation $\sigma_{\langle g \rangle}$ could be used as an adjusting parameter [10], since it is not necessary to be interpreted as the experimental error σ_e^2 . We believe that is the reason of the exit of the original inverse MC version of Kaplow [11] and of the extensively used improved MP-RMC versions [1,9].

Acknowledgements

We thank Drs. Bo Jönsson, T.Åkesson, Bo Svensson and Léo Degrève for stimulating discussions. We also acknowledge the *CNPq*, *FAPESP*/Brazil, CDCHT-ULA and FONACIT (G9700741)/Venezuela for the financial support during the development of this work.

Appendix A

In this appendix, we show that the functional derivative of the variational functional $\mathcal{F}[g(r)]$, given by equation (29), recovers the desired relationship, equation (23). We assume that g(r) is gaussian distributed around its average value $\bar{g}(r)$, with standard deviation σ_o^2 , i.e.:

$$P(g) = \frac{1}{\sqrt{2\pi}\sigma_g} \exp\left[-\frac{(g-\bar{g})^2}{2\sigma_g^2}\right]. \tag{A1}$$

Thus, from equation (29), the functional is

$$\mathcal{F}[g(r)] = \frac{1}{2} \int d\vec{\mathbf{r}} f(r) \left\{ \int \left[g(r) - \overline{g^{\ell}(r)} \right]^2 P(g) dg + \sigma_{\ell}^2 \right\}. \tag{A2}$$

The functional variation $\delta \mathcal{F}[g(r)]$ when the g(r) trial function changes in a differential $\delta g(r)$, is by definition equal to

$$\begin{split} \delta \mathcal{F}[g(r)] &= \mathcal{F}[g(r) + \delta g(r)] - \mathcal{F}[g(r)] \\ &= \frac{1}{2} \int \mathrm{d}\vec{\mathbf{r}} f(r) \bigg\{ \int \big[g(r) + \delta g - \overline{g^e(r)} \big]^2 \\ &\times P(g + \delta g) \mathrm{d}(g + \delta g) + \sigma_e^2 \bigg\} - \mathcal{F}[g(r)]. \end{split} \tag{A3}$$

The different terms in equation (A3), can be expanded to $O[\delta g]$

$$d(g + \delta g) = dg + d(\delta g) \equiv dg, \tag{A4}$$

$$\left[g(r) + \delta g - \overline{g^{e}(r)}\right]^{2} \approx \left[g(r) - \overline{g^{e}(r)}\right]^{2} + 2\left[g(r) - \overline{g^{e}(r)}\right] \delta g,$$
(A5)

$$P(g + \delta g)) \approx \frac{1}{\sqrt{2\pi}\sigma_g} \exp\left[-\frac{(g - \bar{g})^2}{2\sigma_g^2} - \frac{(g - \bar{g})}{\sigma_g^2}\delta g\right]$$
$$\approx P(g)\left(1 - \frac{(g - \bar{g})}{\sigma_g^2}\delta g\right). \tag{A6}$$

Replacing equations (A2) and (A4)–(A6) in equation (A3), the variation $\delta \mathcal{F}[g(r)]$, up to $O[\delta g]$, reads:

$$\delta \mathcal{F}[g(r)] = \int d\vec{r} f(r) \int P(g) dg \, \delta g \left\{ \left(g(r) - \overline{g^e} \right) - \frac{1}{2\sigma_g^2} \left(g(r) - \overline{g^e} \right)^2 (g(r) - \overline{g}) \right\}.$$

In the production cycles the functional $\mathcal{F}[g(r)]$ fluctuates around the minimum value and $\langle g(r) \rangle \approx \overline{g^e(r)}$, so the last term in equation (A7), vanishes, since it would

contain the skewness of the Gaussian variable. We then get

$$\delta \mathcal{F}[g(r)] = \int d\vec{\mathbf{r}} f(r) \langle g(r) - \overline{g^{e}(r)} \rangle_{g} \delta g. \tag{A8}$$

Therefore, the functional derivative of $\mathcal{F}[g(r)]$ vanishes, accordingly with equation (23).

References

- R.L. McGreevy, L. Pusztai. Reverse Monte Carlo simulation: a new technique for the determination of disordered structures. *Mol. Simul.*, 1, 359 (1988).
- [2] L. Degrève, F.L.B. da Silva, C. Quintale Jr., A.R. de Souza. Application of the Reverse Monte Carlo simulations to diatomic molecules. I: The noncomplete radial distribution functions. J. Molec. Struct. (Theochem), 335, 89 (1995).
- [3] F.L.B. da Silva, B. Svensson, T. Åkesson, B. Jönsson. A new algorithm for reverse Monte Carlo simulations. *J. Chem. Phys.*, 109, 2624 (1998).
- [4] F.L.B. da Silva, L. Degrève, W. Olivares-Rivas, T. Åkesson. Application of a new reverse Monte Carlo method to polyatomic molecules. 1. Liquid water case. J. Chem. Phys., 114, 907 (2001).
- [5] N.A. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A. Teller, E. Teller. Equation of state calculations by fast computing machines. J. Chem. Phys., 21, 1087 (1953).
- [6] K. Binder. Applications of the Monte Carlo methods to statistical physics. *Rep. Prog. Phys.*, 60, 487 (1997).
- [7] B.J. Alder, T.E. Wainwright. Phase transition for hard-sphere system. J. Chem. Phys., 27, 1208 (1957).
- [8] D. Frenkel, B. Smith. Understanding Molecular Simulation: From Algorithms to Applications, Academic Press, San Diego (1996).
- [9] R.L. McGreevy. Topical review: Reverse Monte Carlo modelling. J. Phys.: Condens. Matter, 13, R877 (2001).
- [10] J. Pikunic, C. Clinard, N. Cohaut, K.E. Gubbins, J.M. Guet, R.J.M. Pellenq, I. Rannou, J.N. Rouzaud. Structural modeling of porous carbons: Constrained reverse Monte Carlo method. *Langmuir*, 19, 8565 (2003).
- [11] R. Kaplow, T.A. Rowe, B.L. Averbach. Atomic arrangement in vitreous selenium. *Phys. Rev.*, 168, 1068 (1968).
- [12] G. Evrard, L. Pusztai. Reverse Monte Carlo modelling of the structure of disordered material with RMC++: A new implementation of the algorithm in C++. J. Phys.: Condens. Matter, 17, S1 (2005).
- [13] G. Tóth, A. Baranyai. Comparison of reverse Monte Carlo variants via the accuracy of their three-particle correlation functions. *Mol. Phys.*, 97, 339 (1999).
- [14] G. Tóth, A. Baranyai. Conceptual and technical improvement of the reverse Monte Carlo algorithm. J. Chem. Phys., 107, 7402 (1997).
- [15] A.K. Soper. Empirical potential Monte Carlo simulation of fluid structure. Chem. Phys., 202, 295 (1996).
- [16] A.K. Soper. Neutron scattering studies of solvent structure in systems of chemical and biochemical importance. *Faraday Discuss.*, 103, 41 (1996).
- [17] J.H. Finney. Hydration processes in biological and macromolecular systems. *Faraday Discuss.*, 103, 1 (1996).
- [18] A.P. Lyubartsev, A. Laaksonen. Calculation of effective interaction potentials from radial distribution functions: A reverse Monte Carlo approach. *Phys. Rev. E*, **52**, 3730 (1995).
- [19] M. Ostheimer, H. Bertagnolli. Test of the inverse Monte Carlo method for the calculation of the interatomic potential energies in atomic liquids. *Mol. Simul.*, 3, 227 (1989).
- [20] R.L. Henderson. A uniqueness theorem for fluid pair correlation functions. *Phys. Lett.*, 49A(3), 197 (1974).
- [21] M.A. Howe, R.L. McGreevy. Recent Developments in the Physics of Fluids, pp. f305–f312, International Symposium, Oxford (1991).
- [22] G. Tóth, L. Pusztai, A. Baranyai. A new algorithm for Reverse Monte Carlo simulations. J. Chem. Phys., 111, 5620 (1999).
- [23] F.L.B. da Silva, B. Svensson, T. Åkesson, B. Jönsson. A new algorithm for Reverse Monte Carlo simulations. *J. Chem. Phys.*, 111, 5622 (1999).
- [24] F.L.B. da Silva, L. Degrève. Application of a new Reverse Monte Carlo algorithm to polyatomic molecular systems. *Annals of 11th*

- Nordic Symposium on Computer Simulations, p. 34, Hillerød, Denmark (1997).
- [25] D.A. McQuarrie. Statistical Mechanics, Harper Collins, New York (1976).
- [26] W. Olivares-Rivas. General variational principles associated with the Ornstein–Zernike equation. Acta Cient. Venezolana, 29, 453 (1978).
- [27] A.M. Arthurs. Complementary Variational Principles, Clarendon Press, Oxford (1970).
- [28] P. Jedlovszky, I. Bako, G. Palinkas, T. Radnai, A.K. Soper. Investigation of the uniqueness of the reverse Monte Carlo method: Studies on liquid water. J. Chem. Phys., 105(1), 245 (1996).
- [29] S. Eidelman, *et al.* Review of particle physics: Probability and statistics. *Phys. Lett. B.*, **592**, 1 (2004).
- [30] R. Evans. Comment on reverse Monte Carlo simulations. Mol. Simul., 4, 409 (1990).
- [31] R. Webb. Functional approximation by feed forward networks: A least-squares approach to generalization. *IEEE Trans. Neural Netw.*, 5, 363 (1994).
- [32] G. Taguchi, M. El Sayed, C. Hsaing. *Quality Engineering and Production Systems*, McGraw-Hill, New York, USA (1989).
- [33] M.P. Allen, D.J. Tildesley. *Computer Simulation of Liquids*, Oxford University Press, Oxford (1989).