

## Adaptive Simulated Annealing-Useful Lessons Learned

<sup>1</sup>A. Iyem Perumal and <sup>2</sup>S.P. Rajagopalan

<sup>1</sup>Department of Mathematics, <sup>2</sup>School of Computer Science and Engineering,

Dr. M. G. R. University, Chennai-95, India

**Abstract:** Simulated Annealing (SA) is a powerful stochastic search method applicable to a wide range of problems. It can produce very high quality solutions for hard combinatorial optimization problems. SA can be generalized to fit non-convex cost-functions arising in a variety of problems which is known as Boltzmann Annealing (BA). The purpose of describing Simulated Quenching (SQ) and Fast Annealing (FA) is to highlight the importance of Adaptive Simulated Annealing (ASA). ASA is a global optimization algorithm based on an associated proof that the parameter space can be sampled much more efficiently than by using the previous Simulated algorithm.

**Key words:** Simulated Annealing (SA), Boltzmann Annealing (BA), Adaptive Simulated Annealing (ASA), Reannealing (RA), Simulated Quenching (SQ), Very Fast Simulated Reannealing (VFSR)

### INTRODUCTION

Adaptive Simulated Annealing (ASA) is a global optimization algorithm that relies on randomly important sampling the parameter space. The public release of the Very Fast Simulated Reannealing (VFSR) code (Ingber, 1989) later known as ASA code (Ingber, 1993a) is quite stable and free of bugs. This study deals with some lessons which may be useful to other developers of Simulated Annealing (SA) code as well as to many users.

### ASA ALGORITHMS

**“Standard” Simulated Annealing (SA):** The Metropolis Monte Carlo integration algorithm (Metropolis, 1953) was generalized by the Kirkpatrick algorithm to include a temperature schedule for efficient searching (Kirkpatrick *et al.*, 1983). A sufficiency proof was then shown to put a lower bound on that schedule as  $1/\log(t)$ , where  $t$  is an artificial time measure of the annealing schedule (Geman and Geman, 1984). However, independent credit usually goes to several other authors for independently developing the algorithm that is now recognized as simulated annealing (Cerny, 1982; Pincus, 1970).

**Boltzmann Annealing (BA):** Credit for the first simulated annealing is generally given to a Monte Carlo importance-sampling technique for doing large-dimensional path integrals arising in statistical physics problems (Metropolis, 1953). This method was generalized to fitting non-convex cost-functions arising in a variety of problems, e.g., finding the optimal wiring for a densely

wired computer chip (Kirkpatrick *et al.*, 1983). The choices of probability distributions described in this study are generally specified as Boltzmann annealing (Szu and Hartley, 1987).

The method of simulated annealing consists of three functional relationships:

- $g_T(x)$ : Probability density of state-space of  $D$  parameters  $x = \{x^i, i = 1, D\}$ , where the subscript  $T$  signifies a parameterization popularly referred to as the “temperature”.
- $h(\Delta E)$ : Probability for acceptance of new cost-function given the just previous value.
- $T(k)$ : Schedule of “annealing” the “temperature”  $T$  in annealing-time steps  $k$ , i.e., of changing the volatility or fluctuations of one or both of the two previous probability densities.

The acceptance probability is based on the chances of obtaining a new state with “energy”  $E_{k+1}$  relative to a previous state with “energy”  $E_k$ .

$$\begin{aligned} h(\Delta E) &= \frac{\exp(-E_{k+1}/T)}{\exp(-E_{k+1}/T) + \exp(-E_k/T)} \\ &= \frac{\exp(E_{k+1}/T) \exp(-E_{k+1}/T)}{\exp(E_{k+1}/T) [\exp(-E_{k+1}/T) + \exp(-E_k/T)]} \\ &= \frac{1}{1 + \exp[(E_{k+1} - E_k)/T]} \\ &= \frac{1}{1 + \exp(\Delta E/T)} \\ &\approx \exp(-\Delta E/T) \end{aligned} \quad (1)$$

Where  $\Delta E$  represents the “energy” difference between the present and previous values of the energies (considered here as cost functions) appropriate to the physical problem, i.e.,  $\Delta E = E_{k+1} - E_k$ . This essentially is the Boltzmann distribution contributing to the statistical mechanical partition function of the system (Binder and Stauffer, 1985).

This sampling algorithm also can be described by considering: A set of states labeled by  $x$ , each with energy  $e(x)$ ; a set of probability distributions  $p(s)$  and the energy distribution per state  $d(e(x))$ , giving an aggregate energy  $E$ ,

$$\sum_x p(x) d(e(x)) = E \quad (2)$$

The principle of maximizing the entropy,

$$S = -\sum_x p(x) \ln[p(x)/p(\bar{x})] \quad (3)$$

Where  $\bar{x}$  represents a reference state, using Lagrange multipliers (Mathews and Walker, 1970) to constrain the energy to average value  $T$ , leads to the most likely Gibbs distribution  $G(x)$ ,

$$G(x) = \frac{1}{Z} \exp(-H(x)/T) \quad (4)$$

in terms of the normalizing partition function  $Z$  and the Hamiltonian  $H$  operator as the “energy” function,

$$Z = \sum_x \exp(-H(x)/T) \quad (5)$$

For such distributions of states and acceptance probabilities defined by functions such as  $h(\Delta E)$ , the equilibrium principle of detailed balance holds. i.e., the distributions of states before,  $G(x_k)$  and after,  $G(x_{k+1})$ , applying the acceptance criteria,  $h(\Delta E) = h(E_{k+1} - E_k)$  are the same:

$$G(x_k) h(\Delta E(x)) = G(x_{k+1}) \quad (6)$$

This is sufficient to establish that all states of the system can be sampled, in theory. However, the annealing schedule interrupts equilibrium every time the temperature is changed and so, at best, this must be done carefully and gradually.

An important aspect of the SA algorithm is to pick the ranges of the parameters to be searched. In practice, computation of continuous systems requires some discretization, so without loss of much generality for applications described here, the space can be assumed to be discretized. There are additional constraints that are required when dealing with generating and cost functions

with integral values. Many practitioners use novel techniques to narrow the range as the search progresses. For example, based on functional forms derived for many physical systems belonging to the class of Gaussian-Markovian systems, one could choose an algorithm for  $g$ ,

$$g(\Delta x) = (2\pi T)^{-D/2} \exp[-(\Delta x)^2/(2T)] \quad (7)$$

Where  $\Delta x = x - x_0$  is the deviation of  $x$  from  $x_0$  (usually taken to be the just-previously chosen point), proportional to a “momentum” variable and where  $T$  is a measure of the fluctuations of the Boltzmann distribution  $g$  in the  $D$ -dimensional  $x$ -space. Given  $g(\Delta x)$ , it has been proven (Geman and Geman, 1984) that it suffices to obtain a global minimum of  $E(x)$  if  $T$  is selected to be not faster than

$$T(k) = \frac{T_0}{\ln k} \quad (8)$$

with  $T_0$  “large enough”.

For the purposes of this study, a heuristic demonstration follows, to show that Eq. 8 will suffice to give a global minimum of  $E(x)$  (Szu and Hartley, 1987). In order to statistically assure, i.e., requiring many trials, that any point in  $x$ -space can be sampled infinitely often in annealing-time (IOT), it suffices to prove that the products of probabilities of not generating a state  $\times$  IOT for all annealing-times successive to  $k_0$  yield zero,

$$\lim_{k \rightarrow \infty} \prod_k (1 - g_k) = 0 \quad (9)$$

This is equivalent to,

$$\lim_{k \rightarrow \infty} \sum_k g_k = \infty \quad (10)$$

as seen by taking the log of Eq. 9 and Taylor expanding in  $g_k$ . The problem then reduces to finding  $T(k)$  to satisfy Eq. 10. Note that, given a very large space to sample, often at best only a “weak” ergodicity can be assumed for this proof and any such ergodicity even for well-understood physical systems is an open area of research (Ma, 1985).

For BA, if  $T(k)$  is selected to be Eq. 8, then 7 gives

$$\sum_{k=k_0}^{\infty} g_k \geq \sum_{k=k_0}^{\infty} \exp(-\ln k) = \sum_{k=k_0}^{\infty} 1/k = \infty \quad (11)$$

Although there are sound physical principles underlying choices of Eq. 7 and 1 (Metropolis *et al.*, 1953),

it was noted that this method of finding the global minimum in x-space was not limited to physics examples requiring bona fide “temperatures” and “energies”. Rather, this methodology can be readily extended to any problem for which a reasonable probability density  $h(\Delta x)$  can be formulated (Kirkpatrick *et al.*, 1983).

**Simulated Quenching (SQ):** Many researchers have found it very attractive to take advantage of the ease of coding and implementing SA, utilizing its ability to handle quite complex cost functions and constraints. However, the long time of execution of standard Boltzmann-type SA has many times driven these projects to utilize a temperature schedule too fast to satisfy the sufficiency conditions required to establish a true (weak) ergodic search. A logarithmic temperature schedule is consistent with the Boltzmann algorithm, e.g., the temperature schedule is taken to be

$$T_k = T_0 \frac{\ln k_0}{\ln k} \quad (12)$$

Where  $T$  is the “temperature”,  $k$  is the “time” index of annealing and  $k_0$  is some starting index. This can be written for large  $k$  as

$$\Delta T = -T_0 \frac{\ln k_0 \Delta k}{k(\ln k)^2}, k \gg 1$$

$$T_{k+1} = T_k - T_0 \frac{\ln k_0}{k(\ln k)^2} \quad (13)$$

However, some researchers using the Boltzmann algorithm use exponential schedules, e.g.,

$$T_{k+1} = cT_k, 0 < c < 1$$

$$\frac{\Delta T}{T_k} = (c-1)\Delta k, k \gg 1$$

$$T_k = T_0 \exp((c-1)k) \quad (14)$$

with expediency the only reason given.

**Fast Annealing (FA):** Although there are many variants and improvements made on the “standard” Boltzmann algorithm described above, many textbooks finish just about at this point without going into more detail about other algorithms that depart from this explicit algorithm (Van Laarhoven and Aarts, 1987). Specifically, it was noted that the Cauchy distribution has some definite advantages over the Boltzmann form (Szu and Hartley, 1987). The Cauchy distribution they define is

$$g(\Delta x) = \frac{T}{((\Delta x)^2 + T^2)^{(D+1)/2}} \quad (15)$$

which has a “fatter” tail than the Gaussian form of the Boltzmann distribution, permitting easier access to test local minima in the search for the desired global minimum. It is instructive to note the similar corresponding heuristic demonstration, that the Cauchy  $g(\Delta x)$  statistically finds a global minimum. If Eq. 8 is replaced by

$$T(k) = \frac{T_0}{k} \quad (16)$$

then here

$$\sum_{k_0}^{\infty} g_k \approx \frac{T_0}{\Delta x^{D+1}} \sum_{k_0}^{\infty} \frac{1}{k} = \infty \quad (17)$$

Note that the “normalization” of  $g$  has introduced the annealing-time index  $k$ , giving some insights into how to construct other annealing distributions. The method of FA is thus seen to have an annealing schedule exponentially faster than the method of BA. This method has been tested in a variety of problems (Szu and Hartley, 1987).

**Adaptive Simulated Annealing (ASA):** In a variety of physical problems we have a  $D$ -dimensional parameter-space. Different parameters have different finite ranges, fixed by physical considerations and different annealing-time-dependent sensitivities, measured by the curvature of the cost-function at local minima. BA and FA have  $g$  distributions which sample infinite ranges and there is no provision for considering differences in each parameter-dimension, e.g., different sensitivities might require different annealing schedules. These are among several considerations that gave rise to Adaptive Simulated Annealing (ASA). Full details are available by obtaining the publicly available source code (Ingber, 1993a).

ASA considers a parameter  $\alpha_k^i$  in dimensions  $i$  generated at annealing-time  $k$  with the range

$$\alpha_k^i \in [A_i, B_i], \quad (18)$$

calculated with the random variable  $y^i$ ,

$$\alpha_{k+1}^i = \alpha_k^i + y^i(B_i - A_i),$$

$$y_i \in [-1, 1] \quad (19)$$

Define the generating function

$$g_T(y) = \prod_{i=1}^D \frac{1}{2(|y^i| + T_i) \ln(1 + 1/T_i)} = \prod_{i=1}^D g_T^i(y^i) \quad (20)$$

Where the subscript  $i$  on  $T_i$  specifies the parameter index and the  $k$ -dependence in  $T_i(k)$  for the annealing schedule has been dropped for brevity. Its cumulative probability distribution is

$$G_T(y) = \int_{-1}^{y^1} \int_{-1}^{y^2} \dots \int_{-1}^{y^D} g_T(y') = \prod_{i=1}^D G_T^i(y^i) \\ G_T^i(y^i) = \frac{1}{2} + \frac{\text{sgn}(y^i) \ln(1 + |y^i|/T_i)}{2 \ln(1 + 1/T_i)} \quad (21)$$

$y^i$  is generated from a  $u^i$  from the uniform distribution  $u^i \in U[0,1]$

$$y^i = \text{sgn}\left(u^i - \frac{1}{2}\right) T_i [(1 + 1/T_i)^{|2u^i - 1|} - 1] \quad (22)$$

It is straightforward to calculate that for an annealing schedule for  $T_i$

$$T_i(k) = T_{0i} \exp(-c_i k^{1/D}) \quad (23)$$

a global minima statistically can be obtained. I.e.,

$$\sum_{k_0}^{\infty} g_k \approx \sum_{k_0}^{\infty} \left[ \prod_{i=1}^D \frac{1}{2|y^i|c^i} \right] \frac{1}{k} = \infty \quad (24)$$

It seems sensible to choose control over  $c_i$ , such that

$$T_{0i} = T_{0i} \exp(-m_i) \text{ when } k_f = \exp n_i, \\ c_i = m_i \exp(-n_i/D) \quad (25)$$

Where  $m_i$  and  $n_i$  can be considered “free” parameters to help tune ASA for specific problems.

It has proven fruitful to use the same type of annealing schedule for the acceptance function  $h$  as used for the generating function  $g$ , i.e., Eq. 23 and 25, but with the number of acceptance points, instead of the number of generated points, used to determine the  $k$  for the acceptance temperature.

In one implementation of this algorithm, new parameters  $\alpha_{k+1}^i$  are generated from old parameters  $\alpha_k^i$  by generating the  $y$ 's until a set of  $D$  are obtained satisfying the range constraints. In another alternative supported in ASA, useful for some constraint problems, the  $y$ 's are generated sequentially for each test of the cost function.

**Reannealing:** Whenever doing a multi-dimensional search in the course of a real-world nonlinear physical problem, inevitably one must deal with different changing sensitivities of the  $\alpha^i$  in the search. At any given annealing-time, it seems sensible to attempt to “stretch out” the range over which the relatively insensitive parameters are being searched, relative to the ranges of the more sensitive parameters.

It has proven fruitful to accomplish this by periodically rescaling the annealing-time  $k$ , essentially reannealing, every hundred or so acceptance-events (or at some user-defined modulus of the number of accepted or generated states), in terms of the sensitivities  $s_i$ , calculated at the most current minimum value of the cost function,  $\underline{L}$ ,

$$s_i = \partial \underline{L} / \partial \alpha^i \quad (26)$$

In terms of the largest  $s_i = s_{\max}$ , a default rescaling is performed for each  $k_i$  of each parameter dimension, whereby a new index  $k'_i$  is calculated from each  $k_i$ ,

$$k_i \rightarrow k'_i \\ T'_{ikg} = T_{ik} (s_{\max}/s_i), \\ k'_i = (\ln(T_{i0}/T'_{ikg})/c_i)^D \quad (27)$$

$T_{i0}$  is set to unity to begin the search, which is ample to span each parameter dimension.

The acceptance temperature is similarly rescaled. Since the initial acceptance temperature is set equal to an initial trial value of  $L$ , this is typically very large relative to the current best minimum, which may tend to distort the scale of the region currently being sampled. Therefore, when this rescaling is performed, the initial acceptance temperature is reset to the maximum of the most current minimum and the best current minimum of  $\underline{L}$  and the annealing-time index associated with this temperature is reset to give a new temperature equal to the minimum of the current cost-function and the absolute values of the current best and last minima.

Also generated are the “standard deviations” of the theoretical forms, calculated as  $[\partial^2 \underline{L} / (\partial \alpha^i)^2]^{-1/2}$ , for each parameter  $\alpha_i$ . This gives an estimate of the “noise” that accompanies fits to stochastic data or functions. At the end of the run, the off-diagonal elements of the “covariance matrix” are calculated for all parameters. This inverse curvature of the theoretical cost function can provide a quantitative assessment of the relative sensitivity of parameters to statistical errors in fits to stochastic systems.

**Quenching:** Another adaptive feature of ASA is its ability to perform quenching in a methodical fashion. This is applied by noting that the temperature schedule above can be redefined as

$$\begin{aligned} T_i(k_i) &= T_{0i} \exp(-c_i k_i^{Q_i/D}) \\ c_i &= m_i \exp(-n_i Q_i/D) \end{aligned} \quad (28)$$

in terms of the “quenching factor”  $Q_i$ . The above proof fails at Eq. 24 if  $Q_i > 1$  as

$$\sum_k \prod_{i=1}^D 1/k^{Q_i/D} = \sum_k 1/k^{Q_i} < \infty \quad (29)$$

This simple calculation shows how the “curse of dimensionality” arises and also gives a possible way of living with this disease. In ASA, the influence of large dimensions becomes clearly focused on the exponential of the power of  $k$  being  $1/D$ , as the annealing required to properly sample the space becomes prohibitively slow. So, if we cannot commit resources to properly sample the space ergodically, then for some systems perhaps the next best procedure would be to turn on quenching, whereby  $Q_i$  can become on the order of the size of number of dimensions.

The scale of the power of  $1/D$  temperature schedule used for the acceptance function can be altered in a similar fashion. However, this does not affect the annealing proof of ASA and so this may be used without damaging the (weak) ergodicity property.

**ASA applications:** The above defines this method of Adaptive Simulated Annealing (ASA), previously called Very Fast Simulated Reannealing (VFSR) (Ingber, 1989) only named such to contrast it the previous method of fast annealing (Szu and Hartley, 1987). The annealing schedules for the temperatures  $T_i$  decrease exponentially in annealing-time  $k$ , i.e.,  $T_i = T_{i0} \exp(-c_i k^{1/D})$ . Of course, the fatter the tail of the generating function, the smaller the ratio of acceptance to generated points in the fit. However, in practice, it is found that for a given generating function, the ratio is approximately constant as the fit finds a global minimum. Therefore, for a large parameter space, the efficiency of the fit is determined by the annealing schedule of the generating function.

A major difference between ASA and BA algorithms is that the ergodic sampling takes place in an  $n+1$  dimensional space, i.e., in terms of  $n$  parameters and the cost function. In ASA the exponential annealing schedules permit resources to be spent adaptively on reannealing and on pacing the convergence in all

dimensions, ensuring ample global searching in the first phases of search and ample quick convergence in the final phases.

## ASA OPTIONS

ASA likely is the most powerful and flexible SA code presently available, because the code has benefited from the feedback of many users and their feedback has been used to add much to the code beyond the basic ASA algorithm described above.

The code has two basic modules in the ASA C-code, a user and an asa module. All options in the code have been tested to work with templates provided in the user module. Feedback has developed a code which seems to run well across many platforms, e.g., PC's, Macs, Crays, many UNIX workstations, etc.

The emphasis in development of ASA has been to add power and flexibility wherever possible. To make these extra features and code accessible to non-expert programmers, a “meta-language” of options is used. Many of these options can be set in the provided Makefile, an asa\_opt data file from which to read in information, arguments passed to the compilation procedures, or in the user module files.

Adequate investment has been made for continual development of a more powerful and more flexible code. The new user is presented with many options, on the order of a hundred. In many cases, when the ASA default options work fine, only the user's own call to his/her cost function is required. However, if these defaults are not suitable for a particular system, then the user can become bewildered by the many options. If not much is known *a priori* about the system to be optimized, then the task is to try to find the values of the options appropriate to the given system. The less known about the system, the harder is this task.

Experiences support the premise that the output of the code, using the ASA\_PRINT\_MORE OPTIONS to give information at each new best accepted state, often can be used to diagnose problems in annealing. Eventually, we hope that enough experience will be generated, to be able to develop some kind or graphical menu-driven expert system to help guide users to optimize a wide range of cost functions.

**Examples of options:** The following discussion of some of the options available in ASA also serves to illustrate the typical kinds of problems many users have with their particular systems and some of the approaches that SA can offer to face these problems. The options are organized into three groups. The define\_options

comprise two sets of options, the Pre-Compile define\_options and the Printing define\_options, which are called at the time of compilation; these comprise about half of the options. The other Program options are housed in a structure passed with the cost function and together with the other parameters passed in the cost function, these can be modified adaptively. That is, they can be changed within the cost function to take effect upon reentering the ASA program.

**Integer and continuous parameters:** ASA can accommodate mixture of integer and continuous parameters. This is accomplished quite simply, with a small overhead for integers, by truncating generated floating-point numbers within sensible integral windows. There have been many rumours that SA can only handle integer or continuous parameters, but these statements are unsupported.

**Constraints:** One of the immediate attractions of SA to people trying to optimize complex systems is the ease with which SA can accommodate complex constraints. Typically, there is no need for penalty functions, etc. Generated points that do not satisfy the constraints are simply rejected before trying any acceptance test.

Equality constraints, if processed will present a problem for any global optimization that relies on sampling, because the search is being constrained on the surface of some volume and the entire volume is being sampled. Therefore, it is recommended that the user first numerically substitute solution(s) of the equalities for some parameters. For example, if the cost function  $C$  has  $n$  parameters,  $C(p_1, p_2, \dots, p_n)$  and an equality constraint exists between parameters  $p_n$  and  $p_{n-1}$ , then solve this equation for  $p_n$ , numerically or algebraically, redefining the cost function to one with  $n-1$  parameters,  $C'$ . If the solution to this equation, or perhaps a set of  $m$  such equality constraints to reduce the number of parameters actually processed by ASA to  $n-m$ , is not simply written down, then such constraints must be solved with other algorithms within the cost function.

**Annealing scales:** Perhaps the easiest to understand problem that can arise when using SA, also is the most often neglected. A question may arise that why ASA doesn't immediately find the global optimal point? The answer most often lies in the scaling parameters used in annealing the parameter and/or cost temperatures.

For example, if the search is carried out in a system with several local minima, but the temperature is too low so that only rarely can the search sample these minima, it may take an extremely long time with arbitrarily good

numerical precision to eventually sample these minima as normal annealing proceeds to lower and lower temperatures. Clearly, it would be better to have the starting temperature at the scale in question be commensurately larger and perhaps be cooled more slowly.

**Parameter temperatures:** In some SA algorithms, like BA, the starting temperature controls the values of temperature encountered at subsequent stages of search. In ASA, because of the finite ranges of the parameters, the parameter temperatures are started to establish to a fat tail throughout the range; the exponential annealing rates usually permit selecting even quite large initial ranges to be sure of covering all optima. There are free ASA-parameters for each temperature to scale its exponential decrease, without affecting the basic sampling proof.

**Cost temperature:** The annealing scale for the cost temperature, also called the acceptance temperature, affects the rate of narrowing the window of the Boltzmann acceptance test. In ASA, this scale can be adaptively changed and even the Boltzmann test can be changed to a different distribution.

### Reannealing

**Parameter temperatures:** For the parameter temperatures, the tangents (or any other alternative functions that might be defined by the user) are used as a relative measure of the "steepness" of each dimension the most recent best saved state. As demonstrated for the ASA\_TEST problem (Ingber, 1993b), this feature can enhance the efficiency of the search.

**Cost temperature:** For the cost temperature, a separate options permits rescaling of the cost temperature to be set to the scale of the minimum of the current cost temperature and the absolute values of the last and best saved minima, to keep the acceptance test sensitive at a reasonable scale. This can be extremely important of the system's terrain changes with the scale of the search. This procedure also may need to radically altered, possible with other options, if the search early becomes struck in local optima, e.g., because the system's terrain abruptly changes with the scale of the search.

**Quenching:** An SA algorithm loses much of its authority if the search "cheats" by trying to anneal at rates faster than permitted by its associated proof, e.g., Simulated Quenching (SQ). However, this can be useful in a number of circumstances (Ingber, 1993b).

When the dimension of a parameter space, each parameter having a continuous or large integral set of values, reaches 15-20, the volume of search typically becomes quite large and this can severely tax most present-day workstations. Instead of just giving up on SA and trying a different “greedy” and/or quasi-Newton algorithm, ASA provides a methodical way to deviate from SA into SQ algorithms.

As another use of quenching, one that does not necessarily violate any sampling proof, it may be useful in the course of search to adaptively drop subsets of parameters that seem to have been reasonably optimized relative to other parameters. The remaining parameters can then be more efficiently searched within their smaller dimensional space, by adjusting the dependence of the annealing to the new dimension. This can be accomplished conveniently with the quenching options.

**ASA sampling:** Since ASA accomplishes its fit by importance sampling the space of parameters, it would seem that this process should provide a good sampling technique for other purposes, e.g., performing integrals. As stated above, the use of Monte Carlo techniques for performing integrals (Metropolis *et al.*, 1953) is generally credited to be the origin of the development of SA (Kirkpatrick *et al.*, 1983). However, importance sampling with the fastest permitted temperature schedules often can lead to quite poor resolutions of local minima which may substantially contribute to integrals. Then, the rates of annealing must be slowed down, e.g., using inverse quenching, to get better resolution. The ASA\_SAMPLE options collect the generating and acceptance biases incurred during importance sampling, so that this information can be used more generally than for just finding the optimal point of the fit.

**Self optimization:** An advantage of C code over some other languages is the relative ease by which recursive calls can be implemented. Some care must be taken to keep variables local to each subroutine. In its current form ASA can recursively call itself. Some complex problems, possessing nests of optimized systems, require this.

If not much information is known about a particular system, if the ASA defaults do not seem to work very well and if after a bit of experimentation it still is not clear how to select values for some of the ASA options, then the self-optimize options can be very useful. This sets up a top level search on the ASA options themselves, using criteria of the system as its own cost function, e.g., the best attained optimal value of the system’s cost

function (the cost function for the actual problem to be solved) for each given set of top level options, or the number of generated states required to reach a given value of the system’s cost function, etc. Since this can consume a lot of CPU resources, it is recommended that only a few ASA options and a scaled down system cost function or system data be selected for this options.

Even if good results are being attained by ASA, self\_optimize can be used to find a more efficient set of ASA options. I think that this kind of options would be useful for many non-linear optimization algorithms. Many of the options broken out in clear view in ASA are similarly represented but “hidden” within the code of other algorithms. Self optimization of such parameters can be very useful for production runs of complex systems.

**Alternative distributions/functions:** There are options to permit replacing or modifying the functions and distributions used in the ASA module. For example, modifications can be made of the generating function (e.g., variants of the Boltzmann and Cauchy distributions are given in the user module), the acceptance function (e.g., a class of functions that asymptotically approach the Boltzmann function is given in the user module) and the reannealing functions used to rescale the parameter and cost temperatures.

**Parallel code:** It is quite difficult to directly parallelize an SA algorithm (Ingber, 1993b) e.g., without incurring very restrictive constraints on temperature schedules (Kimura and Taki, 1991) or violating an associated sampling proof (Frost, 1993). However, the fat tail of ASA permits parallelization of developing generated states prior to subjecting them to the acceptance test (Ingber, 1992). The ASA\_parallel options provide parameters to easily parallelize the code, using various implementations, e.g., shared memory.

The scale of parallelization afforded by ASA, without violating its sampling proof, is given by a typical ratio of the number of generated to accepted states. Several experts in parallelization suggest that massive parallelization e.g., on the order of the human brain, may take place quite far into the future, that this might be somewhat less useful for many applications than previously thought and that most useful scales of parallelization might be on scales of order 10-1000. Depending on the specific problem, such scales are common in ASA optimization and the current ASA code can implement such parallelization.

## CONCLUSION

If asked to state one major common feature of nonlinear system in the context of optimization, the feature most likely should be given is that nonlinear system typically are non-typical. It is unlikely that any “canned black-box” code can be developed, requiring no or few minor adjustments, that will usefully guarantee efficient global optimization for severely nonlinear system, e.g., similar to what might be expected for many quasi-linear system.

In the absence of knowledge about a particular system, given that only SA can offer at least a “statistical” proof of global optimization, then the first algorithm of choice clearly is SA. Modification of SA, e.g., SQ quenching algorithms, may be competitive with other techniques, e.g., simplex or genetic algorithms, but among these the best choice is not so clear. SQ does offer a relatively simple approach to quickly writing code for optimization, but ultimately the end results must justify this means.

This argument for the use of SA has an opposite side. If some information about a system can be incorporated into some other global optimization technique and it can be determined that the technique can deliver the global optimum point, often that technique will be more efficient than SA. E.g., a quasi-Newton algorithm will be more efficient than SA for parabolic systems.

For many researchers, the first choice of algorithm to use for a nonlinear or stochastic problem likely will be one with which they already are familiar, if that fails, SA is an option to try next. More research needs to be done to see if a more objective classification of nonlinear system can be developed to help guide a given researcher to a given algorithm for a given problem. As the examples included in the documentation of the ASA code illustrate, there have been “surprises” whereby some very difficult problems have been quickly solved by ASA, while others have required quite a bit of “tuning” to establish a good set of starting options.

Especially among first users of SA, often there is much misunderstanding and lack of appreciation of just what an SA code can immediately do for a particular problem. Some education is necessary to make users aware of the potential problems that may arise and what remedies the particular algorithm can offer to overcome these obstacles.

## REFERENCES

- Binder, K. and D. Stauffer, 1985. A Simple Introduction to Monte Carlo Simulations and Some Specialized Topics, In: Applications of the Monte Carlo Method in Statistical Physics, K. Binder (Ed.), Berlin, Springer-Verlag.
- Cerny, V., 1982. A thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm, Report, Bratislava, Czechoslovakia, Comenius University.
- Frost, R., 1993. Ensemble Based Simulated Annealing (EBSA), [ftp.sdsc.edu/pub/sdsc/math/Ebsa](http://ftp.sdsc.edu/pub/sdsc/math/Ebsa), La Jolla, CA, University of California San Diego.
- Geman, S. and D. Geman, 1984. Stochastic relaxation, Gibbs distribution and the Bayesian restoration in images, *IEEE. Trans. Patt. Anal. Mac. Int.*, 6: 721-741.
- Ingber, L., 1989. Very fast simulated re-annealing, *Mathl. Comput. Modelling*, 12: 967-973.
- Ingber, L., 1992. Generic mesoscopic neural networks based on statistical mechanics of neocortical interactions, *Phys. Rev. A.*, 45: 2183-2186.
- Ingber, L., 1993a. Adaptive Simulated Annealing (ASA), [[ftp.alumni.Caltech.edu/pub/ingber/ASA-shar](http://ftp.alumni.Caltech.edu/pub/ingber/ASA-shar), [ASA-shar.Z](http://ftp.alumni.Caltech.edu/pub/ingber/ASA-shar.Z), [ASA-shar.tar.Z](http://ftp.alumni.Caltech.edu/pub/ingber/ASA-shar.tar.Z), [ASA-shar.tar.gz](http://ftp.alumni.Caltech.edu/pub/ingber/ASA-shar.tar.gz), [ASA-shar.zip](http://ftp.alumni.Caltech.edu/pub/ingber/ASA-shar.zip)], McLean, VA, Lester Ingber Research.
- Ingber, L., 1993b. Simulated annealing: Practice versus theory, *Mathl. Comput. Modelling*, 18: 29-57.
- Kimura, K. and K. Taki, 1991. Time-homogeneous parallel annealing algorithm, Report TR-673, Tokyo, Japan, Institute for New Generation Computer Technology.
- Kirkpatrick, S., C.D.J.R. Gelatt and M.P. Vecchi, 1983. Optimization by simulated annealing, *Science*, 220: 4598, 671-680.
- MA, S.K., 1985. Statistical Mechanics, Philadelphia, World Scientific.
- Mathews, J. and R.L. Walker, 1970. Mathematical Methods of Physics, (2nd Edn.), New York, Benjamin.
- Metropolis, N., 1953. Equation of state calculations by fast computing machines, *J. Chem. Phys.*, 21: 1087-1092.
- Pincus, M., 1970. A Monte Carlo method for the approximate solution of certain types of constraint optimization problems, *Oper. Res.*, 18: 1225-1228.
- Szu, H. and R. Hartley, 1987. Fast simulated annealing, *Phys. Lett. A.*, 122: 157-162.
- Van Laarhoven, P.J.M. and E.H.L. Aarts, 1987. Simulated Annealing: Theory and Applications, Dordrecht, The Netherlands, D. Reidel.