

Stochastic Optimisation

March 15, 2011

1 Introduction

In this coursework, the problem of global optimisation of Keane's Bump was solved with Biased Monte Carlo Sampling (BMCS) and with Simulated Annealing (SA) algorithms. The performance of the algorithms was investigated, as was the effect of different parameters on the performance.

2 Problem Description

$$\begin{aligned} \text{Maximise} \quad & f(\mathbf{x}) = \left| \frac{(\cos x_1)^4 + (\cos x_2)^4 - 2(\cos x_1)^2 (\cos x_2)^2}{\sqrt{x_1^2 + 2x_2^2}} \right| \\ \text{subject to} \quad & x_1, x_2 \in [0, 10] \\ & x_1 x_2 > 0.75 \\ & x_1 + x_2 < 15 \end{aligned}$$

The surface of the feasible region is shown in Figure 1. Since we can see that the global optimum is on the $x_1 x_2 > 0.75$ boundary, we can locate it by searching along that curve. This tells us that the optimum is located at $[1.6009 \ 0.4685]^T$, where the objective function is 0.36498. There are many other local optima, both on and off the constraint boundaries.

3 Techniques

3.1 Quantifying Performance

To investigate the performance of an algorithm with given parameters, it is run n times with the random seeds $1, 2, \dots, n$, and each time the best solution found is recorded. To get an impression of the performance, we can consider the mean μ and standard deviation σ of the results. This can be likened to approximating the distribution of the results as a Gaussian, and matching the first and second moments. This approximation can be very bad, as shown in Figure 2, but it serves to indicate the performance of the algorithm. In the cases shown, the Gaussians predict 6% and 38% chances of a result exceeding the global optimum.

If we trust the Gaussian approximation, μ and σ are sufficient to calculate probabilistic bounds for the performance of the algorithm. For example, in 95% of cases the result would be above $\mu - 1.645\sigma$. It would be more accurate to directly estimate the 5th percentile from the data and perhaps this should have been used instead. However, one feature of the Gaussian approach is that it penalises rare, catastrophic failures heavily. For the SA results in Figure 2, the approximated and actual 5th percentiles are 0.338 and 0.358, respectively.

Another advantage of calculating the standard deviation is that it can be used to determine the accuracy of the calculated mean. By the Central Limit Theorem, as $n \rightarrow \infty$ the distribution of μ tends to a normal distribution, $N(\mu_{\text{true}}, n^{-1}\sigma_{\text{true}}^2)$. Thus, with a probability of around 95%, $\mu_{\text{true}} \in [\mu - 2n^{-1/2}\sigma, \mu + 2n^{-1/2}\sigma]$. This gives us a statistical confidence bound when we are drawing conclusions about performance.

Two concerns are ignored in the confidence bound above. Firstly, we assume $\sigma_{\text{true}} \approx \sigma$ —this error will hopefully become negligible for large n . Secondly, if we are selecting the results with the highest mean, this will tend to select evaluations with $\mu > \mu_{\text{true}}$ and will skew the distributions of the estimated μ .

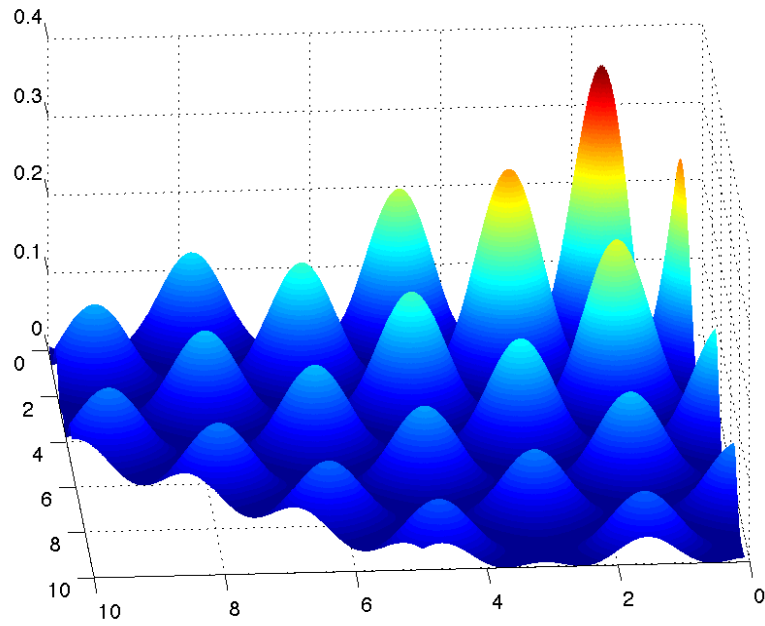


Figure 1: Keane's 2-D Bump

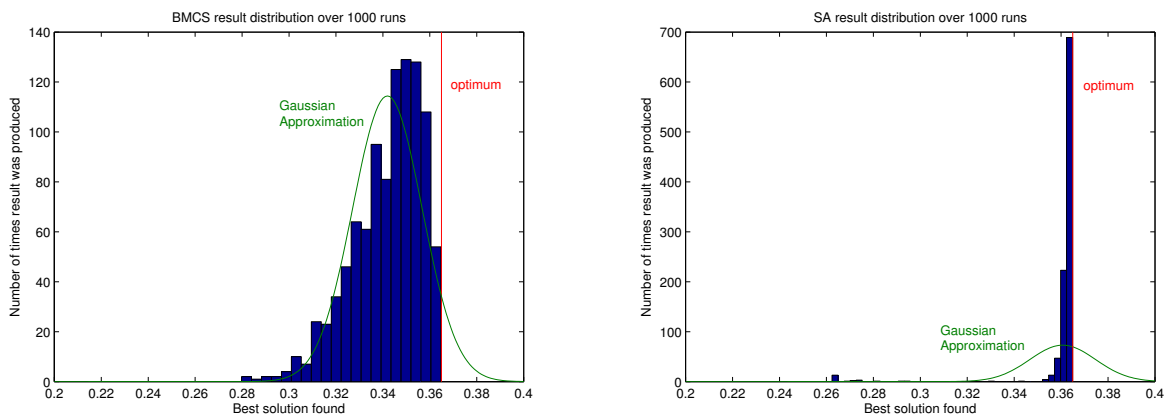


Figure 2: Result distributions and Gaussian approximations

3.2 Visualising Performance

By estimating μ and σ for different parameter settings, we can see how the parameter settings affect the performance, and we can determine the parameters that result in the best performance. This is effectively a multidimensional optimisation of a noisy objective. However, the purpose of this report is not to find the best parameters for this problem, but to gain an understanding of how they affect the performance. For this purpose, we need to be able to see how a single parameter changes the performance, but this is complicated by the settings of the other parameters—Figure 3 shows how complicated this is. To take this into account, for a given setting of one parameter we can choose the other parameters to optimise the performance. This allows us to view the performance as a function of a single variable, which makes it easier to visualise and explain.

More formally, given a performance measure $o(\theta_1, \theta_2, \dots, \theta_k)$, we consider the k different functions:

$$o_i(\theta_i) = \max_{\theta_1 \dots \theta_{i-1}, \theta_{i+1} \dots \theta_k} o(\theta_1 \dots \theta_k)$$

3.3 Penalty Function

The same penalty function was used in both the BMCS and SA algorithms:

$$\begin{aligned} c(\mathbf{x}) = & \min(x_1 - 0, 0) + \\ & \min(10 - x_1, 0) + \\ & \min(x_2 - 0, 0) + \\ & \min(10 - x_2, 0) + \\ & \min(15 - x_1 - x_2, 0) + \\ & \min(x_1 x_2 - 0.75, 0) \end{aligned}$$

It is scaled by a configurable weight, given by the `penalty_factor` variable. In general, we could use separate weights for every constraint, but since the variable scales are so similar in this problem, this possibility was not investigated.

3.4 Archiving

Archiving was implemented exactly as described in the lecture notes [Parks, 2011]. The same archiving code was used for the BMCS and SA implementations, and the algorithm implementations themselves are independent of the number and type of archives used. One downside of this low-coupling approach is that an implementation might want to use a dissimilarity archive when restarting the search. This was avoided by having the SA implementation keep track of the single best solution seen so far, and restart from there if necessary.

4 Biased Monte Carlo Sampling

The BMCS algorithm is much like that shown in lectures [Parks, 2011]. Each axis range is divided into m equally sized ranges, for a total of m^2 regions. To determine initial region probabilities it samples a certain number of times per region, then uses the average objective observed in each region to rank them. The linear selection probability relationship is used to translate the ranks to probabilities. 1000 random samples are made according to these probabilities, then the probabilities are recalculated.

The penalty function is used in two ways: firstly to detect invalid solutions and to avoid archiving them, and secondly when calculating the average objective in a region (for determining selection probability).

The following parameters of the algorithm were investigated:

Parameter	Description
<code>penalty_factor</code>	A constant weight used for the penalty function
<code>m</code>	The number of divisions of each axis
<code>initial_samples</code>	The number of samples to take, per region, before calculating region probabilities
<code>pressure</code>	The selection pressure, S , used when calculating region probabilities.

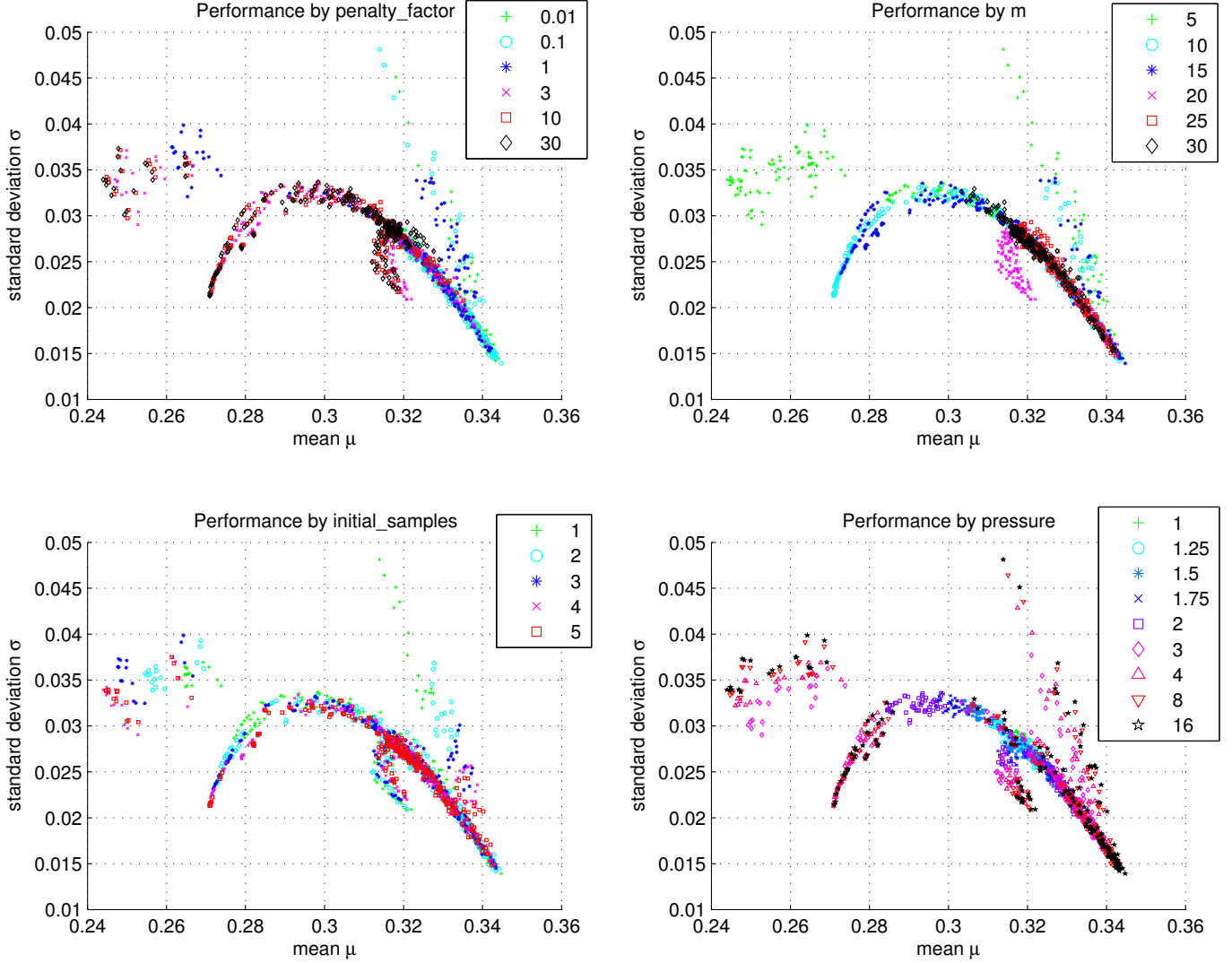


Figure 3: BMCS performance with varying parameters

4.1 Searching the Parameter Space

To investigate the effect of various parameters, the following parameter values were chosen for investigation:

Parameter	Values					
penalty_factor	0.01	0.1	1	3	10	30
m	5	10	15	20	25	30
initial_samples	1	2	3	4	5	
pressure	1	1.25	1.5	1.75	2	3 4 8 16

The algorithm was run for every combination of these values: a total of 1,620 combinations. For each combination, it was run $n = 1000$ times, as described in Section 3.1. Figure 3 shows the results.

These plots are rather complicated. They show some striking patterns (especially the variation with m) but it is not clear what is causing these. Some insights that can be drawn from these plots:

- The large bump visible can be explained by analogy to a Bernoulli random variable, with $\mu = p$ and $\sigma = \sqrt{p(1-p)}$, where the two values of the variable correspond to the two largest optima of the function, peaking at 0.263 and 0.365.

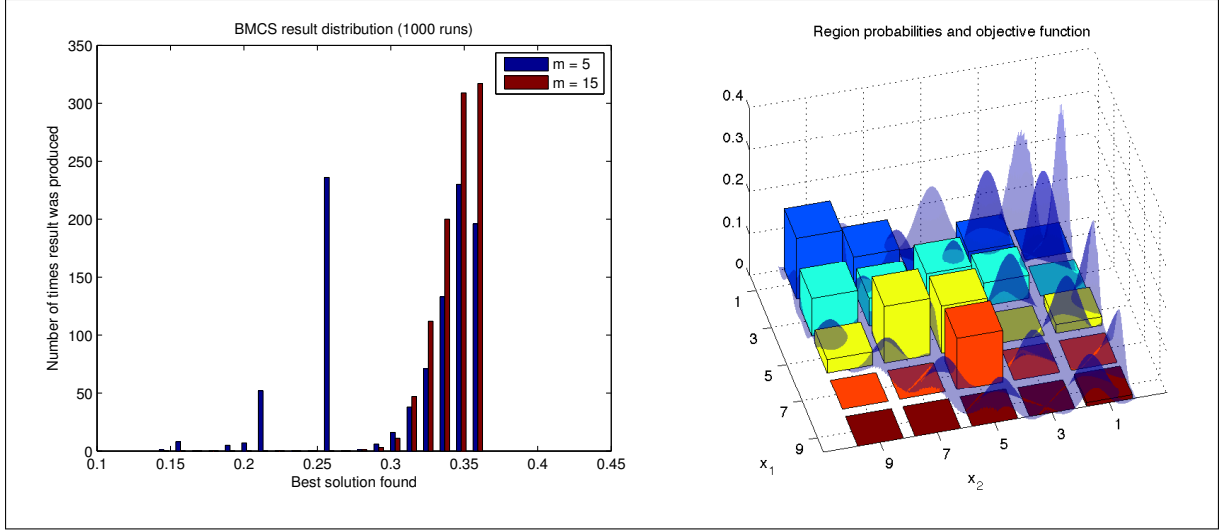


Figure 4: Small regions leading to an unreliable BMCS

- The graph of variation with selection pressure shows that at low pressures the performance varies very little when changing other parameters (the green and cyan dots are all hidden in the densely packed blob).

At larger selection pressures, (pink, red and black) we start to see interesting patterns like the patch on the left, the main curve, curl coming off it and the “jet” of high variance results above the main curve.

- Similarly, results for smaller regions (larger m) are more tightly clustered on the graph. In this case, however, the regions can be made small without sacrificing performance.

For larger regions (especially $m = 5$) the performance is very variable, and most of the strange features of the graphs are composed of results for small m .

- Smaller penalty factors appear to give better performance. However, the combination of:
 - Small penalty factor
 - Small m (and thus large regions)
 - Few initial samples
 - High selection pressure

leads to the very high result variances seen in the “jet”. Figure 4 shows why this is: with large regions and few initial samples, there is a high probability that the initial samples represent the objective function in the region badly—the regions contain both optima and zeroes of the function. The high selection pressure means that regions with low average objective functions are completely ignored, and the search never discovers the optimum in the region.

- Interestingly, when the penalty factor is higher, instead of seeing very high variances with high means, we see the lowest means observed but not such high variances. This is the patch of results on the left. This is because the penalty reduces the chance that sampling will occur near the constraints, and concentrates the search on the local optima near the centre of the space. Fewer, lower optima lead to a smaller, lower range of possible results, and as such smaller μ and σ .

4.2 Effects of Individual Parameters

Figure 6 shows how the performance is affected by difference parameters. These plots use the method described in Section 3.2—for each plot, the named parameter is varied and the others are chosen separately

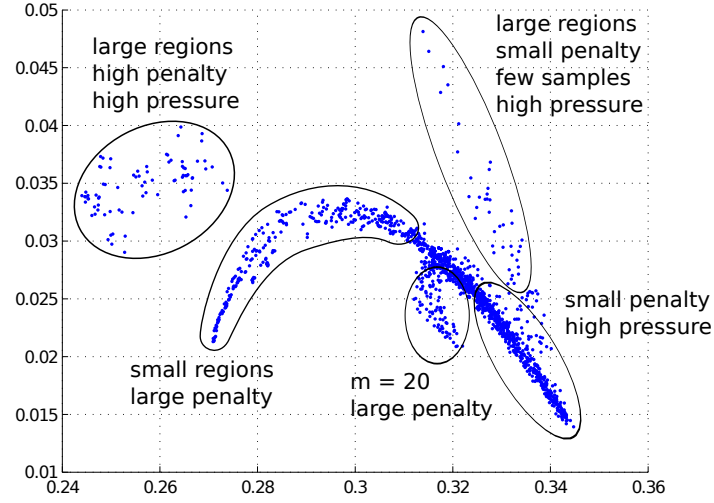


Figure 5: Summary of patterns in BMCS performance

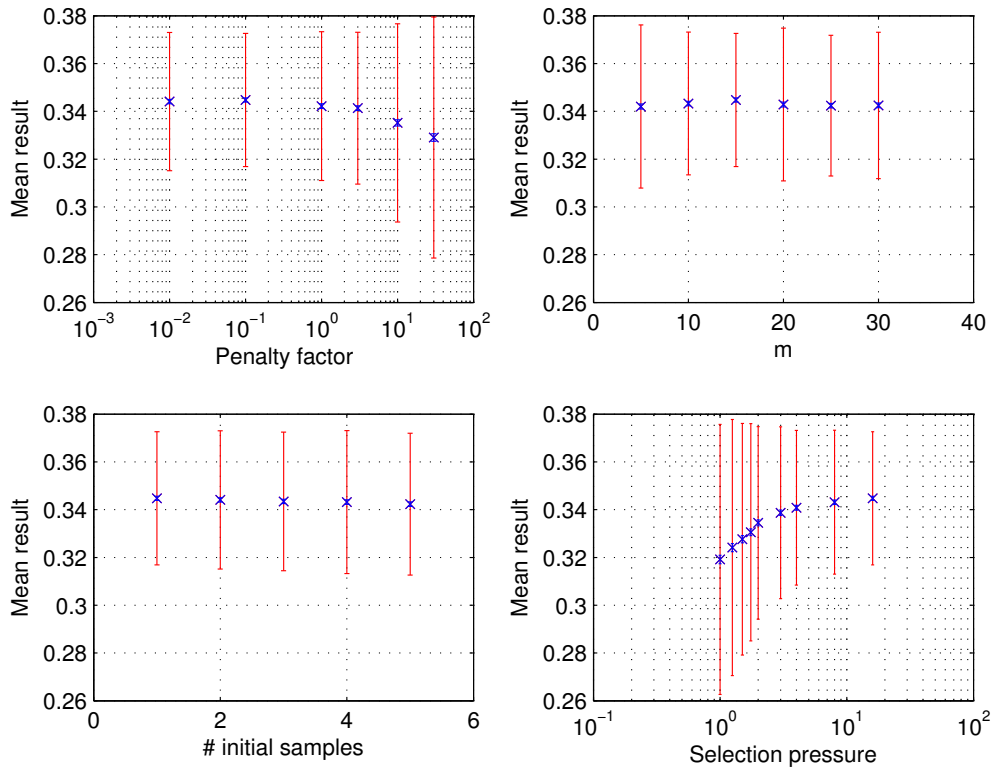


Figure 6: The effects of constraining individual parameters on best performance

for each point to maximise the mean result. The blue points are shown with very small blue error bars that indicate the uncertainty in the estimation of μ . The red error bars stretch 2σ in both directions: they are meant to indicate the variability of the algorithm’s performance. Because the result is not normally distributed, they stretch beyond the maximum possible result.

4.3 Summary of BMCS Parameters

As shown in Figures 5 & 6, the parameters that are critical to good performance on this problem are the penalty weighting factor and the selection pressure. However, appropriate region sizes and initial samples make the algorithm more robust against changes in other parameters.

Penalty factor

If this is much too small, the algorithm may waste time sampling in infeasible space. However, if it is too high and the selection pressure is high, the algorithm may completely fail to sample optima near constraints.

Region size

If this is much too small (large m), the algorithm may waste time in the initial survey. If regions are too large, the algorithm becomes very sensitive to bad settings of other parameters. Figure 3 shows how small regions ($m \geq 20$) guarantee reasonable performance with the tested parameter ranges.

Number of initial samples

This does not have a particularly great effect on the algorithm’s performance. If the other parameters are set well then just one sample per region suffices for good performance, and allows more samples to be directed at promising regions. However, setting this ≥ 4 provides protection against the failure shown in Figure 4.

Selection pressure

This plays an important part in making the algorithm efficient. Figure 6 shows that with $S \leq 3$, the algorithm’s performance is significantly worse, as it doesn’t focus its efforts near the optima. Setting it very high allows the best performance with well-chosen parameters, but causes problems when the regions are too small.

The best performing parameter settings tested:

penalty_factor	m	initial_samples	pressure
0.1	15	1	16

4.4 Dissimilarity Archiving

Since it is harder to quantify the quality of a dissimilarity archive, the above approach focused on the best solution found. Figure 7 shows the contents of the dissimilarity archive after a run of the BMCS algorithm with the parameters in Section 4.3. It also shows the results with a smaller region size, which appears to do a better job of locating the 2nd highest peak (although not much can be drawn from a single run of the algorithm). Both results appear reasonable.

A more full investigation of the problem could develop objective functions to characterise the quality of the dissimilarity archive, such as the sum of the values, and then investigate how this varies with the parameters.

4.5 Computational Cost

In order to run such an intensive evaluation of the algorithm, the implementation needed to be very efficient. 1620 parameter combinations were tested, each requiring 1000 runs of the algorithm to reduce the impact of noise, resulting in over 8 billion function evaluations.

By fully vectorising the MATLAB code for the BMCS implementation, the time required for a single run of the BMCS algorithm was reduced to 5ms, and the evaluation ran in about 2 hours on a notebook computer.

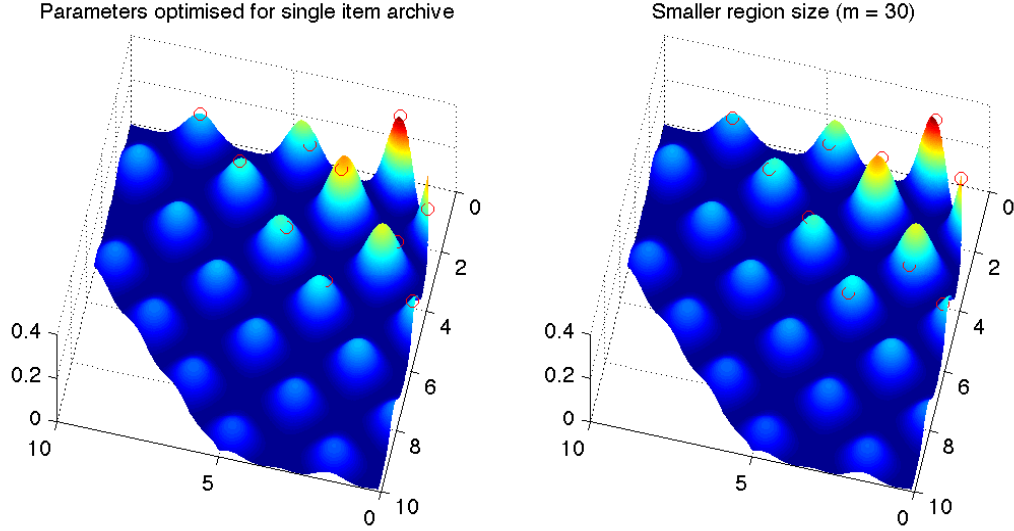


Figure 7: Dissimilarity archives from the BMCS algorithm

5 Simulated Annealing

Since the SA algorithm is more complicated, there were more design decisions made in the implementation, and more parameters could be varied.

Initial Temperature Setting

The algorithm can set the initial temperature one of three ways:

- T_0 given as a parameter to the algorithm.
- $T_0 = \frac{\delta f^-}{\ln 0.8}$ where δf^- is the mean value of objective reductions encountered in an initial survey [Kirkpatrick, 1984].
- $T_0 = \sigma_0$ where σ_0 is the standard deviation of objective changes encountered in the initial survey [White, 1984].

In both adaptive cases the length of the initial survey was chosen as 500 function evaluations.

Solution Generation

The algorithm generates new candidate solutions by adding a random offset to the current position. The random offset can be chosen in one of three ways. Each way takes a step size parameter, S .

- Each dimension change is uniformly selected from the interval $(-S, S)$.
- Each dimension change is selected from a normal distribution with zero mean and variance S^2 .
- Adaptive step sizes [Parks, 1990], as described in the lecture notes. Here, S is the initial maximum step size.

Constraint Handling

As with BMCS, the penalty function described in Section 3.3 is used, with a constant weight parameter, and invalid solutions are not archived. The penalty is also multiplied by T^{-1} as recommended in the lecture notes.

Annealing Schedule

After an optional initial survey of 500 evaluations, the algorithm will perform up to L_k evaluations and up to $\lfloor 0.6L_k \rfloor$ accepted solutions before reducing the temperature. The temperature is then reduced

either by a constant factor [Kirkpatrick et al., 1982] or using a scheme [Huang et al., 1986], as given in the lecture notes.

Restarts

The algorithm was found to frequently get stuck in local optima, especially with adaptive step sizes. To work around this, if the algorithm runs 500 evaluations without improving the best value seen so far, it restarts from the best value, and resets the maximum step size to the initial value.

5.1 Searching the Parameter Space

A similar search to that for BMCS was performed, although considerably more intensive:

Solution generation

Uniform, Gaussian or Adaptive (Parks)

Initial maximum step size

0.5, 0.75, 1.0, 1.5 or 2.0

Penalty weight

0.01 0.03 0.1 0.3 1.0

Initial temperature

Adaptive (Kirkpatrick), Adaptive (White), 0.07, 0.1, 0.2, 0.3, 0.5 or 1.0

Temperature length L_k

50, 100, 200, or 300

Temperature decay

Adaptive (Huang), 0.9, 0.92, 0.95, 0.97, 0.98 and 0.99

This time, there is a total of 16,800 combinations. For each combination, it was run $n = 3000$ times, as described in Section 3.1. Figure 8 shows the results. Note that because so many more datapoints were collected, these plots show random subsets of the data.

Unfortunately, there are no patterns of the same clarity in this format, and so it is presented mainly for completeness. A few things can be seen, though:

- Gaussian solution generation seems to give better performance in more cases than the other techniques. However, top performance can be achieved with all 3 methods.
- Larger maximum step sizes (1.5 or 2) seem to give more consistent results. It is possible that this is because they get stuck in local optima less frequently, as they are able to step between the peaks of the function.
Figure 9 suggests that this intuition is correct, but is of limited use since the two parameter sets used differ greatly.
- Similarly, higher initial temperatures seem to give more consistent results than lower temperatures (or the adaptive methods). Again, this is probably due to lower probability of getting caught in local optima, as suggested by Figure 9.
- Smaller penalty weights once again give the best performance, it is not such an important parameter as for BMCS. This is possibly because the penalty is multiplied by T^{-1} .
- The adaptive temperature reduction seems to perform worse in general than the constant exponential approach. However, it is still capable of achieving good performance.

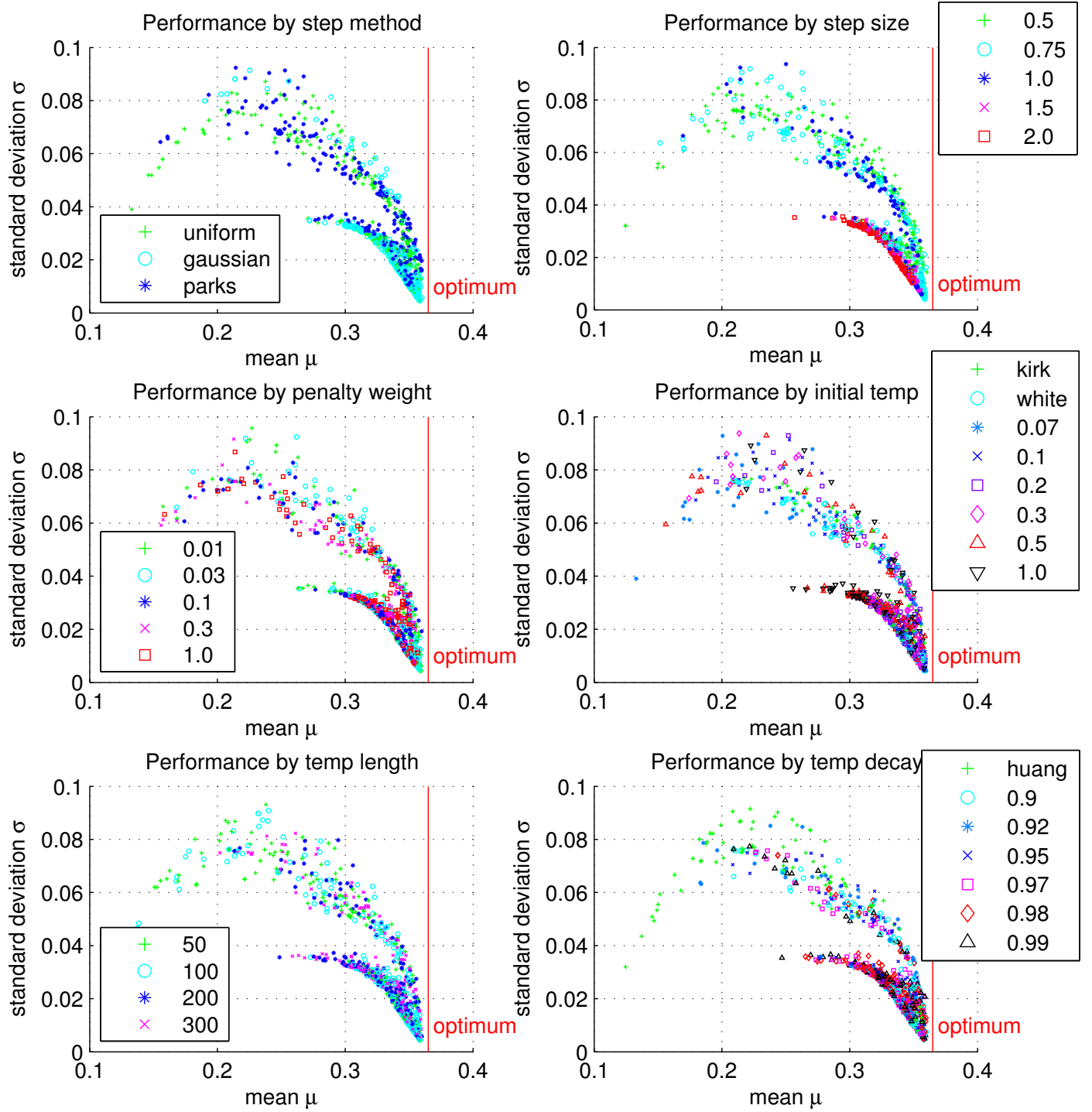


Figure 8: SA performance with varying parameters

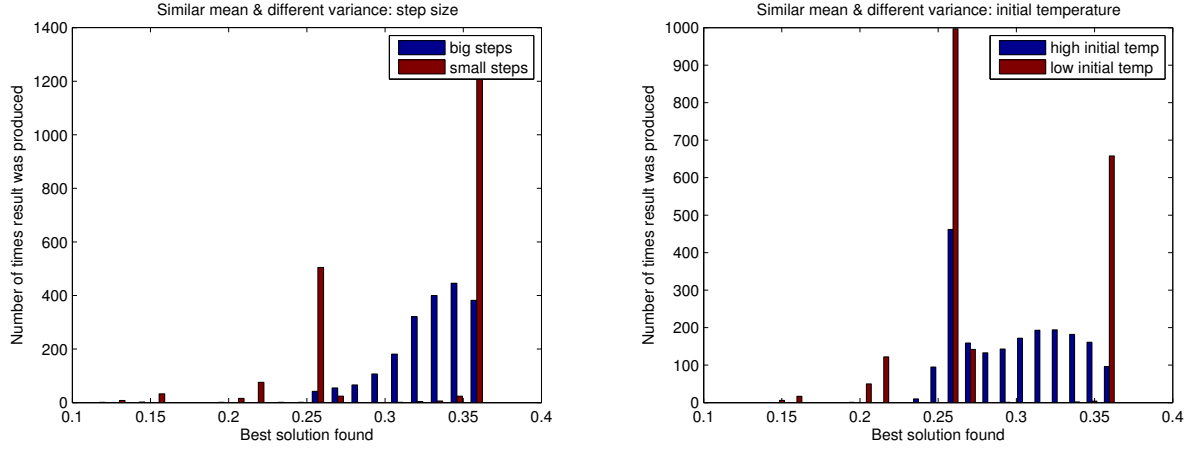


Figure 9: Comparison of result distributions with similar means

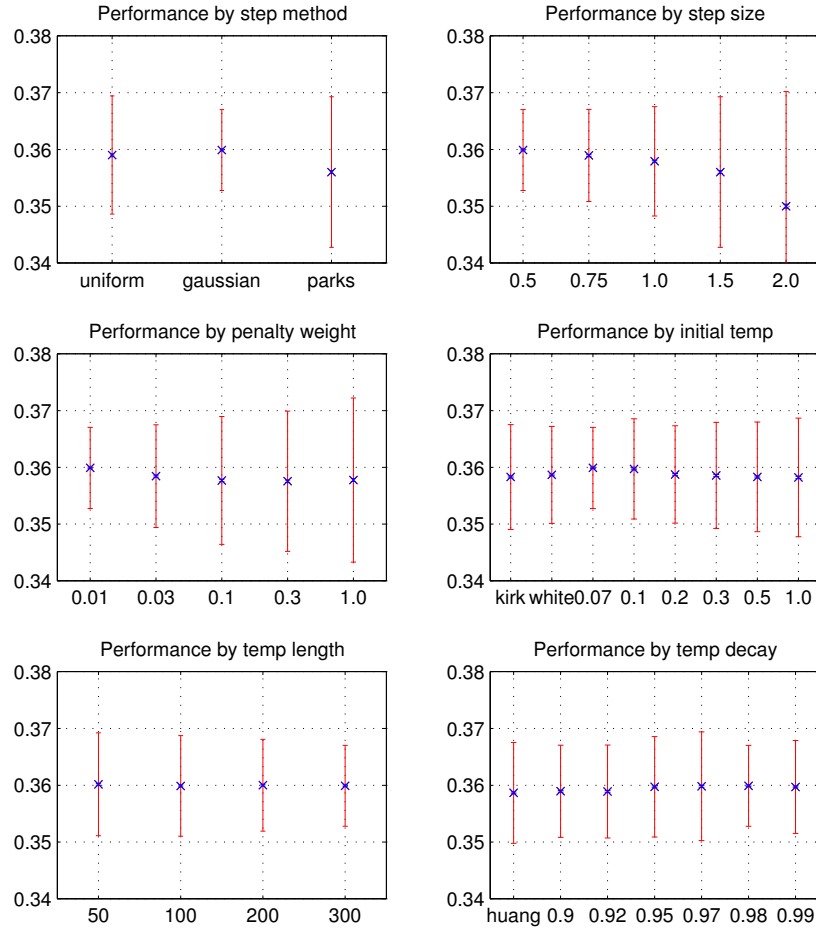


Figure 10: The effects of constraining individual parameters on best performance

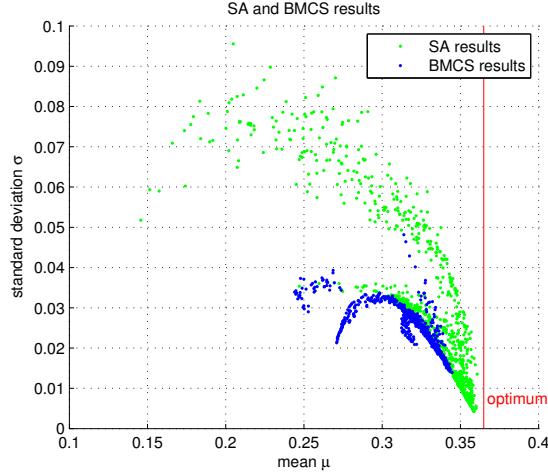


Figure 11: Comparison of SA and BMCS results

5.2 Effects of Individual Parameters

Figure 10 shows the results of the “best performing” parameter sets, as described in Section 3.2—here, the performance metric used to choose which result to display is $\mu - 1.645\sigma$. These plots appear to back up some trends noticed in Figure 8. In particular, Gaussian solution generation and small penalty weightings seem to perform better. However, it also appears that low initial temperatures and small step sizes enable the best performance, which disagrees with what we observed earlier.

It is possible that this approach to the problem is overfitting the parameters: that when the best performing selection out of thousands of different parameter combinations is chosen, an algorithm that regularly finds the optimum in this problem but might perform poorly in a similar problem can be considered best.

5.3 Summary of SA Parameters

Figure 11 makes it clear that SA has the potential to perform better and more reliably than BMCS, if the parameters are chosen well. However, it also has the potential to perform very badly - much worse than the simplest random sampling approach. Unfortunately, it is much harder to see patterns in the performance, or to give general advice on maximising it.

Solution generation

As Figure 12 shows, adaptive solution generation can get closer to the peak more often, because it can reduced the step size and search area near the end of the search. However, it is more prone to only finding the lower local optima, whereas non-adaptive solution generation is less likely to fall into this trap.

Initial maximum step size

A trade-off must be made here: small steps allow pinpointing the optimum, but large steps avoid becoming trapped in a local maximum. This parameter has much less importance when adaptive step sizing is used, but that brings its own problems.

Penalty weight

This parameter is not so important, perhaps since the penalty weight changes as the temperature decreases, but as before lower weights seem to perform slightly better.

Initial temperature

The two methods of adaptive temperature selection appear to perform similarly. Better performance can be achieved by tuning the initial temperature to fit with the other parameters.

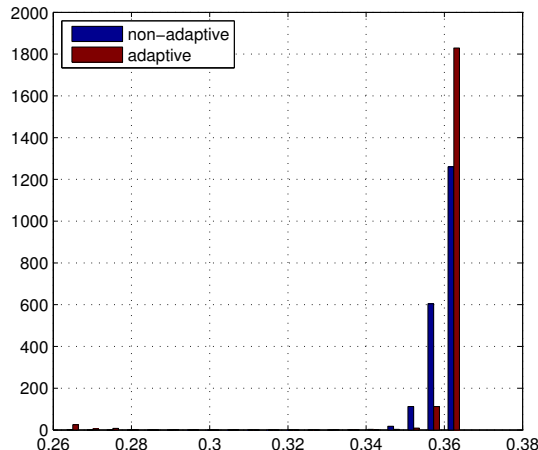


Figure 12: Comparison of adaptive and non-adaptive solution generation

Temperature length L_k

No clear effect is visible here—it appears to be highly dependent on other parameters.

Temperature decay

The adaptive temperature reduction method appears to perform poorly on this problem (unless other parameters are carefully chosen). The exponential cooling scheme [Kirkpatrick et al., 1982] with $\alpha = 0.95$ appears to perform satisfactorily.

5.4 Dissimilarity Archiving

The comments of Section 4.4 apply here too: all that will be done to evaluate the quality of the dissimilarity archive is a visual examination. Figure 13 shows the dissimilarity archive for the best performing SA algorithm, and that for the best performing BMCS algorithm for comparison.

The SA algorithm leaves a worse dissimilarity archive: some local optima are completely ignored, and others are less well maximised. This is not surprising, since it focuses on finding the global optimum, rather than exploring any promising parts of the solution space.

5.5 Computational Cost

The SA algorithm was first implemented in MATLAB. This made it very easy to stop the execution part-way through with the debugger, and to visualise intermediate stages. A debugging tool was written that allows the user to easily scroll through the different stages of the annealing schedule, to give an idea of how the search progresses as the temperature reduces.

However, the performance of the MATLAB implementation is very poor. A single run of the algorithm takes around 0.6 seconds—the examination of the parameter space would have taken almost a year to execute with this implementation. To enable the investigation, the implementation was translated to C++, decreasing the time for a single optimisation to 0.6 milliseconds. With this implementation, the 252 billion function evaluations were conducted in about 10 hours on a notebook computer.

6 Conclusions

The Biased Monte Carlo Sampling algorithm was investigated, and found to perform best with small regions and high selection pressures. However, it was still unable to reliably get results close to the optimum. The

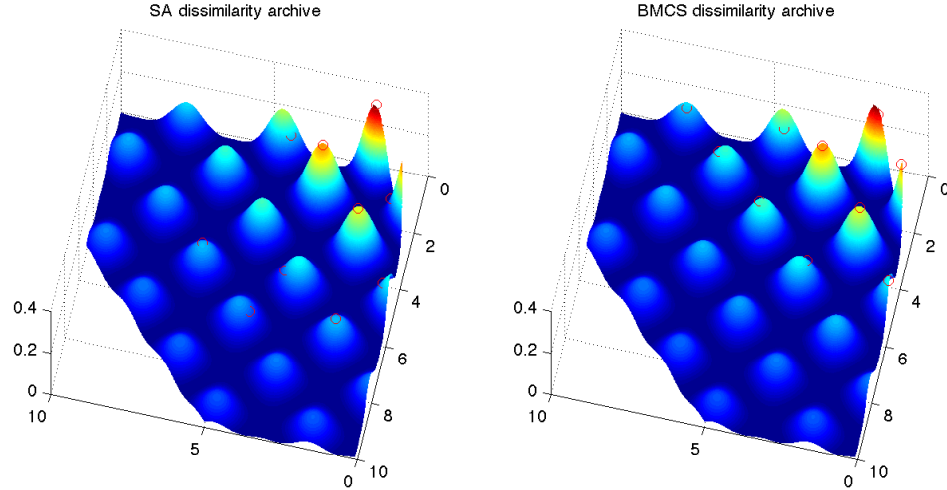


Figure 13: Dissimilarity archives for best performing SA and BMCS algorithms

Simulated Annealing algorithm was found to be able achieve much better performance, and with the right parameters, it could get within a small distance of the global optimum in a large majority of cases. Adaptive step sizing [Parks, 1990] and adaptive temperature initialisation [Kirkpatrick, 1984] [White, 1984] were found to perform well, although not quite as reliably as the best-tuned initialisations.

References

- [Huang et al., 1986] Huang, M. D., Romeo, F., and Sangiovanni-Vincentelli, A. (1986). An efficient general cooling schedule for simulated annealing. In *International Conference on Computer Aided Design*, pages 381–384.
- [Kirkpatrick, 1984] Kirkpatrick, S. (1984). Optimization by simulated annealing: Quantitative studies. *Journal of Statistical Physics*, 34(5-6):975–986.
- [Kirkpatrick et al., 1982] Kirkpatrick, S., C.D. Gelatt, J., and Vecchi, M. (1982). Optimization by simulated annealing. *IBM Research Report*, RC 9355.
- [Parks, 2011] Parks, G. (2011). 5R1 lecture notes.
- [Parks, 1990] Parks, G. T. (1990). An intelligent stochastic optimization routine for nuclear fuel cycle design. *Nucl. Technol.*, (89):233–246.
- [White, 1984] White, S. R. (1984). Concepts of scale in simulated annealing. *AIP Conference Proceedings*, 122(1):261–270.