

# Solving Maximum-Entropy Sampling Problems Using Factored Masks

February 28, 2005; revised October 18, 2005

**Abstract.** We present a practical approach to Anstreicher and Lee’s masked spectral bound for maximum-entropy sampling, and we describe favorable results that we have obtained with a Branch-&-Bound algorithm based on our approach. By representing masks in factored form, we are able to easily satisfy a semidefiniteness constraint. Moreover, this representation allows us to restrict the rank of the mask as a means for attempting to practically incorporate second-order information.

---

## Introduction

Let  $n$  be a natural number, and let  $N := \{1, 2, \dots, n\}$ . Let  $C$  be an order- $n$  symmetric positive definite matrix. Let  $s$  be an integer between 0 and  $n$ . For subsets  $S$  and  $T$  of  $N$ , let  $C[S, T]$  denote the submatrix of  $C$  having rows indexed by  $S$  and columns indexed by  $T$ . The *maximum-entropy sampling problem* is to calculate

$$z(C, s) := \max \{ \ln \det C[S, S] : S \subset N, |S| = s \}.$$

This fundamental problem in the *design of experiments* was introduced in [19] and first used in a monitoring design context in [8] (also see [10, 17, 18, 21]).

From the perspective of algorithmic optimization, the problem has been studied extensively; see [1–3, 11–13, 16] and the surveys [14, 15]. Exact algorithms are based on the Branch-&-Bound framework. Crucial to such an approach is a good upper bound on  $z(C, s)$ . Previous upper bounds are based on either eigenvalues or nonlinear programming relaxations.

In §1, we describe the masked spectral bound of [3]. In §2, we describe two first-order rank-restricted algorithms for quickly obtaining a good mask. In §3, we compare our methods with one another and with the affine-scaling algorithm of [3]. We also examine the issue of how to choose the rank. In §4, we describe our Branch-&-Bound implementation and results. Finally, in §5, we mention possibilities for incorporating second-order information, and in §6, we discuss some future directions.

Notation: For a square symmetric matrix  $A$ ,  $A \succeq 0$  (resp.,  $A \succ 0$ ) means that  $A$  is positive semidefinite (resp., definite). We use  $e$  to denote a vector with

---

Samuel Burer: Department of Management Sciences, University of Iowa, e-mail: [samuel-burer@uiowa.edu](mailto:samuel-burer@uiowa.edu). Supported in part by NSF Grant CCR-0203426.

Jon Lee: Department of Mathematical Sciences, T.J. Watson Research Center, IBM, e-mail: [jonlee@us.ibm.com](mailto:jonlee@us.ibm.com).

all components equal to 1. For a square matrix  $A$ ,  $\text{diag}(A)$  denotes the vector of diagonal entries from  $A$ . For a vector  $x$ ,  $\text{Diag}(x)$  denotes the square matrix having diagonal  $x$  and off-diagonal entries of zero. For a vector  $x$ ,  $x^{-1/2}$  denotes the entrywise reciprocal square roots. For  $n \times n$  matrices  $A$  and  $B$ ,  $A \circ B$  denotes Hadamard (i.e., elementwise) product, while  $A \bullet B := \text{trace}(AB')$ . For a matrix  $A$ ,  $A_i$  denotes the  $i$ -th row.

## 1. The masked spectral bound

Anstreicher and Lee introduced the masked spectral (upper) bound (see [3]) for  $z(C, s)$ . A *mask* is any symmetric, positive semidefinite matrix having ones on the diagonal, i.e., any  $X \succeq 0$  having  $\text{diag}(X) = e$ . We define the associated *masked spectral bound* as

$$\xi_{C,s}(X) := \sum_{l=1}^s \ln(\lambda_l(C \circ X)) ,$$

where  $\circ$  represents the Hadamard product of matrices and eigenvalues  $\lambda_1, \dots, \lambda_n$  are in nonincreasing order. Validity of the masked spectral bound is based on: (i) *Oppenheim's inequality*, i.e.,  $\det A \leq \det A \circ B / \prod_{j=1}^n B_{jj}$ , where  $A \succeq 0$  and  $B \succ 0$ ; and (ii) the eigenvalue inequalities  $\lambda_l(A) \geq \lambda_l(B)$ , where  $A \succeq 0$ , and  $B$  is a principal submatrix of  $A$ .

The masked spectral bound is a generalization of the *spectral partition bound* of [11] (take  $X$  to be block diagonal, with blocks of 1's), which itself is a generalization of both the *eigenvalue bound* of [12] (take  $X = ee'$ ) and the *diagonal bound* of [11] (take  $X = I$ ). The spectral partition bound can produce much better bounds than the eigenvalue and diagonal bounds, but there is no known practical methodology for efficiently choosing a near-optimal partition of  $N$ , which describes the block structure of the mask (see [11]). Some success has been reported using the following method for calculating the so-called “one-big partition” of  $N$  (see [11] and [3]): (i) let  $S$  be a heuristic solution of the maximum-entropy sampling problem; (ii) the associated *one-big partition* of  $N$  has one block  $S$  and the elements of  $N \setminus S$  as singletons.

The motivation for working with the masked spectral bound is to try to use methods of continuous optimization to get the strength of the combinatorial bounds (i.e., the spectral partition bounds and its specializations) but more efficiently. Specifically, we try to find a mask yielding a good bound. Finding an optimal mask is a problem of minimizing a nondifferentiable nonconvex function subject to a semidefiniteness constraint. That sounds daunting, except for the following inter-related points: (i) for our purposes, we do not need a global minimum; (ii) in our experience the nondifferentiability is not so serious; (iii) in our experience, local solutions obtained from different starting points are not wildly different in value; and (iv) by suitably initializing the search for a mask, we can do at least as well as fixed masks.

An important point is that function, gradient and Hessian evaluations for  $\xi$  are expensive. All of these are based on a single (expensive) eigensystem calculation. Function evaluations require eigenvalues, gradients require the eigenvectors

as well, and Hessians further require some non-negligible arithmetic. Considering that we are embedding these calculations in Branch-&-Bound, with the goal for each subproblem of trying to push its upper bound below the value of a global lower bound, we cannot afford to spend much time minimizing. So we need to find a way to descend quickly in just a few steps.

In [3], Anstreicher and Lee proposed a method for minimizing  $\xi$  which was based on ideas coming from the affine scaling algorithm for linear programming. They considered different algorithmic variants, such as short- and long-step approaches, and they demonstrated the ability to achieve good bounds (relative to the one-big, eigenvalue, and diagonal bounds), even in the face of the non-convexity and nondifferentiability of  $\xi$ . They did not, however, consider efficient evaluation of the bound or use of the bound in a complete Branch-&-Bound algorithm.

## 2. Rank-restricted masks

Our idea is to work with rank-restricted masks in factored form. That is, we consider masks of the form  $X := VV'$ , where  $V$  is  $n \times k$  and  $1 \leq k \leq n$ . As we try to find a mask yielding a low bound  $\xi$ , we can hope to make more rapid progress owing to two factors: (i) the semidefiniteness restriction is handled implicitly by the factored form; and (ii) by choosing smaller  $k$ , we work with many fewer variables ( $kn$  rather than  $n(n-1)/2$ ), which may lead to faster convergence of a steepest-descent type method and also opens up second-order methods as a realistic possibility. As it turns out in our experiments, (i) proves to be more beneficial than (ii).

We note that the restriction to low-rank PSD matrices has been used successfully in other contexts [4–6]. In fact, the first algorithm that we describe borrows its key ideas from [4].

Our usage of rank-restricted masks will be based on three functions. The first normalizes the rows of  $V$ , i.e.,

$$h(V) := \text{Diag} \left( [\text{diag}(VV')]^{-1/2} \right) V;$$

the second takes a normalized  $\bar{V}$  and calculates the corresponding masked spectral bound, i.e.,

$$g(\bar{V}) := \xi_{C,s}(\bar{V}\bar{V}');$$

and the third is the composite function  $f(V) := g(h(V))$ . Note that  $h$  is only well-defined if all rows of  $V$  are nonzero.

We propose a steepest-descent type algorithm for minimizing  $f(V)$ , and for this, we consider the gradient of  $f$ . Let  $u_l(C \circ \bar{V}\bar{V}')$  be an eigenvector, of unit Euclidean norm, associated with  $\lambda_l(C \circ \bar{V}\bar{V}')$ . Let  $U$  be the  $n \times s$  matrix having columns  $u_l$  ( $l = 1, \dots, s$ ), and let  $\Sigma$  be the order- $s$  square diagonal matrix with  $\Sigma_{ll} = \lambda_l$  ( $l = 1, \dots, s$ ). Then, as long as  $\lambda_s > \lambda_{s+1}$ , we have that the gradient of  $f$  at  $V$  is the matrix

$$\nabla f(V) = D \nabla g(\bar{V}) - (\nabla g(\bar{V})V' \circ D^3) V;$$

where

$$\begin{aligned} d &= [\text{diag}(VV')]^{-1/2}, \\ D &= \text{Diag}(d), \\ \bar{V} &= h(V), \\ \nabla g(\bar{V}) &= 2(C \circ U \Sigma^{-1} U') \bar{V}. \end{aligned}$$

This can be derived using standard results concerning symmetric functions of eigenvalues (see [20], for example). Note that we must define  $u_l$  properly when  $\lambda_l = \lambda_{l+1}$  for any  $l = 1, \dots, s-1$ ; in such situations, we just take care that the associated  $\{u_l\}$  form an orthonormal basis for each of the eigenspaces corresponding to distinct eigenvalues.

Our steepest-descent algorithm (Algorithm 1) is described in Figure 1. Due to

- 
0. Initialize  $V$ ,  $t_{\max}$ ,  $k_{\max}$  and  $0 < \beta < 1$ . Set  $t = 0$ .
  1. Let  $k := 0$  and  $t := t + 1$ . If  $t > t_{\max}$  then STOP.
  2. Let  $W := V - \beta^k \nabla f(V)$ .
  3. If  $f(WW') < f(VV')$  or  $k = k_{\max}$ , then let  $V := W$  and GOTO 1. Otherwise, let  $k := k+1$  and GOTO 2.
- 

**Fig. 1.** Basic Descent Algorithm

the possible non-differentiability of  $f$ , an important ingredient of the algorithm is the acceptance of a non-improving step if descent is not detected after a certain number ( $k_{\max}$ ) of trials. Values such as  $k_{\max} := 2$  and  $\beta := 1/3$  are practical. (Anstreicher and Lee [3] used a similar back-tracking approach for their affine-scaling algorithm. In particular, they also used parameters  $k_{\max} := 2$  and  $\beta := 1/3$ .) In what follows, we will describe the results of our experiments with various values for  $t_{\max}$  and various choices of the initial  $V$ . For example, in the Branch-&-Bound results of §4, we will take  $t_{\max} \leq 5$ .

Several alternatives to Algorithm 1 can be considered. For example, we can define  $W := h(V - \beta^k \nabla f(V))$  in Step 2, which normalizes the algorithm's iterates and simplifies the gradient formulas because  $d = e$ . Another possibility is to move in the direction  $-\nabla g(\bar{V})$ , i.e., to change Step 2 to read  $W := V - \beta^k \nabla g(\bar{V})$ . This is valid because, in fact,  $-\nabla g(\bar{V})$  is a descent direction for  $f$  at  $V$ :

**Theorem 1.** *Suppose that  $f$  is differentiable at  $V$ , and define  $\bar{V} := h(V)$ . Then  $-\nabla g(\bar{V})$  is a descent direction for  $f$  at  $V$ .*

*Proof.* Define  $M := \nabla g(\bar{V})$ , and let  $\theta_i$  be the angle between  $M_{i\cdot}$  and  $V_{i\cdot}$ . Then

$$\begin{aligned} \nabla f(V) \bullet \nabla g(\bar{V}) &= [DM - (MV' \circ D^3)V] \bullet M \\ &= DM \bullet M - (MV' \circ D^3)V \bullet M \\ &= \sum_{i=1}^n d_i \|M_{i\cdot}\|^2 - \sum_{i=1}^n d_i^3 (V_{i\cdot} M_{i\cdot}')^2 \\ &= \sum_{i=1}^n (d_i \|M_{i\cdot}\|^2 - d_i \|M_{i\cdot}\|^2 \cos^2 \theta_i) \\ &= \sum_{i=1}^n d_i \|M_{i\cdot}\|^2 (1 - \cos^2 \theta_i) \\ &\geq 0. \end{aligned}$$

A final possibility is to combine both of these alternatives to define  $W := h(V - \beta^k \nabla g(\bar{V}))$ . We have experimented with this last alternative, which we refer to as Algorithm 2. It is given in Figure 2. (The notation  $\bar{V}$  and  $\bar{W}$  emphasizes that all iterates have unit-norm rows.)

- 
0. Initialize  $\bar{V}$  with unit-length rows,  $t_{\max}$ ,  $k_{\max}$  and  $0 < \beta < 1$ . Set  $t = 0$ .
  1. Let  $k := 0$  and  $t := t + 1$ . If  $t > t_{\max}$  then STOP.
  2. Let  $\bar{W} := h(\bar{V} - \beta^k \nabla g(\bar{V}))$ .
  3. If  $f(\bar{W}\bar{W}') < f(\bar{V}\bar{V}')$  or  $k = k_{\max}$ , then let  $\bar{V} := \bar{W}$  and GOTO 1. Otherwise, let  $k := k + 1$  and GOTO 2.
- 

**Fig. 2.** Alternative Descent Algorithm

In the next section, we will present evidence to illustrate that the two algorithms above perform quite well, even though they are somewhat simplistic approaches to the minimization of the non-convex, non-differentiable function  $\xi$ . Overall, we believe that the non-differentiability of  $\xi$  is actually not so bad that it precludes the use of simple gradient-based descent methods.

### 3. Experiments with rank-restricted masks

In this section, we detail our experience with Algorithms 1 and 2. A few preliminary remarks are in order. First, the data for our experiments, which has  $n = 63$  and was also used in [3], comes from an environmental monitoring application (see [10]).

Second, we will often compare bounds on the *original* and *complementary* problems. This terminology refers to the fact that, using the identity

$$\ln \det C[S, S] = \ln \det C + \ln \det C^{-1}[N \setminus S, N \setminus S], \quad (1)$$

any bound for the complementary problem of choosing a maximum entropy set of  $n - s$  points with respect to the covariance matrix  $C^{-1}$  translates to a bound

for the original problem (see [1,2]). In practice, the same class of bound (e.g., the one-big bound) can yield very different values when calculated with respect to the original and complementary problems.

Finally, a primary basis for comparison will be the (absolute) gap between a calculated bound and the entropy value of a heuristic solution. This measure, which is invariant under matrix scaling, has been used in all of the computational experiments reported in previous work on bounds.

The heuristic solution used to calculate the gap is obtained via a routine introduced in [12], which consists of two stages: the greedy construction of a candidate  $S$  with  $|S| = s$  and a straightforward pairwise-interchange (“2-opt”) local improvement of  $S$ . More specifically, the greedy scheme is as follows, where  $H(S) := \ln \det C[S, S]$  is the entropy function: initialize  $S = \emptyset$ , and then, for  $j = 1, 2, \dots, s$ , choose  $k \in N \setminus S$  so as to maximize  $H(S \cup \{k\})$  and adjoin  $k$  to  $S$ . Then beginning from the output set  $S$  of the greedy method described above, repeating while possible, choose  $k \in N \setminus S$  and  $l \in S$  so that  $H(S \cup \{k\} \setminus \{l\}) > H(S)$ , and replace  $S$  with  $S \cup \{k\} \setminus \{l\}$ .

In fact, this routine can also be applied to the complementary problem via (1) to obtain a second heuristic solution for the original problem. In [12], it has been observed that this second solution is typically of better quality when  $s > n/2$ , and so we adopt the following scheme: if  $s \leq n/2$ , then we calculate the heuristic solution via the original problem; otherwise, we calculate it via the complementary problem.

This same scheme for calculating a heuristic solution will also be used in §4 to generate an initial global lower bound for use in Branch-&-Bound. Also, we did not attempt to improve this heuristic or construct better alternatives since the main goal of this paper was to evaluate the relative quality of various upper bounds.

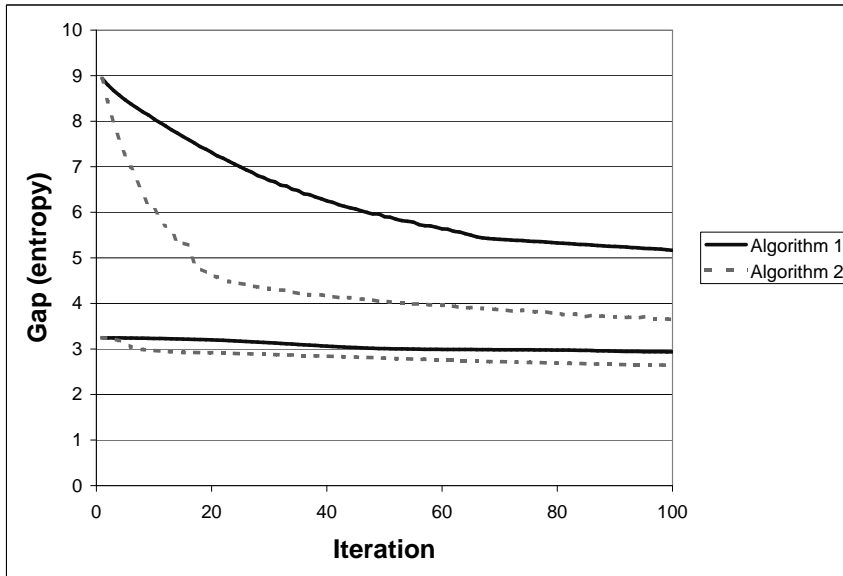
### 3.1. Comparison of Algorithms 1 and 2

We first compare Algorithm 1 with Algorithm 2. Because both algorithms require roughly the same amount of work per iteration, we have found that the running times of the algorithms do not differ significantly. Accordingly, we are primarily interested in comparing convergence properties. (For the interested reader, all runs depicted in Figure 3 below took less than 5 seconds.)

Although Algorithms 1 and 2 can exhibit different convergence rates or patterns on different instances of the maximum-entropy sampling problem with different values of  $k$ , we found that, on any specific instance, Algorithm 2 typically converged more quickly and reliably than Algorithm 1. Figure 3 depicts a typical outcome. We ran both algorithms from the same, random starting point (entries uniformly in  $[-1, 1]$ ) for 100 iterations with  $(n, s, k) = (63, 31, 63)$  on both the original and complementary problems. The curves indicate convergence by depicting entropy gap versus iteration number.

As Figure 3 demonstrates, Algorithm 2 converges more quickly than Algorithm 1 and also achieves better gaps. This is intriguing behavior in light of the

**Fig. 3.** Comparison of the gaps produced by Algorithm 1 and Algorithm 2 when each is run for 100 iterations from the *same, random* starting point on both the *original* and *complementary* problems. Relevant parameters are  $(n, s, k) = (63, 31, 63)$ . The top two curves represent the original problem; the bottom two curves represent the complementary problem.



fact that Algorithm 1 employs the steepest descent direction for  $f$ , while Algorithm 2 uses a kind of projected search path. One hypothesis as to why we see this behavior with Algorithm 2 is as follows: its search path may be perturbing us away from points where  $\xi$  is not differentiable (e.g., Lee and Overton (unpublished) have observed nondifferentiability at true local minimizers obtained by the gradient-sampling algorithm (see [7])).

Another comparison of Algorithms 1 and 2, which supports a similar conclusion as Figure 3, will be given in §3.3.1. Nonetheless, we will advocate the use of Algorithm 1 in §4. Specific reasons for doing so will be explained in that section.

### 3.2. Comparison with affine scaling

We next compare Algorithm 1 with the affine scaling algorithm of Anstreicher and Lee [3], which we refer to as Algorithm AS. In order to obtain as fair a comparison as possible, we obtained their code, which was written in Matlab, and

constructed our own Matlab code for Algorithm 1. All tests for both algorithms were initialized with the same starting point as described in [3] (i.e., a positive definite perturbation of a fixed mask). In particular, since Algorithm AS is an interior-point method, we took  $k = n$  for Algorithm 1 in all comparisons between the two methods.

Our comparisons are based on the following three criteria, each of which has two possible realizations (for an overall total of eight comparisons):

- (a) *initial  $V$*  — a full-rank perturbation of either the eigenvalue or one-big mask;
- (b)  $t_{\max}$  — either 5 or 100;
- (c) *problem type* — either original or complementary problem.

In each of the eight comparisons, we ran the algorithms for all  $s$  between 3 and 60, inclusive.

The purpose of comparing on the basis of  $t_{\max}$  is to observe performance early and late in the algorithm. Early performance (e.g.,  $t_{\max} = 5$ ) gives an indication of how the algorithms would perform in a Branch-&-Bound algorithm, where we can only afford to do a few steps to obtain a bound at each node in the tree. On the other hand, late performance (e.g.,  $t_{\max} = 100$ ) gives an indication of the overall performance and robustness of the algorithms.

Over all eight comparisons, we found similar behavior, namely that Algorithm 1 did at least as well as — and often significantly better than — Algorithm AS. Accordingly, in the interest of space, we provide two representative comparisons of the eight total in Figures 4 and 5.

One difference between Algorithm 1 and Algorithm AS, which we observed but is not immediately evident from Figure 5, is the ability of Algorithm 1 to achieve lower gaps than Algorithm AS — even when Algorithm AS is allowed to run for a large number of iterations (for example, 1000 iterations). Said differently, Algorithm AS seems to stall at higher gaps than Algorithm 1.

Another important criterion for comparison is running time. For both Algorithm AS and Algorithm 1, the work per iteration is dominated by the eigenvalue and eigenvector calculations associated with the function  $\xi$ . Recall also that both algorithms use the same line search back-tracking technique. As such, the running times for the two methods (for a fixed number of iterations) are quite comparable and so are not a source of differentiation. (For example, each individual run represented in Figure 5 took less than 5 seconds.)

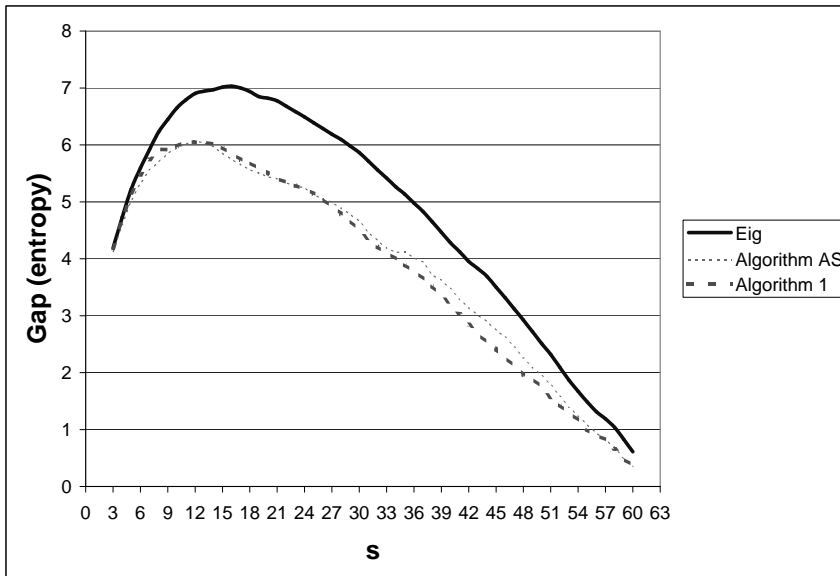
Overall, our experiments lead us to favor Algorithm 1 over Algorithm AS. In particular, we will use Algorithm 1 exclusively for our Branch-&-Bound experiments in §4.

### 3.3. Choosing the rank

An important parametrization issue concerns how the rank  $k$  affects the masked spectral bound. In particular, if we can achieve nearly the same bound by running Algorithm 1 with a smaller  $k$ , we could hope that this would lead to decreased computational requirements. Of course, practically speaking, we need to know



**Fig. 4.** Comparison of the gaps produced by Algorithm 1 and Algorithm AS after 5 steps when initialized near the *eigenvalue* mask on the *complementary* problem, for  $s = 3, \dots, 60$ . The gap given by the eigenvalue mask is also shown for reference.



precisely which rank to choose and how much computational savings can be expected.

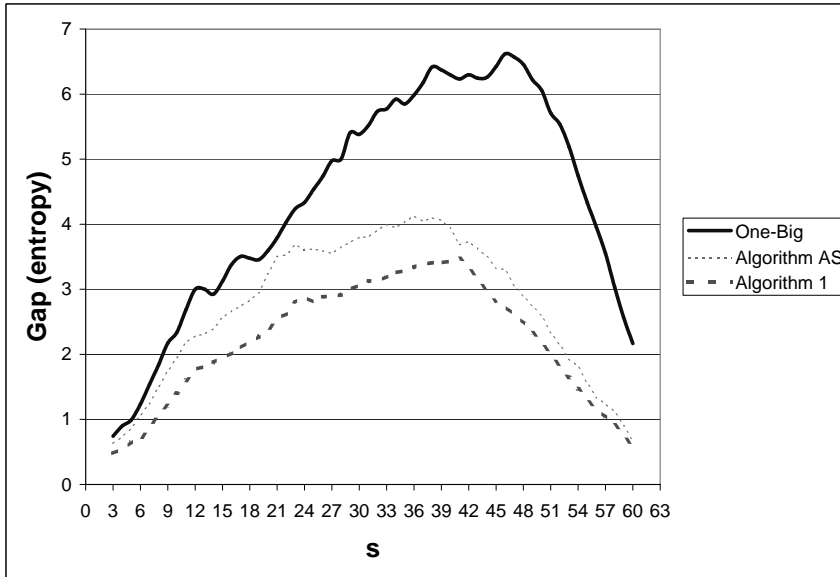
Unfortunately, we have found it difficult to provide any guidelines on which rank to choose. Moreover, the computational savings that we have been able to achieve from reducing the rank are insignificant due to the fact that the eigensystems calculations associated with the function  $\xi$  dominate the running of Algorithm 1. (The eigensystems were calculated using Matlab's `eig` command, which in turn calls standard LAPACK subroutines.) We illustrate these points by means of some example runs.

For each  $k$  between 2 and 63, we ran Algorithm 1 twice with  $s = 31$ , once on the original problem and once on the complementary problem. For both runs, we initialized  $V$  randomly (entries uniformly in  $[-1, 1]$ ) and took  $t_{\max} = 1000$ , that is, we did 1000 iterations. The purpose of so many iterations was to give Algorithm 1 sufficient opportunity to reduce the bound as far as possible. As before, our basis of comparison was the entropy gap. Figure 6 depicts the results.

From Figure 6, it appears that the optimal rank (that is, the smallest rank that yields the best possible bound) for the original problem is approximately  $k = 30$ . On the other hand, for the complementary problem, the optimal rank appears to be around  $k = 5$ . This example illustrates the difficulty in choosing the optimal rank *a priori*.

For the original problem, the average running time over all  $k = 2, \dots, 63$  was 10.8 seconds with a standard deviation of 0.3 seconds; for the complementary

**Fig. 5.** Comparison of the gaps produced by Algorithm 1 and Algorithm AS after 100 steps when initialized near the *one-big* mask on the *original* problem, for  $s = 3, \dots, 60$ . The gap given by the spectral partition mask is also shown for reference.



problem, the numbers were 24.7 and 1.4, respectively. So, the timings showed little variability with respect to the rank  $k$ . As mentioned above, this is due to the dominance of the eigensystem calculations.

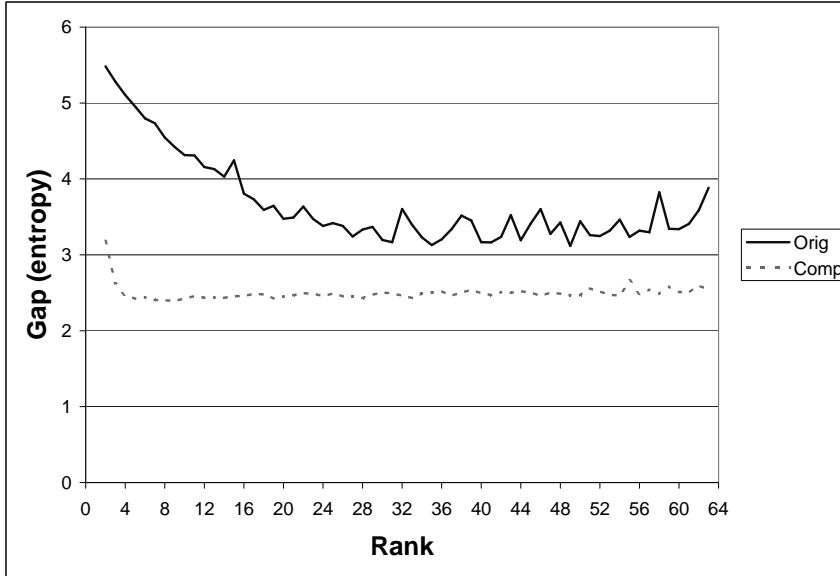
Overall, a conclusion of our experiments is that, when running Algorithm 1, it is reasonable to take  $k = n$  since the resulting running time will not be significantly more than with smaller  $k$ . (In particular, we did not search further for different strategies of choosing  $k$ .) On the other hand, our experiments support the choice of lower  $k$  to get comparable bounds — if a way can be found to exploit the lower rank for better convergence and/or running times.

*3.3.1. Another comparison of Algorithms 1 and 2* The above experiments give us an additional opportunity to compare Algorithm 1 with Algorithm 2. In Figure 7, we replicate the experiments of Figure 6 except this time with Algorithm 2. Notice that, overall, Algorithm 2 converges more reliably than Algorithm 1 for varying ranks. It is also clear that Algorithm 2 achieves slightly better gaps. This evidence provides additional support for the earlier conclusion that Algorithm 2 is more robust than Algorithm 1.

#### 4. Incorporating in Branch-&-Bound

The Branch-&-Bound approach to maximum-entropy sampling was first described in [12]. Branching is done on single elements of  $N$  — either forcing

**Fig. 6.** Gaps produced by Algorithm 1 after 1000 steps when initialized *randomly* on the *original* and *complementary* problems, for ranks  $k = 2, \dots, 63$ .



a single element into the set  $S$  or barring a single element from  $S$ . Thus, at any node in the Branch-&-Bound tree, the corresponding subproblem is determined by forcing a particular set  $F$  of  $f$  elements into  $S$  and barring a particular set  $U$  of  $u$  elements from  $S$ . It remains then to optimize

$$\max \{ \ln \det C[S, S] : S \subset N \setminus U, F \subset S, |S| = s \}.$$

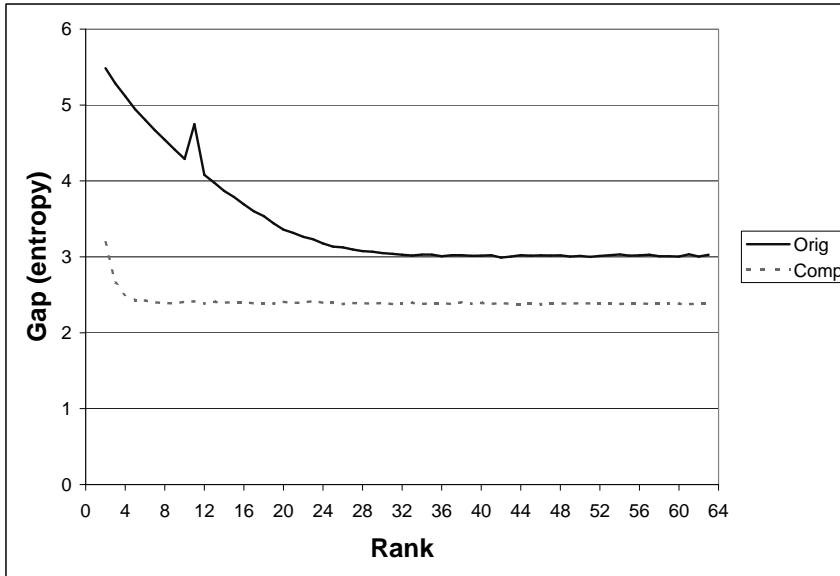
By the Schur complement theorem, this problem is in turn equivalent to choosing a set  $T$  of  $s-f$  elements from the set  $N \setminus (U \cup F)$ , so as to maximize the conditional entropy  $\ln \det C_{F,U}[T, T]$  (plus the constant  $\ln \det C[F, F]$ ), where

$$C_{F,U} := C[N \setminus (F \cup U), N \setminus (F \cup U)] - C[N \setminus (F \cup U), F] (C[F, F])^{-1} C[F, N \setminus (F \cup U)].$$

In other words, the task at each node is to determine  $z(C_{F,U}, s-f) + \ln \det C[F, F]$ , which is itself an instance of the maximum-entropy sampling problem. Hence, any bound developed for the maximum-entropy sampling problem may be used throughout the Branch-&-Bound tree.

We adapted the Branch-&-Bound implementation of [1, 2], which was written in the C programming language and uses LAPACK eigensystem calculations. We kept all default options, including the decision rules for selecting the next node in the tree for branching (the node with the largest upper bound) and for selecting the specific index to branch on (the largest index not already in  $F$  or  $U$ ). One enhancement that we did make, however, was the calculation of an initial global

**Fig. 7.** Gaps produced by Algorithm 2 after 1000 steps when initialized *randomly* on the *original* and *complementary* problems, for ranks  $k = 2, \dots, 63$ .



lower bound via a heuristic solution of the maximum-entropy sampling problem, which has been discussed in §3.

Our goal was to determine whether optimizing the masked spectral bound at each node of the tree is a useful bounding strategy. We compare with the following bounding strategy (for simplicity, *orig* and *comp* refer to bounds coming from the original and complementary problems, respectively):

*Fixed bounding strategy.* For a fixed type of masked spectral bound (i.e., eigenvalue, diagonal, or one-big), compute *orig* and/or *comp* according to the following steps:

1. If  $orig < comp$  for the parent, then calculate *orig*; otherwise, *comp*.
2. If the calculated bound does not fathom, then also calculate the remaining bound.

Note that, for the eigenvalue bound,  $orig = comp$ , so the above steps simplify.

One consequence of the default bounding strategy is that, high in the tree where fathoming is less likely to occur, it is likely that both *orig* and *comp* will be unnecessarily calculated at each node. We found this drawback to be outweighed by the benefit of incorporating both *orig* and *comp*. In particular, the single-minded strategies of just computing either *orig* or *comp* throughout the tree did not work as well.

In developing our strategy for optimizing the masked spectral bound, we felt that calculating bounds for both the original and complementary problems at each node would be too expensive. On the other hand, we knew it would be

highly beneficial to compute the better of the two. After some experimentation, we arrived at a compromise that was based on the following: we often noticed that different subproblems with the same number of elements fixed in and out of  $S$  (i.e., the same  $f$  and  $u$ ) behaved similarly in terms of which bound, *orig* or *comp*, was stronger. (Though we are unable to provide a full explanation for this behavior, we believe it is not surprising that certain bounds behave similarly when, say, a few elements of  $N$  are fixed or when many elements are fixed.) So throughout the Branch-&-Bound calculations, we kept track of all pairs  $(f, u)$ . The first time that a subproblem with a particular  $(f, u)$  was encountered, we calculated both *orig* and *comp* and took note of which was stronger. Afterwards, every time the same  $(f, u)$  was observed, we would calculate *orig* or *comp* according to our first experience.

Our overall strategy for optimizing the masked spectral bound at each node is as follows:

*Optimization bounding strategy.* Determine whether to compute *orig* and/or *comp* (as described above). Then, for each bound to be calculated, do the following:

1. Run Algorithm 1 with  $t_{\max} = 2$ , terminating immediately if fathoming occurs.
2. Using the points generated by Algorithm 1 so far (three points, including the starting point), use quadratic interpolation to judge whether fathoming will occur in the next three iterations.
3. If so, then set  $t_{\max} = 5$  and continue Algorithm 1, terminating immediately if fathoming occurs; otherwise, stop.

We note that the use of quadratic interpolation is a simple, reasonably effective way to judge how quickly the objective function is descending.

A couple of other details concerning the optimization bounding strategy are worth mentioning. First, in accordance with our discussion in §3, we take full rank, i.e.,  $k = s - f$ , in Algorithm 1 at each node of the tree. Second, instead of initializing  $V$  at each node randomly, we attempted to simulate the warm starting of a child node by its parent. We did this by keeping a single, shared-memory copy of  $V$ , which was initialized randomly and made available to all nodes in the tree. A subproblem of size  $s - f$  corresponding to the matrix  $C_{F,U}$  (defined above) was then initialized from the submatrix  $V[N \setminus (U \cup F), N \setminus (U \cup F)]$ , and its final iterate was stored back into the same entries of  $V$ . Although heuristic, we found this warm-start approach to be much better than our original strategy of initializing each node randomly. (The root node was initialized with  $V$  having entries uniform in  $[-1, 1]$ .)

Even though the above strategy uses Algorithm 1, it could also be based on Algorithm 2. In Branch-&-Bound, however, we advocate the use of Algorithm 1, even though §3 suggests the superiority of Algorithm 2. Our decision is based on experiments, which showed that the performance of Algorithm 2 in Branch-&-Bound was worse in almost all cases. (One exception having  $(n, s) = (63, 31)$  is given in the next paragraph.) Although the underlying reasons for this behavior are not entirely clear to us, we did notice certain characteristics of the Branch-&-

Bound runs, which indicated that Algorithm 2 was doing more work on average than Algorithm 1 (apparently without significant improvement in bounds). For example, Algorithm 2 required about 2.5 trial stepsizes per linesearch on average, while Algorithm 1 required about 1.05.

During early experiments with Branch-&-Bound, it became clear to us that none of the bounding strategies that we have proposed would be able to solve the hardest instances (e.g.,  $s = 31$  for the data with  $n = 63$ ) in a reasonably short period of time (say, a few hours). In this sense, we cannot claim that optimizing the masked spectral bound is the “silver bullet” for solving extremely large and difficult instances. (It is worth mentioning that we were able to find and prove optimality for the  $(n, s) = (63, 31)$  problem in 84 hours using Algorithm 1—and 65 hours using Algorithm 2—on a Pentium 4 2.4 GHz running Linux. This problem is out of reach using a fixed mask, and in fact, this is the first time that this instance has been solved using a strategy based completely on eigenvalue-based bounds.)

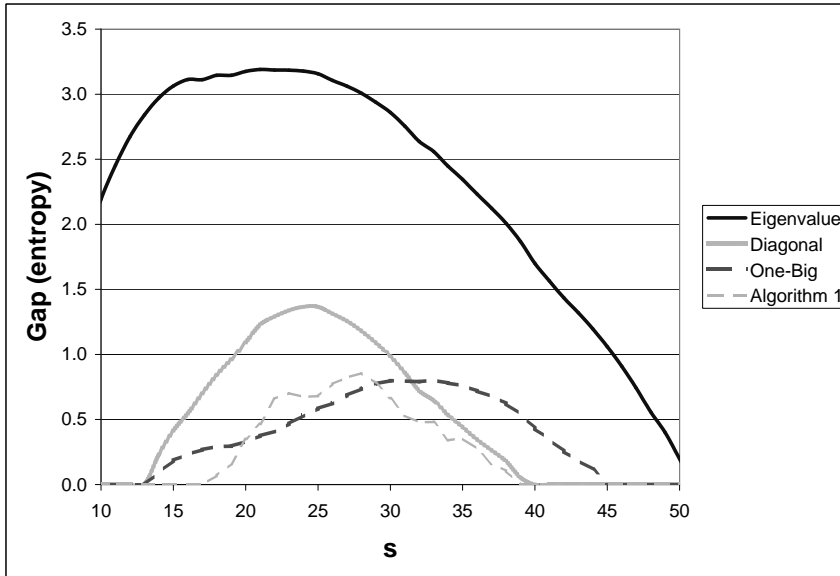
Nevertheless, it was also clear that the various bounding strategies behaved very differently. In order to highlight these differences while maintaining a reasonable testing environment, we settled upon the following testing scheme. For a particular instance of the maximum-entropy sampling problem, we ran each of the four bounding strategies — the three fixed bounding strategies and the optimized bounding strategy — for at most 3600 seconds. We then calculate the gap between the global upper bound (i.e., the largest upper bound of all nodes remaining in the tree) and the global lower bound (i.e., the entropy of the best feasible solution found so far).

Since we limit the Branch-&-Bound calculations to 3600 seconds, we stress that the computational results shown in Figures 8 and 9 below should be interpreted simply as an indication that our approach is more effective in reducing the gaps when compared with other bounding strategies based on masks. This is not to say that it is the most effective among all bounding strategies or that all problems considered here will be solved to optimality with modest increases in time. For example, a completely different class of bounds based on nonlinear programming techniques have been shown in [2] to be quite effective on the  $n = 63$  instances, for example solving the  $s = 31$  instance in 3300 seconds (when the computer clock speed of 125 MHz in [2] is normalized to our clock speed of 2.4 GHz). For the  $n = 124$  instances shown in Figure 9, it should also be mentioned that no known bounding technique has been able to solve the hardest of these instances.

Figure 8 gives Branch-&-Bound results for the data with  $n = 63$  introduced in §3. The results for each of the four bounding strategies are graphed for  $s = 10, \dots, 50$ . From this figure, we can see that the optimization bound strategy did uniformly better than the fixed bound strategies based on the eigenvalue and diagonal masks. Furthermore, compared to the one-big spectral partition bound, our method performed better in the ranges  $s = 14, \dots, 19$  and  $s = 30, \dots, 44$ .

In Figure 9, we give similar results for a data set with  $n = 124$ , which was first used in [16]. Again, we see that optimizing the masked spectral bound is quite competitive with the other bounding strategies.

**Fig. 8.** Branch-&-Bound results with  $n = 63$  after 3600 seconds for the four bounding strategies. Gap refers to the difference between global upper and lower bounds at termination.



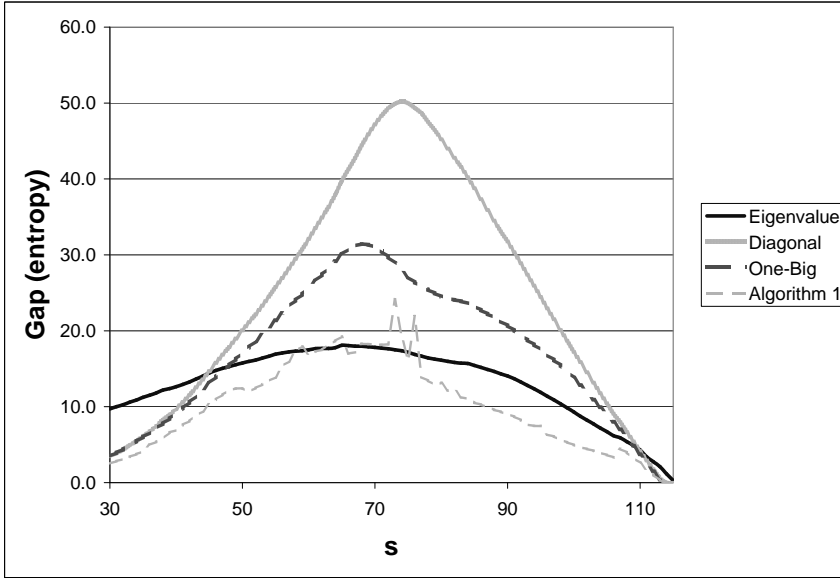
## 5. Exploiting second-order information

As mentioned at the beginning of §2, one of our initial motivations for considering rank-restricted masks was the opportunity for reducing the dimension of the problem so that second-derivative knowledge of  $\xi$  could be incorporated efficiently.

In particular, we were interested in developing a variant of Algorithm 1 with search directions produced by a modified Newton's method with exact Hessians. We did successfully implement just such a method and tested it with small rank, e.g.,  $k \approx 5$ , but unfortunately, it did not outperform the steepest-descent version of Algorithm 1. Downsides to the method included the time required to form the Hessian and to factor the modified Hessian as well as the reduced bound quality caused by taking  $k$  small. In addition, the strong nonconvexity displayed by  $\xi$  resulted in search directions that were not much better than steepest descent.

We also tried other approaches for incorporating second-order information, including BFGS and limited-memory BFGS search directions. Here again, however, our attempts at improving steepest descent were unsuccessful. An added complication was a more stringent strong-Wolfe linesearch (see [9], for example) required by BFGS-type directions. (Since the first version of this paper appeared, recent results obtained by Overton (private communication) suggest that the strong-Wolfe linesearch should not be used when applied within BFGS for the optimization of nondifferentiable functions.)

**Fig. 9.** Branch-&-Bound results with  $n = 124$  after 3600 seconds for the four bounding strategies. Gap refers to the difference between global upper and lower bounds at termination.



We include our experiences here in order to provide a complete picture of our efforts and a starting point for additional research. We also include below the exact Hessian formulae for  $\xi$  since these have not appeared before in the literature and since we are hopeful that they will be of use to others. The formulae were developed from [20].

In the results below, all eigenvectors are taken to have unit length. Note that, as in the development of the gradient of  $f(V)$ , we must define  $u_l$  properly when  $\lambda_l = \lambda_{l+1}$ .

We first examine the Hessian of  $\xi_{C,s}(X)$ . For simplicity, and with an eye toward implementation, we consider the symmetric variable  $X$  to be encoded as a column vector of dimension  $n(n+1)/2$ , where the lower-triangular part of  $X$  is stored in column-major format. In the theorem below, the vectors  $\sigma_{kl}$  are indexed by the entries of  $X$  in the same fashion.

**Proposition 1.** *Suppose  $X \succeq 0$ , and let  $\{(\lambda_k, u_k)\}_{k=1}^n$  be the eigenvalues and eigenvectors of  $C \circ X$ . Suppose also that  $\lambda_s(C \circ X) > \lambda_{s+1}(C \circ X)$ . Then  $\xi$  is analytic at  $X$  with Hessian*

$$\nabla^2 \xi_{C,s}(X) = \sum_{k=1}^s \sum_{l=1}^s \left( -\frac{1}{\lambda_k \lambda_l} \right) \sigma_{kl} \sigma'_{kl} + 2 \sum_{k=1}^s \sum_{l=s+1}^n \left( \frac{1}{\lambda_k} \frac{1}{\lambda_k - \lambda_l} \right) \sigma_{kl} \sigma'_{kl},$$



where, for each  $(k, l)$ ,

$$[\sigma_{kl}]_{ij} = \begin{cases} [C \circ u_l u'_k]_{ij} + [C \circ u_l u'_k]_{ji} & \text{if } i > j \\ [C \circ u_l u'_k]_{ij} & \text{if } i = j. \end{cases}$$

We next examine the Hessian of  $f(V)$ . It is helpful here as well to consider  $V$  to be encoded as a column vector in column-major format.

**Proposition 2.** *Suppose each row of  $V$  is nonzero, and define*

$$\begin{aligned} d &:= [\text{diag}(VV')]^{-1/2}, \\ \bar{V} &:= h(V), \\ X &:= \bar{V}\bar{V}'. \end{aligned}$$

Let  $\{(\lambda_k, u_k)\}_{k=1}^n$  be the eigenvalues and eigenvectors of  $C \circ X$ , and define

$$H := \sum_{k=1}^s \lambda_k^{-1} u_k u'_k.$$

Suppose also that  $\lambda_s(C \circ X) > \lambda_{s+1}(C \circ X)$ . Then  $f$  is analytic at  $V$  with Hessian

$$\nabla^2 f(V) = \sum_{k=1}^s \sum_{l=1}^s \left( -\frac{1}{\lambda_k \lambda_l} \right) \sigma_{kl} \sigma'_{kl} + 2 \sum_{k=1}^s \sum_{l=s+1}^n \left( \frac{1}{\lambda_k} \frac{1}{\lambda_k - \lambda_l} \right) \sigma_{kl} \sigma'_{kl} + \mathcal{M} + \mathcal{N} + \mathcal{P},$$

where, for all  $(k, l)$  [and defining  $J := C \circ (u_l u'_k + u_k u'_l)$  and  $\bar{x} := \bar{V}_p$  locally],

$$[\sigma_{kl}]_{ip} = [d \circ (J\bar{x} - \bar{x} \circ (J \circ X) e)]_i$$

and where, for all  $(ip, jq)$  [and defining  $\bar{x} := \bar{V}_p$  and  $\bar{y} := \bar{V}_q$  locally],

$$\mathcal{M}_{(ip)(jq)} = 2 [dd' \circ ((e'_p e_q) e e' - (\bar{x} \circ \bar{y}) e' - e (\bar{x} \circ \bar{y})' + \bar{x} \bar{y}' \circ X) \circ H]_{ij},$$

$$\mathcal{N}_{(ip)(jq)} = \begin{cases} 2 [d^2 \circ (3\bar{x} \circ \bar{y} \circ (X \circ H) e - \bar{x} \circ H\bar{y} - \bar{y} \circ H\bar{x})]_i & \text{if } i = j \\ 0 & \text{otherwise,} \end{cases}$$

and

$$\mathcal{P}_{(ip)(jq)} = \begin{cases} -2 [d^2 \circ (X \circ H) e]_i & \text{if } (i, j) = (p, q) \\ 0 & \text{otherwise.} \end{cases}$$

## 6. Further directions

One possible extension would be to combine our low-rank approach with gradient-sampling ideas (see [7]) when nondifferentiability becomes an issue at a mask that nearly fathoms, though it may well be that the time would be just as well spent on further children.

We feel that exploiting second-order information, particularly when working with low-rank masks, is still an interesting avenue of research. On some smaller instances (e.g.,  $n \approx 50$ ) of Branch-&-Bound with Algorithm 1, we noticed that taking low rank produced overall faster running times than taking full rank. At this time, we are unable to reconcile this behavior with our observations concerning rank in §3, but it would be interesting if one could determine effective strategies for using small rank — and then effectively incorporate second-order information as well.

*Acknowledgements.* The authors are grateful to Kurt Anstreicher for introducing them to one another and for providing many insightful comments on the paper. Michael Overton provided invaluable input as well.

## References

1. Kurt M. Anstreicher, Marcia Fampa, Jon Lee, and Joy Williams, *Continuous relaxations for constrained maximum-entropy sampling*, Integer programming and combinatorial optimization (Vancouver, BC, 1996), Lecture Notes in Computer Science, vol. 1084, Springer, Berlin, 1996, pp. 234–248.
2. ———, *Using continuous nonlinear relaxations to solve constrained maximum-entropy sampling problems*, Mathematical Programming, Series A **85** (1999), no. 2, 221–240.
3. Kurt M. Anstreicher and Jon Lee, *A masked spectral bound for maximum-entropy sampling*, mODA 7—Advances in Model-Oriented Design and Analysis, Contributions to Statistics, Physica-Verlag, Heidelberg, 2004, pp. 1–10.
4. Samuel Burer and Renato D. C. Monteiro, *A projected gradient algorithm for solving the Maxcut SDP relaxation*, Optimization Methods and Software **15** (2001), 175–200.
5. ———, *A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*, Mathematical Programming, Series B **95** (2003), no. 2, 329–357.
6. Samuel Burer, Renato D. C. Monteiro, and Yin Zhang, *Rank-two relaxation heuristics for max-cut and other binary quadratic programs*, SIAM Journal on Optimization **12** (2001), no. 2, 503–521.
7. James V. Burke, Adrian S. Lewis, and Michael L. Overton, *A robust gradient sampling algorithm for nonsmooth, nonconvex optimization*, SIAM Journal on Optimization (to appear) (2003).
8. William F. Caselton and James V. Zidek, *Optimal monitoring networks*, Statistics and Probability Letters **2** (1984), 129–178.
9. R. Fletcher, *Practical methods of optimization*, second ed., John Wiley & Sons, New York, 1987.
10. Peter Guttorp, Nhu D. Le, Paul D. Sampson, and James V. Zidek, *Using entropy in the redesign of an environmental monitoring network*, Tech. Report 116, Department of Statistics, University of British Columbia, 1992.
11. Alan Hoffman, Jon Lee, and Joy Williams, *New upper bounds for maximum-entropy sampling*, MODA 6 — Advances in model-oriented design and analysis (Anthony C. Atkinson, Peter Hackl, and Werner G. Müller, eds.), Springer, 2001, pp. 143–153.
12. Chun-Wa Ko, Jon Lee, and Maurice Queyranne, *An exact algorithm for maximum entropy sampling*, Operations Research **43** (1995), no. 4, 684–691.
13. Jon Lee, *Constrained maximum-entropy sampling*, Operations Research **46** (1998), no. 5, 655–664.

14. ———, *Semidefinite programming in experimental design*, Handbook of Semidefinite Programming (Romesh Saigal Henry Wolkowicz and Lieven Vandenberghe, eds.), International Series in Operations Research and Management Science, vol. 27, Kluwer, Boston, 2000, pp. 528–532.
15. ———, *Maximum-entropy sampling*, Encyclopedia of Environmetrics (Abdel H. El-Shaarawi and Walter W. Piegorsch, eds.), vol. 3, Wiley, 2001, pp. 1229–1234.
16. Jon Lee and Joy Williams, *A linear integer programming bound for maximum-entropy sampling*, Mathematical Programming, Series B **94** (2003), no. 2-3, 247–256.
17. Werner G. Müller, *Collecting spatial data*, revised ed., Contributions to Statistics, Physica-Verlag, Heidelberg, 2001, Optimum design of experiments for random fields.
18. Paola Sebastiani and Henry P. Wynn, *Maximum entropy sampling and optimal Bayesian experimental design*, Journal of the Royal Statistical Society, Series B (Statistical Methodology) **62** (2000), no. 1, 145–157.
19. Michael C. Shewry and Henry P. Wynn, *Maximum entropy sampling*, Journal of Applied Statistics **46** (1987), 165–170.
20. Nam-Kiu Tsing, Michael K. H. Fan, and Erik I. Verriest, *On analyticity of functions involving eigenvalues*, Linear Algebra and its Application **207** (1994), 159–180.
21. Shiyong Wu and James V. Zidek, *An entropy based review of selected NADP/NTN network sites for 1983-86*, Atmospheric Environment **26A** (1992), 2089–2103.