

# Bucket and mini-bucket Schemes for M Best Solutions over Graphical Models

Natalia Flerova, Emma Rollon and Rina Dechter

## Abstract

The paper focuses on finding the  $m$  best solutions of a combinatorial optimization problem defined over a graphical model (e.g., the  $m$  most probable explanations for a Bayesian network). We describe *elim-m-opt*, a new bucket elimination algorithm for solving the  $m$ -best task, provide efficient implementation of its defining combination and marginalization operators, analyze its worst-case performance, and compare it with that of recent related algorithms. An extension to the mini-bucket framework, yielding a collection of bounds for each of the  $m$ -best solutions is discussed and empirically evaluated. We also formulate the  $m$ -best task as a regular reasoning task over general graphical models defined axiomatically, which makes all other inference algorithms applicable.

## 1 Introduction

Given an optimization problem, the objective typically is to find an optimal solution, i.e., a solution that provides the best value of the objective function. However, in many applications it is desirable to obtain not just a single optimal solution but a set (of a given size  $m$ ) of the best possible solutions. Such a set can be useful, for example, in assessing the sensitivity of the optimal solution to variation of the parameters of the problem, or when a set of diverse assignments with approximately the same cost is wanted.

Lawler [Lawler, 1972] provided a general scheme for using any optimization algorithm to solve the  $m$ -best task. Its main idea is to compute the  $m$ -best solutions by successively computing the best solution, each time using a slightly different reformulation of the original problem. This approach has been extended and improved over the years and is still one of the primary strategies to date for finding the  $m$ -best solutions. The approach used in this paper is to develop direct algorithms that avoid the repeated computation inherent in Lawler's scheme. The main idea is to integrate the  $m$ -best task into existing optimization schemes such as inference or search. In the paper we focus on graphical models and show how the well-known Bucket Elimination (BE) framework can be extended to compute the  $m$ -best solutions by a relatively simple modification of its underlying combination

and marginalization operators [Dechter, 1999] yielding algorithm *elim-m-opt*.

We analyze the complexity of *elim-m-opt*, compare it to previously developed schemes, and discuss extensions to Mini-Bucket Elimination to compute bounds on each of the  $m$ -best solutions, yielding algorithm *mbe-m-opt*. We then show that the  $m$ -best task can be formalized more broadly as a reasoning task over a general graphical model as defined axiomatically by [Shenoy and Shafer, 1990]. As a consequence, any inference algorithm for solving a reasoning task over this representation (including Bucket Elimination) is immediately applicable as a sound and complete algorithm for the  $m$ -best task. In particular, *elim-m-opt* is therefore applicable to a wide range of graphical models.

We also provide empirical analysis for *mbe-m-opt* demonstrating its effectiveness both as an exact scheme as well as for approximation.

## 2 Background

We consider problems expressed as graphical models. The graphical model framework provides a common formalism to model a broad spectrum of problems, and a collection of general algorithms to efficiently solve them. Examples of graphical models are Markov and Bayesian networks [Pearl, 1988], constraint networks and influence diagrams. We next follow definitions as in [Kask *et al.*, 2005].

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be an ordered set of variables and  $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_n)$  an ordered set of domains. Domain  $\mathbf{D}_i$  is a finite set of potential values for  $X_i$ . The assignment of variable  $X_i$  with  $a \in \mathbf{D}_i$  is noted  $(X_i = a)$ . A tuple is an ordered set of assignments to different variables  $(X_{i_1} = a_{i_1}, \dots, X_{i_k} = a_{i_k})$ . A *complete assignment* to all the variables in  $\mathbf{X}$  is called a *solution*. Let  $t$  and  $s$  be two tuples having the same instantiations to the common variables. Their join, noted  $t \cdot s$ , is a new tuple which contains the assignments of both  $t$  and  $s$  (this notation is also used for multiplication, so we assume the meaning will be clear from the context). If  $t$  is a tuple over a set  $T \subseteq X$  and  $\mathbf{S}$  is a set of variables, then  $t_{[S]}$  is a relational projection of  $t$  on  $\mathbf{S}$ .

A valuation set  $\mathbf{A}$  admits two binary operations over valuations:  $\otimes : \mathbf{A} \times \mathbf{A} \rightarrow \mathbf{A}$  called combination and  $\oplus : \mathbf{A} \times \mathbf{A} \rightarrow \mathbf{A}$  called addition. Both operators are associative and commutative. Typical combination operators are sum and product over numbers, and logical AND (i.e.,  $\wedge$ ) over booleans. Typical

addition operators are min, max and sum over numbers and logical OR (i.e.,  $\vee$ ) over booleans.

We denote by  $\mathbf{D}_Y$  the set of tuples over a subset of variables  $Y$ , also called the *domain of Y*. Functions are defined on subsets of variables of  $\mathbf{X}$ , called scopes and their range is a set  $\mathbf{A}$  whose elements are called *valuations*. If  $f : \mathbf{D}_Y \rightarrow \mathbf{A}$  is a function the scope of  $f$ , denoted  $\text{var}(f)$ , is  $Y$ . In the following, we will use  $\mathbf{D}_f$  as a shorthand for  $\mathbf{D}_{\text{var}(f)}$ .

**DEFINITION 1** (combination operator). *Let  $f : \mathbf{D}_f \rightarrow \mathbf{A}$  and  $g : \mathbf{D}_g \rightarrow \mathbf{A}$  be two functions. Their combination, noted  $f \otimes g$  is a new function with scope  $\text{var}(f) \cup \text{var}(g)$ , s.t.*

$$\forall t \in \mathbf{D}_{\text{var}(f) \cup \text{var}(g)}, (f \otimes g)(t) = f(t) \otimes g(t)$$

**DEFINITION 2** (marginalization operator). *Let  $f : \mathbf{D}_f \rightarrow \mathbf{A}$  be a function and  $\mathbf{W} \subseteq \mathbf{X}$  be a set of variables. The marginalization of  $f$  over  $\mathbf{W}$ , noted  $\Downarrow_{\mathbf{W}} f$ , is a function whose scope is  $\text{var}(f) - \mathbf{W}$ , s.t.*

$$\forall t \in \mathbf{D}_{\text{var}(f) - \mathbf{W}}, (\Downarrow_{\mathbf{W}} f)(t) = \bigoplus_{t' \in \mathbf{D}_{\mathbf{W}}} (t \cdot t')$$

**DEFINITION 3** (graphical model). A **graphical model** is a tuple  $\mathcal{M} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes)$ , where:  $\mathbf{X} = \{X_1, \dots, X_n\}$  is a set of variables;  $\mathbf{D} = \{D_1, \dots, D_n\}$  is the set of their finite domains of values;  $\mathbf{A}$  is a set of valuations ( $\mathbf{A}, \otimes, \oplus$ );  $\mathbf{F} = \{f_1, \dots, f_r\}$  is a set of discrete functions, where  $\text{var}(f_j) \subseteq \mathbf{X}$  and  $f_j : \mathbf{D}_{f_j} \rightarrow \mathbf{A}$ ; and  $\otimes$  is the combination operator over functions (see Definition 1). The graphical model  $\mathcal{M}$  represents the function  $C(\mathbf{X}) = \bigotimes_{f \in \mathbf{F}} f$ .

**DEFINITION 4** (reasoning task). A **reasoning task** is a tuple  $\mathcal{P} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \Downarrow)$ , where  $(\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes)$  is a graphical model and  $\Downarrow$  is a marginalization operator (Definition 2). The reasoning task is to compute  $\Downarrow_{\mathbf{X}} C(\mathbf{X})$ .

For a reasoning task  $\mathcal{M} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \Downarrow)$  the choice of  $(\mathbf{A}, \otimes, \oplus)$  determines the combination  $\otimes$  and marginalization  $\Downarrow$  operators over functions, and thus the nature of the graphical model and its reasoning task. For example, if  $\mathbf{A}$  is the set of non-negative reals and  $\otimes$  is product, the graphical model is a Markov network or a Bayesian network. If  $\Downarrow$  is max, the task is to compute the Most Probable Explanation (MPE), while if  $\Downarrow$  is sum, the task is to compute the Probability of the Evidence.

It was shown in several works (e.g. [Shenoy and Shafer, 1990]) that, when the operators of a graphical models satisfy certain axioms, inference algorithms, such as variable elimination and join tree schemes, are sound and complete for the reasoning task. In this paper we assume graphical models obeying the Shenoy-Shafer axioms.

**DEFINITION 5** (A proper reasoning task). *Given a graphical model  $\mathcal{M} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes)$ , a reasoning task  $\mathcal{P} = (\mathcal{M}, \Downarrow)$  is **proper** if its  $\otimes$  and  $\Downarrow$  operators satisfy the following 3 axioms:*

*Axiom A1. Order of marginalization does not matter. Namely,  $\Downarrow_{X,Y} f(x, y) = \Downarrow_{Y,X} f(x, y)$ .*

*Axiom A2.  $\otimes$  is associative and commutative.*

*Axiom A3.  $\Downarrow$  distributes over  $\otimes$ . Namely,  $\Downarrow_{X,Y} f(x) \otimes f(y) = \Downarrow_X f(x) \otimes \Downarrow_Y f(y)$ .*

An important parameter of a graphical model, characterizing its complexity, is the *induced width*.

**DEFINITION 6** (induced graph, induced width). *The **induced graph** of a graphical model relative to ordering  $o$  is an undirected graph that has variables as its vertices. The edges of the graph are added by: 1) connecting all variables that are in the scope of the same function, 2) processing nodes from last to first in  $o$ , recursively connecting preceding neighbours of each node. The **induced width**  $w(o)$  relative to ordering  $o$  is the maximum number of preceding neighbours across all variables in the induced graph. The **induced width of the graphical model**  $w^*$  is the minimum induced width of all orderings.*

**DEFINITION 7** (bucket elimination). *Bucket elimination (BE) is an algorithmic framework that generalizes dynamic programming for many reasoning tasks [Dechter, 1999]. The input of BE (see Algorithm 1) is a reasoning task  $\mathcal{P} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \Downarrow)$  and an ordering  $o = (X_1, X_2, \dots, X_n)$ , dictating an elimination order for BE, from last to first. Each function from  $\mathbf{F}$  is placed in the bucket of its latest variable in  $o$ . The algorithm processes the buckets from  $X_n$  to  $X_1$ , computing for each  $\text{Bucket}_{X_i}$ , noted  $\mathbf{B}_i$ ,  $\Downarrow_{X_i} \bigotimes_{j=1}^n \lambda_j$ , where  $\lambda_j$  are the function in the  $\mathbf{B}_i$ , some of which are original  $f'_i$ s and some are earlier computed messages. The result of the computation is a new function, also called message, that is placed in the bucket of its latest variable in the ordering  $o$ .*

For example, algorithm *elim-opt*, which solves the optimization task, is obtained by substitution of the operators  $\Downarrow_X f = \max_{S-X} f$  and  $\otimes_j = \prod_j$ .

The message passing between buckets follows a bucket-tree structure.

**DEFINITION 8** (bucket tree). *Bucket elimination defines a bucket tree, where the bucket of each  $X_i$  is linked to the destination bucket of its message (called the parent bucket). A node of the bucket is associated with its bucket variable.*

---

#### Algorithm 1 Bucket elimination

---

**Input:** A reasoning task  $\mathcal{P} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \Downarrow)$ ; An ordering of variables  $o = \{X_1, \dots, X_n\}$ ;

**Output:** A zero-arity function  $\lambda_1 : \emptyset \rightarrow \mathbf{A}$  containing the solution of the reasoning task.

- 1: **Initialize:** Generate an ordered partition of functions in buckets  $\mathbf{B}_1, \dots, \mathbf{B}_n$ , where  $\mathbf{B}_i$  contains all the functions whose highest variable in their scope is  $X_i$ .
  - 2: **Backward:**
  - 3: **for**  $i \leftarrow n$  **down to** 1 **do**
  - 4:   Generate  $\lambda_i = (\bigotimes_{f \in \mathbf{B}_i} f) \Downarrow_{X_i}$
  - 5:   Place  $\lambda_i$  in the bucket  $\mathbf{B}_j$  where  $j$  is the largest-index variable in  $\text{var}(\lambda_i)$
  - 6: **end for**
  - 7: **Return:**  $\lambda_1$
- 

**Theorem 1.** [Dechter, 1999] *Given a graphical model and a proper reasoning task  $\mathcal{P} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \Downarrow)$  BE is sound and complete. Given an ordering  $o$ , the time and space complexity of  $BE(\mathcal{P})$  is exponential in the induced width of the ordering.*

The  $m$ -best task over a graphical model is defined formally next.

**DEFINITION 9** (m-best task). *The  $m$ -best task over a graphical model  $\mathcal{M}$  is to find  $m$  complete assignments  $\mathbf{T} = \{t_1, \dots, t_m\}$ , such that  $\forall t' \notin \mathbf{T} \forall t \in \mathbf{T}, C(t') \leq C(t)$ . The solution is the set of valuations  $\{C(t_1), \dots, C(t_m)\}$ , called  $m$ -best solutions.*

Most proofs are omitted for lack of space.

### 3 Algorithm *elim-m-opt*

This section provides the main contribution of our paper, presenting BE algorithm for the  $m$ -best task. We derive the algorithm through an example, having in mind the MPE (most probable explanation) task in probabilistic networks.

#### 3.1 Deriving the algorithm using an example

Consider a graphical model with four variables  $\{X, Y, Z, T\}$  having the following functions (for simplicity we use unnormalizes functions):

| $x$ | $z$ | $f_1(z, x)$ | $y$ | $z$ | $f_2(z, y)$ | $z$ | $t$ | $f_3(t, z)$ |
|-----|-----|-------------|-----|-----|-------------|-----|-----|-------------|
| 0   | 0   | 2           | 0   | 0   | 6           | 0   | 0   | 1           |
| 0   | 1   | 2           | 0   | 1   | 7           | 0   | 1   | 2           |
| 1   | 0   | 5           | 1   | 0   | 2           | 1   | 0   | 4           |
| 1   | 1   | 1           | 1   | 1   | 4           | 1   | 1   | 3           |
| 2   | 0   | 4           | 2   | 0   | 8           |     |     |             |
| 2   | 1   | 3           | 2   | 1   | 2           |     |     |             |

Finding the  $m$  best solutions to  $P(t, z, x, y) = f_3(t, z) \cdot f_1(z, x) \cdot f_2(z, y)$  can be expressed as finding  $Sol$ , defined by:

$$Sol = \underset{t, x, z, y}{\text{sort}}^m \left( f_3(t, z) \cdot f_1(z, x) \cdot f_2(z, y) \right) \quad (1)$$

where operator  $\underset{s}{\text{sort}}^m f(s)$  returns the first  $m$  elements of the set  $\{f(s)_{s \in D_s}\}$  or the entire set, if its size is less than  $m$ . When the argument of the operator is a subset of the function's scope, like  $Y$  in  $\underset{y}{\text{sort}}^m f(x, y)$ , then  $\underset{y}{\text{sort}}^m$  yields a set of the  $m$  best values of  $f(x, y)$  for every fixed  $x$  (by default kept ordered)  $\langle f(x, y^1(x)), \dots, f(x, y^m(x)) \rangle$  where  $f(x, y^j)$  is the  $j^{\text{th}}$  largest  $f$  value relative to  $D_Y$ , for a fixed  $X = x$ .

Since the output of the sort operators are vectors or sets, we will use *vector function* to denote a function whose range is a set (ordered set in our exposition and implementation). The *dimension of a vector function* is the size of its range vector.

Operator  $\underset{x}{\text{sort}}^m f(x)$  can be viewed as applying operator  $\max f(x)$   $m$  times, while at each step generating and then removing the best element from  $\{f(s)_{s \in D_s}\}$ , and it inherits, therefore, its distributive properties over multiplication. Due to this distributivity, we can apply symbolic manipulation and migrate each of the functions to the left of the  $\underset{x}{\text{sort}}^m$  operator over variables that are not in its scope. In our example we rewrite as:

$$Sol = \underset{t}{\text{sort}}^m \underset{z}{\text{sort}}^m \left( f_3(t, z) \underset{x}{\text{sort}}^m f_1(z, x) \right) \left( \underset{y}{\text{sort}}^m f_2(z, y) \right) \quad (2)$$

However, since the output of  $\underset{x}{\text{sort}}^m$  is a set, to make equation 2 well defined, we need to extend the multiplication operation to vector functions as a Cartesian product between all

pair-wise elements. This yields a new binary *combination* operator denoted by  $\overline{\otimes}$  between vector functions (generalizing combination with a scalar function in the obvious way), with which we rewrite expression 2 as

$$Sol = \underset{t}{\text{sort}}^m \underset{z}{\text{sort}}^m (f_3(t, z) \overline{\otimes} \underset{x}{\text{sort}}^m f_1(z, x)) \overline{\otimes} \underset{y}{\text{sort}}^m f_2(z, y) \quad (3)$$

The right to left computation in the order dictated by expression 3 can be carried out by a BE algorithm (see Figure 1) as follows. We assume that original input functions extends to vector functions e.g.,  $\overline{f}_i$  as  $\overline{f}_i(t) = \{f_i(t)\}$ . Then, the algorithm processes bucket  $\mathbf{B}_Y$  generating vector function  $\overline{\lambda}_Y(z)$  and bucket  $\mathbf{B}_X$  generating vector function  $\overline{\lambda}_X(z)$ :

| $z$ | $\overline{\lambda}_X(z)$ | $\overline{x}$ | $\overline{\lambda}_Y(z)$ | $\overline{y}$ |
|-----|---------------------------|----------------|---------------------------|----------------|
| 0   | {5,4,2}                   | {1,2,0}        | {8,6,2}                   | {2,0,1}        |
| 1   | {3,2,1}                   | {2,0,1}        | {7,4,2}                   | {0,1,2}        |

When processing  $\mathbf{B}_Z$ , we compute  $\overline{\lambda}_Z(t) = \underset{z}{\text{sort}}^m [f_3(t, z) \overline{\otimes} \overline{\lambda}_X(z) \overline{\otimes} \overline{\lambda}_Y(z)]$ . The result is a new vector function that has  $m^2$  elements for each tuple  $(t, z)$  as shown below.

| $t$ | $z$ | $f_3(t, z) \overline{\otimes} \overline{\lambda}_X(z) \overline{\otimes} \overline{\lambda}_Y(z)$ |
|-----|-----|---|
| 0   | 0   | {40, 32, 30, 16, 24, 12, 10, 8, 4}  |
| 0   | 1   | {84, 56, 48, 32, 28, 24, 16, 16, 8}   |
| 1   | 0   | {80, 64, 60, 48, 32, 24, 20, 16, 8}   |
| 1   | 1   | {63, 42, 36, 24, 21, 18, 12, 12, 6}   |

Applying  $\underset{z}{\text{sort}}^m$  to the resulting combination we get vector function  $\overline{\lambda}_Z(t)$  along with its variable assignments:

| $t$ | $\overline{\lambda}_Z(t)$ | $\langle \overline{x}, \overline{y}, \overline{z} \rangle$ |
|-----|---------------------------|--|
| 0   | {84, 56, 48}              | {(2, 0, 1), (0, 0, 1), (2, 1, 1)}                          |
| 1   | {80, 64, 63}              | {(1, 2, 0), (2, 2, 0), (2, 0, 1)}                          |

Finally, processing the last bucket yields the vector of  $m$  best solution costs for the entire problem and the corresponding assignments.

| $\overline{\lambda}_Z(t)$ | $\langle \overline{t}, \overline{z}, \overline{x}, \overline{y} \rangle$ |
|---------------------------|--|
| {84, 80, 64}              | {(0, 1, 2, 0), (1, 0, 1, 2), (1, 0, 2, 2)}                               |

Since we are interested in recovering at least one complete assignment for each  $m$ -best solution, the algorithm propagates the variable assignments along with the vector messages when processing each bucket. These variable assignments are generated using the *argsort* operator defined as follows.

**DEFINITION 10.** *Operator  $\text{argsort}_{X_i}^m f$  returns a vector function  $\overline{x}_i(t)$  such that  $\forall t \in D_{\text{var}(f) \setminus X_i}$ , where  $\langle f(t \cdot x_i^1), \dots, f(t \cdot x_i^m) \rangle$ , are the  $m$ -best costs extending  $t$  to  $X_i$ .*

In words,  $\overline{x}_i(t)$  is the vector of assignments to  $X_i$  that yields the  $m$ -best extensions to  $t$ .

#### 3.2 The Algorithm definition

We next more formally define the generalized operators to be used in the bucket-elimination algorithm *elim-m-opt*. Let  $\mathcal{M} = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \overline{\otimes})$  be a graphical model, over which we want to solve the  $m$ -best task. Let  $2^{\mathbf{A}}$  be the set of subsets of  $\mathbf{A}$ , and let us call a function with range in  $2^{\mathbf{A}}$  a vector function.

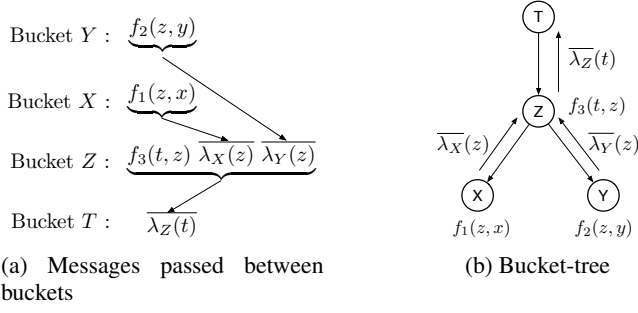


Figure 1: Example of applying *elim-m-opt*

**DEFINITION 11** (combination and addition over sets). Let  $S, T \in 2^A$ . Their combination, noted  $S \otimes T$ , is the set  $\{a \otimes b \mid a \in S, b \in T\}$ . Their addition, noted  $\text{sort}\{S, T\}$ , is the set of the  $m$ -best elements in the set  $S \cup T$ .

**DEFINITION 12** (combination and marginalization over vector functions). Let  $f : D_f \rightarrow 2^A$  and  $g : D_g \rightarrow 2^A$  be two vector functions. Their **combination**, noted  $f \otimes g$ , is a new vector function defined on scope  $\text{var}(f) \cup \text{var}(g)$  s.t.  $\forall t \in D_{\text{var}(f) \cup \text{var}(g)}, f \otimes g(t) = f(t_{\text{var}(f)}) \otimes g(t_{\text{var}(g)})$ . The **marginalization** of  $f$  over  $X_i \in \text{var}(f)$ , noted  $\text{sort}_{X_i}^m f$ , is a function over scope  $\text{var}(f) - \{X_i\}$  such that  $\forall t \in D_{\text{var}(f) - X_i}, (\text{sort}_{X_i}^m f)(t) = \text{sort}^m \{ \bigcup_{x \in D_{X_i}} f(t \cdot x) \}$

The bucket-elimination algorithm *elim-m-opt* is described in Algorithm 2 using the two new combination and marginalization operators of  $\otimes$  and  $\text{sort}^m$ . The algorithm processes the buckets from last to first as usual. The bucket computation of a message-function (step 2) is detailed in the next subsection. The message-function associates each tuple in its domain with the  $m$ -best costs-to-go restricted to the subproblem below the bucket variable in the bucket tree. For clarity we omit the generation of actual  $m$ -best solution assignments.

While the correctness of the algorithm can be given directly, we defer to Section 4, where we show that the combination and marginalization operators just defined are an instance of a more general formulation of the  $m$ -best task as a *proper reasoning task over a proper graphical model* [Shenoy and Shafer, 1990].

### 3.3 The bucket processing algorithm

We will next show that the messages computed in a bucket can be obtained more efficiently than through a brute-force application of  $\otimes$  followed by  $\text{sort}^m$ . Consider processing  $\mathbf{B}_Z$  (see Figure 1a). A brute-force computation of  $\bar{\lambda}_Z(t) = \text{sort}^m_z (f_3(z, t) \otimes \bar{\lambda}_Y(z) \otimes \bar{\lambda}_X(z))$  for each  $t$  combines  $f_3(z, t)$ ,  $\bar{\lambda}_Y(z)$  and  $\bar{\lambda}_X(z)$  for  $\forall z \in D_Z$  first. This results in a vector function with scope  $\{T, Z\}$  having  $m^2$  elements that we call *candidate elements* and denote by  $E(t, z)$ . The second step is to apply  $\text{sort}^m_z E(t, z)$  yielding the desired  $m$  best elements  $\bar{\lambda}_Z(t)$ .

However, since  $\bar{\lambda}_Y(z)$  and  $\bar{\lambda}_X(z)$  can be kept sorted, we can generate only a small subset of these  $m^2$  candidates as

### Algorithm 2 elim-m-opt algorithm

**Input:** A set of functions  $\mathbf{F} = \{f_1, \dots, f_n\}$  over scopes  $\{S_1, \dots, S_n\}$ ; An ordering of variables  $o = \{X_1, \dots, X_n\}$ ,

**Output:** A zero-arity function  $\lambda_1 : \emptyset \rightarrow 2^A$  containing the solution of the  $m$ -best optimization task.

- 1: **Initialize:** Transform each function  $f \in \mathbf{F}$  into a singleton vector function  $h(t) = \{f(t)\}$ ; Generate an ordered partition of vector functions  $h$  in buckets  $\mathbf{B}_1, \dots, \mathbf{B}_n$ , where  $\mathbf{B}_i$  contains all the functions whose highest variable in their scope is  $X_i$ .
- 2: **Backward:**
- 3: **for**  $i \leftarrow n$  **down to** 1 **do**
- 4:   Generate  $\lambda_i = \text{sort}_{X_i}^m (\bigotimes_{f \in \mathbf{B}_i} f)$
- 5:   Generate assignment  $\bar{x}_i = \text{argsort}_{X_i}^m (\bigotimes_{f \in \mathbf{B}_i} f)$ , concatenate with relevant elements of the previously generated assignment messages.
- 6:   Place  $\lambda_i$  and corresponding assignments in the bucket of the largest-index variable in  $\text{var}(\lambda_i)$
- 7: **end for**
- 8: **Return:**  $\lambda_1$

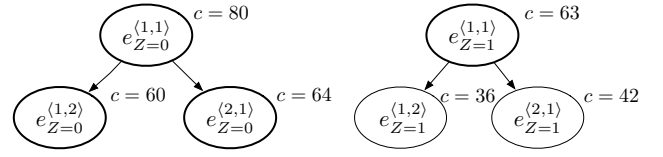


Figure 2: The explored search space for  $T = 0$  and  $m = 3$ . The resulting message is  $\bar{\lambda}_Z(1) = \{80, 64, 63\}$ .

follows. We denote by  $e_z^{(i,j)}(t)$  the candidate element obtained by the product of the scalar function value  $f_3(t, z)$  with the  $i^{\text{th}}$  element of  $\bar{\lambda}_Y(z)$  and  $j^{\text{th}}$  element of  $\bar{\lambda}_X(z)$ , having cost  $c_z^{(i,j)}(t) = f_3(t, z) \cdot \lambda_Y^i(z) \cdot \lambda_X^j(z)$ . We would like to generate the candidates  $e_z^{(i,j)}$  in decreasing order of their costs while taking their respective indices  $i$  and  $j$  into account.

The *child elements* of  $e_z^{(i,j)}(t)$ ,  $\text{children}(e_z^{(i,j)}(t))$  are obtained by replacing in the product either an element  $\lambda_Y^i(z)$  with  $\lambda_Y^{i+1}(z)$ , or  $\lambda_X^j(z)$  with  $\lambda_X^{j+1}(z)$ , but not both.

This leads to a forest-like search graph whose nodes are the candidate elements, where each search subspace corresponds to a different value of  $z$  denoted by  $G_{Z=z}$  and rooted in  $e_{Z=z}^{(1,1)}(t)$ . Clearly, the cost along any path from a node to its descendants is non-increasing. It is easy to see that the  $m$  best elements  $\bar{\lambda}_Z(t)$  can then be generated using a greedy best-first search across the forest search space  $G_{Z=0} \cup G_{Z=1}$ . It is easy to show that we do not need to keep more than  $m$  nodes on the OPEN list (the fringe of the search) at the same time. The general algorithm is described in Algorithm 3. The trace of the search for the elements of cost message  $\bar{\lambda}_Z(t = 1)$  for our running example is shown in Figure 2.

**Proposition 1.** Given a bucket of a variable  $X$  over scope  $S$  having  $j$  functions  $\{\lambda_1, \dots, \lambda_j\}$  of dimension  $m$ , where  $m$  is the number of best-solutions sought and  $k$  bounds the domain size, the complexity of bucket processing is  $O(k^{|S|} \cdot m \cdot j \log m)$ , where  $|S|$  is the scope size of  $S$ .

**Complexity of elim-m-opt:** Given  $n$  buckets, one for each

---

**Algorithm 3** Bucket processing

---

**Input:**  $\mathbf{B}_X$  of variable  $X$  containing a set of ordered  $m$ -vector functions  $\{\lambda_1(S_1, X), \dots, \lambda_d(S_d, X)\}$   
**Output:**  $m$ -vector function  $\bar{\lambda}_X(S)$ , where  $S = \cup_{i=1}^d S_i - X$ .

- 1: **for all**  $t \in D_S$  **do**
- 2:   **for all**  $x \in D_X$  **do**
- 3:      $OPEN \leftarrow e_{X=x}^{(1, \dots, 1)}(t)$ ; Sort OPEN;
- 4:   **end for**
- 5:   **while**  $j \leq m$ , by +1 **do**
- 6:      $n \leftarrow$  first element  $e_{X=x}^{(i_1, \dots, i_d)}(t)$  in OPEN. Remove  $n$  from OPEN;
- 7:      $\lambda_X^j(s) \leftarrow n$ ; {the  $j^{th}$  element is selected}
- 8:      $C \leftarrow$  children( $n$ ) =  $\{e_{X=x}^{(i_1, \dots, i_r+1, \dots, i_d)}(t) | r = 1..d\}$ ;
- 9:     Insert each  $c \in C$  into OPEN maintaining order based on its computed value. Check for duplicates; Retain the  $m$  best nodes in OPEN, discard the rest.
- 10:   **end while**
- 11: **end for**

---

variable  $X_i$ ,  $\mathbf{B}_i$  containing  $deg_i$  (i.e., the degree of the respective node in the bucket-tree) functions and at most  $w^*$  different variables, the total time complexity of *elim-m-opt* is  $\sum_{i=1}^n O(k^{w^*} m \cdot deg_i \log m)$ . Assuming  $deg_i \leq deg$  and since  $\sum_{i=1}^n deg_i \leq 2n$ , we get complexity expression:

$$O(k^{w^*} m \cdot \log m \cdot 2n) = O(nmk^{w^*} \log m)$$

The space complexity is dominated by the size of the messages between buckets, each containing  $m$  costs-to-go for each of  $O(k^{w^*})$  tuples. Having at most  $n$  such message yields the total space complexity of  $O(mnk^{w^*})$ . In summary:

**Theorem 2** (complexity of *elim-m-opt*). *Given a graphical model  $(X, \mathbf{D}, \mathbf{F}, \otimes)$  having  $n$  variables, an ordering  $o$ , induced-width of  $w^*$ , a bucket-tree degree bounded by  $deg$  and the domain size bounded by  $k$ , and an operator  $\Downarrow = \text{max}$ , the time complexity of *elim-m-opt* is  $O(k^{w^*} nm \cdot \log m)$  and its space complexity is  $O(mnk^{w^*})$ .*

Note that *elim-m-opt* is superior to applying Lawler's scheme [Lawler, 1972] to BE, which would lead to complexity of  $O(mn^2 k^{w^*})$ .

## 4 The M-best reasoning task

The  $m$ -best task defined via an optimization reasoning task is not itself in the form of a proper reasoning task over its graphical model. However, by some formal manipulation we can phrase it as a proper optimization reasoning task  $\tilde{\mathcal{P}}(m)$  over an associated graphical model  $\tilde{\mathcal{M}}(m)$ . This formal exercise is worthwhile, because by doing so all inference algorithms, and in particular BE, are immediately applicable. We use definition of combination given in Definition 12.

**DEFINITION 13** (m-best reasoning task). *Let  $\mathcal{M} = (X, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes)$  be a graphical model and  $\mathcal{P} = (\mathcal{M}, \Downarrow)$  be an optimization reasoning task. The m-best graphical model relative to  $\mathcal{M}$  is  $\tilde{\mathcal{M}}(m) = (X, \mathbf{D}, \mathbf{A}_m, \bar{\mathbf{F}}, \otimes)$ , where each function  $f : \mathbf{D}_f \rightarrow \mathbf{A}$  in  $\mathbf{F}$  is transformed into vector function:  $\bar{f}(t) = \{f(t)\}$ . The m-best reasoning task over  $\tilde{\mathcal{M}}(m)$*

is  $\tilde{\mathcal{P}}(m) = (\tilde{\mathcal{M}}(m), \text{sort}^m)$ , where  $\bar{\otimes}$  and  $\text{sort}^m$  are as in Definition 12.

We now can show that our original  $m$ -best task over a graphical model  $\mathcal{M}$  can be captured by  $\tilde{\mathcal{P}}(m)$  over  $\tilde{\mathcal{M}}(m)$ . Namely,

**Theorem 3.** *Let  $\mathcal{M} = (X, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes)$  be a graphical model and let  $\mathcal{P} = (\mathcal{M}, \text{max})$  be its proper optimization reasoning task. Then, the reasoning task  $\tilde{\mathcal{P}}(m) = (\tilde{\mathcal{M}}(m), \text{sort}^m)$  is proper and it coincides with the  $m$ -best task over  $\mathcal{M}$ .*

We can now conclude that the  $m$ -best task can be solved via traditional inference algorithms for reasoning tasks, and, in particular, by Bucket Elimination.

**Theorem 4.** *The m-best task over a graphical model  $\mathcal{M}$  can be solved by bucket-elimination as defined by  $BE(\tilde{\mathcal{P}}(m))$  and it coincides with  $\text{elim-m-opt}(\mathcal{M})$ .*

*Proof.* By Theorem 3, the  $m$ -best reasoning task  $\tilde{\mathcal{P}}(m)$  computes the  $m$ -best solutions of  $\mathcal{M}$  and since  $\tilde{\mathcal{P}}_m$  is a proper reasoning task, it can be solved exactly by BE. It is easy to see that  $\text{elim-m-opt}(\mathcal{M}) = BE(\tilde{\mathcal{P}})$ .  $\square$

## 5 Mini-Bucket Elimination for m-best Task

We now demonstrate the power of *elim-m-opt* in yielding useful bounds via mini-bucket-elimination.

### 5.1 The Mini-Bucket for the m-best

Mini-bucket Elimination (MBE) [Dechter and Rish, 2003] is an approximation designed to avoid the space and time complexity of BE. Consider a bucket  $\mathbf{B}_i$  and an integer bounding parameter  $z$ . MBE creates a  $z$ -partition  $Q = \{Q_1, \dots, Q_p\}$  of  $\mathbf{B}_i$ , where each set  $Q_j \in Q$ , called *mini-bucket*, includes no more than  $z$  variables. Then, each mini-bucket is processed separately, thus computing a set of messages  $\{\lambda_{ij}\}_{j=1}^p$ , where  $\lambda_{ij} = \Downarrow_{X_i} (\otimes_{f \in Q_j} f)$ . In general, greater values of  $z$  increase the quality of the bound.

**Theorem 5.** [Dechter and Rish, 2003] *Given a reasoning task  $\mathcal{P}$ , MBE computes a bound on  $\mathcal{P}$ . Given an integer control parameter  $z$ , the time and space complexity of MBE is exponential in  $z$ .*

Algorithm *m-best MBE* (*mbe-m-opt*) (given in Figure 4) is a straightforward extension of MBE to  $m$ -best reasoning task, where the combination and marginalization operators are the ones defined over vector functions. Algorithm *mbe-m-opt* solves  $m$ -best reasoning task  $\tilde{\mathcal{P}}(m)$  and its output is a *m-best bound* on the  $m$ -best solutions.

**DEFINITION 14** (m-best bound). *Let  $S = \{a_1, \dots, a_j\}$  and  $T = \{b_1, \dots, b_k\}$  be two sets (i.e.,  $S, T \in 2^A$ ).  $S$  is a m-best bound of  $T$  iff  $\forall 1 \leq i \leq |T|, b_i \leq a_i$ .*

**Theorem 6** (*mbe-m-opt* bound and complexity). *Given an m-best reasoning task  $\tilde{\mathcal{P}}(m)$ , mbe-m-opt computes an m-best bound on  $\tilde{\mathcal{P}}(m)$ . Given an integer control parameter  $z$ , the time and space complexity of mbe-m-opt is  $O(mnk^z \log(m))$  and  $O(mnk^z)$ , respectively, where  $k$  is the maximum domain size and  $n$  is the number of variables.*

---

**Algorithm 4** MBE-m-opt algorithm
 

---

**Input:** An  $m$ -best reasoning task  $\tilde{\mathcal{P}}(m) = (\mathbf{X}, \mathbf{D}, \mathbf{A}, \mathbf{F}, \otimes, \text{sort}^m)$ ; An ordering of variables  $o = \{X_1, \dots, X_n\}$ ; parameter  $z$ .

**Output:** bounds on each of the  $m$ -best solution costs and the corresponding assignments for the expanded set of variables (i.e., node duplication).

- 1: **Initialize:** Generate an ordered partition of functions  $\bar{f}(t) = \{f(t)\}$  into buckets  $\mathbf{B}_1, \dots, \mathbf{B}_n$ , where  $\mathbf{B}_i$  along  $o$ .
  - 2: **Backward:**
  - 3: **for**  $i \leftarrow n$  down to 1 (Processing bucket  $B_i$ ) **do**
  - 4: Partition functions in bucket  $B_i$  into  $\{Q_{i_1}, \dots, Q_{i_l}\}$ , where each  $Q_{i_j}$  has no more than  $z$  variables.
  - 5: Generate cost messages  $\lambda_{i_j} = \text{sort}_{X_{i_j}}^m(\overline{\otimes}_{f \in Q_{i_j}} f)$  and place each in the largest index variable in  $\text{var}(Q_{i_j})$
  - 6: **end for**
  - 7: **Return:** The set of all buckets, and the vector of  $m$ -best costs bounds in the first bucket.
- 

## 5.2 Using the $m$ -best bound to tighten the first-best bound

Here is a simple, but quite fundamental observation: whenever upper or lower bounds are generated by solving a relaxed version of a problem, the relaxed problem's solution set contains all the solutions to the original problem. We next discuss the ramification of this observation.

**Proposition 2.** *Given the  $m$ -best solutions costs generated by mbe-m-opt (for clarity we consider MPE problem, the results can be extended for other reasoning tasks)  $\tilde{C} = \{\tilde{p}_1 \geq \tilde{p}_2 \geq \dots, \geq \tilde{p}_m\}$ , let  $p^{\text{opt}}$  be the optimal value (the probability of the most probable explanation) and let  $j_0$  be the first index such that  $\tilde{p}_{j_0} = p^{\text{opt}}$ , or else we assign  $j_0 = m + 1$ . Then, if  $j_0 > m$ ,  $\tilde{p}_m$  is an upper bound on  $p^{\text{opt}}$ , which is as tight or tighter than all other  $\tilde{p}_1, \dots, \tilde{p}_{m-1}$ . In particular  $\tilde{p}_m$  is tighter than the bound  $\tilde{p}_1$ .*

*Proof.* Let  $\tilde{C} = \{\tilde{p}_1 \geq \tilde{p}_2 \geq \dots, \geq \tilde{p}_{N_1}\}$  be an ordered set of probabilities of all tuples over the relaxed problem (with duplicate variables). By the nature of any relaxation,  $\tilde{C}$  must also contain all the probability values associated with solutions of the original problem denoted by  $C = \{p_1 \geq \dots \geq p_{N_2}\}$ . Therefore, if  $j_0$  is the first index such that  $\tilde{p}_{j_0}$  coincides with  $p^{\text{opt}}$ , then clearly for all  $i < j_0$ ,  $p^{\text{opt}} \leq \tilde{p}_i$  with  $\tilde{p}_{j_0-1}$  being the tightest upper-bound. Also, when  $j_0 > m$  we have  $\tilde{p}_m \geq p^{\text{opt}}$   $\square$

In other words if  $j \leq m$ , we already have optimal value, otherwise we can use  $\tilde{p}_m$  as our better upper bound. Such tighter bounds would be useful during search algorithm such as A\*. It is essential therefore to decide efficiently if a bound coincides with the exact optimal cost. Luckily the nature of the MBE relaxation supplies us with an efficient decision scheme.

**Proposition 3.** *Given a set of bounds produced by mbe-m-opt  $\tilde{p}_1 \geq \tilde{p}_2 \geq \dots \geq \tilde{p}_m$ , deciding if  $\tilde{p}_j = p^{\text{opt}}$  can be done efficiently.*

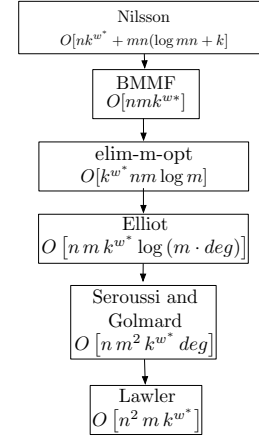


Figure 3: Complexity comparison. A parent node in the graph has a better complexity than its children.

*Proof.* mbe-m-opt provides both the bounds on the  $m$ -best costs and for each bound a corresponding tuple, where assignments to duplicated variables is maintained. The first assignment from these  $m$ -best bounds (going from largest to smallest) that corresponds to a tuple whose duplicate variables are assigned identical value, is optimal. And, if no such tuple is observed the optimal value is smaller than  $\tilde{p}_m$ . Since the above tests require just  $O(nm)$  steps applied to  $m$ -best assignments already obtained in polytime, the claim follows.  $\square$

## 6 Related work

**Comparing with exact schemes.** Lawler's approach, whose complexity is  $O(nmT(n))$ , where  $T(n)$  is the complexity of finding a single best solution, was applied by [Nilsson, 1998] to a join-tree. Instead of solving reformulated problems from scratch in each iteration, as Lawler does, Nilsson utilizes the results from previous computations, achieving worst case complexity of  $O(mT(n))$ . If applied to a bucket-tree Nilsson's algorithm is superior to all other schemes mentioned here, with run time of  $O(nk^{w*} + mn \log(mn) + mnk)$ . More recently Yanover and Weiss [Yanover and Weiss, 2004] developed a belief propagation approximation scheme for loopy graphs, called BMMF which also finds solutions iteratively. At each iteration BMMF uses loopy Belief Propagation to solve two new problems obtained by restricting the values of certain variables. When applied to junction tree it can function as an exact algorithm with complexity  $O(mnk^{w*})$ .

Two algorithms based on dynamic programming, similar to elim-m-opt, are [Seroussi and Golmard, 1994] and [Elliott, 2007]. Seroussi and Golmard extract the  $m$  solutions directly, by propagating the  $m$  best partial solutions along a junction tree. Given a junction tree with  $p$  cliques and branching degree  $deg$ , the complexity of the algorithm is  $O(m^2 p \cdot k^{w*} deg)$ . Elliott [Elliott, 2007] propagates the  $m$  best partial solutions along a representation called Valued And-Or Acyclic Graph, i.e., smooth deterministic decomposable negation normal form (sd-DNNF) [Darwiche, 2001]. The

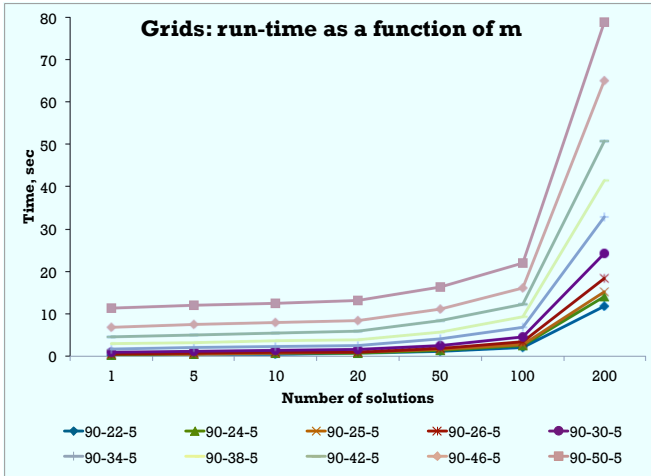


Figure 4: Run time as a function of number of solutions on grid instances

complexity of Elliot’s algorithm is  $O(nk^{w^*} m \log(mdeg))$ .

Figure 3 depicts the dominance relationships between the time complexities of the exact algorithms discussed, when specialized to a bucket tree. Clearly our *elim-m-opt* algorithm does not boast the best complexity, being dominated by the algorithms [Nilsson, 1998; Yanover and Weiss, 2004]. However, it demonstrates the direct applicability of established inference schemes to the generalized formulation of the  $m$  best solution problem as the *m-best reasoning problem*. Moreover, the main significance *elim-m-opt* is in the natural extension to an approximation scheme with guarantees on the solution quality that provides flexible trade off between accuracy and complexity.

**Comparing with approximation schemes.** In addition to BMMF, another extension of Nilsson’s and Lawler’s idea that yields an approximation scheme is an algorithm called STRIPES by [Fromer and Globerson, 2009]. They focus on  $m$ -MAP problem over binary Markov networks, solving each new subproblem by an LP relaxation. The algorithm solves the task exactly if the solutions to all LP relaxations are integral, and provides an upper bound of each  $m$  MAP assignments otherwise. In contrast, our algorithm *mbe-m-opt* can compute bounds over any graphical model (not only binary) and over a variety of  $m$ -best optimization tasks.

## 7 Empirical demonstrations

All experiments assume solving  $m$ -best MPE task. We evaluated empirically algorithm *mbe-m-opt* with  $m = \{1, 5, 10, 20, 50, 100, 200\}$  and with  $z$ -bound 10 on two sets of instances. The first set contained grid instances with a hundred to 2.5 thousand variables and tree-width from 12 to 50, the second - pedigree instances with several hundred variables and tree width from 15 to 30. Those instances were taken from the UAI 2008 evaluation. For clarity and space reasons we present only a subset of instances illustrating typical behaviour.

Figures 4 and 6 present the dependence of the run-time on

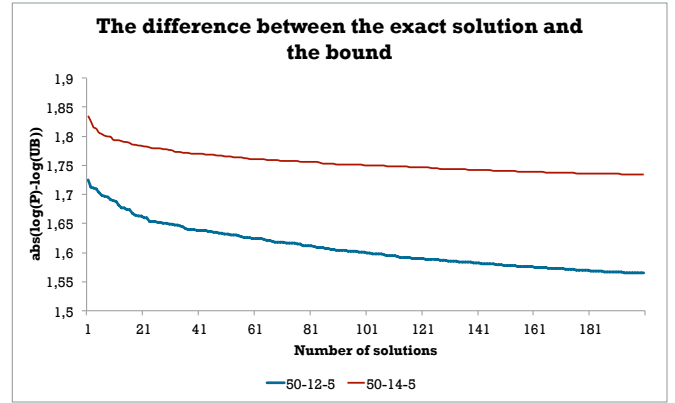


Figure 5: Accuracy of upper bound for different solutions on grid instances

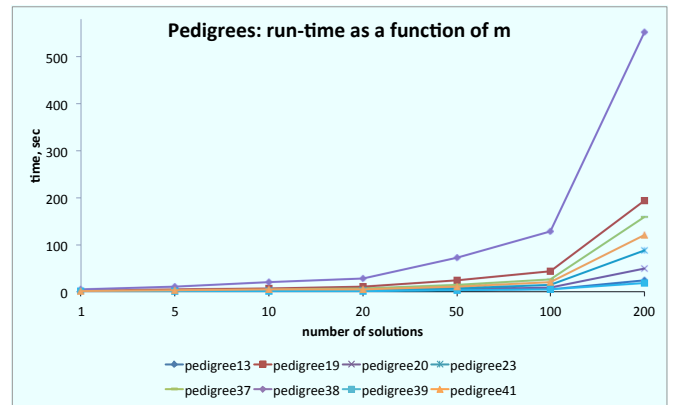


Figure 6: Run-time performance of *mbe-m-opt* on pedigrees

$m$ , for a few selected instances. Figure 5 shows the change in accuracy as a function of  $m$ . Accuracy is measured by the absolute values of the difference between the optimal solution and the bounds on the  $j^{th}$  solution up to  $m = 200$  solutions. For these grid instances as  $j$  increases, the bound on the cost of the  $j^{th}$  solution slowly approaches the exact best solution. This demonstrates that there is a potential of improving the bound on the optimal assignment using the  $m$ -best bounds as discussed in Section 5.2.

We carried some comparison with BMMF by [Yanover and Weiss, 2004] on randomly generated 10 by 10 grids. The run times of the algorithms are not comparable since our algorithm is implemented in C and BMMF in Matlab, which is inherently slower. For most instances that *mbe-m-opt* can solve exactly in under a second, BMMF takes more than 5 minutes. The algorithms also differ in the nature of the outputs: BMMF provides approximate solutions with no guarantees while *mbe-m-opt* generates bounds on all the  $m$ -best solutions. Still some information can be learned from viewing the two algorithms side by side as is demonstrated by a typical result in Figure 7. We know that in this case the solutions obtained with  $z$ -bound equal to 1000 are exact, while

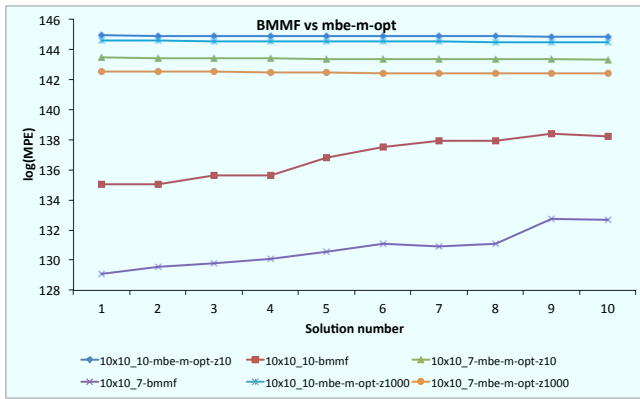


Figure 7: Comparison of  $mbe-m-opt$  with  $z$ -bounds 10 and 1000 vs. BMMF on random 10 by 10 grids

$z$ -bound equal to 10 yields an upper bound. BMMF outputs significantly less accurate results than  $mbe-m-opt$  with even a low  $z$ -bound. Admittedly, these experiments are quite preliminary and not conclusive.

## 8 Conclusions

We presented a new bucket-elimination algorithm for solving the  $m$ -best task over a graphical model, analyzed its performance and related it to other approaches in the literature. We proved the algorithm’s correctness by formally showing that the  $m$ -best task can be formulated as a reasoning task satisfying the Shenoy-Shafer axioms and, therefore, all inference algorithms and in particular our algorithm,  $elim-m-opt$ , are immediately applicable, sound and complete.

The significance of the proposed algorithm is primarily in providing an inference framework for the  $m$ -best task that can both suggest approximation schemes and yield heuristic advice. Indeed, optimization tasks that seek a single optimal solution are solved far more effectively by search (e.g., branch and bound and best-first search), than by variable elimination, because they can benefit from the bounding power of the guiding cost function. It is also likely that search will be more effective for  $m$ -best task. The promise of the  $elim-m-opt$  inference algorithm is in its potential to yield viable lower- and upper-bounds for the  $m$ -best solutions via the mini-bucket algorithm, as we discussed.

Furthermore, it could also lead to loopy propagation message-passing schemes that are now the most common way for approximations in graphical models, since those schemes are relaxation of exact message-passing schemes such as bucket-elimination. In particular, our algorithm can be extended into a loopy max-prod for the  $m$ -best task. This approach will yield a direct loopy-propagation for the  $m$ -best reasoning problem, while the approach by Yanover and Weiss uses loopy max-prod for solving a sequence of optimization problems in the style of Lawler’s approach. Moreover, all such approximation extensions would be applicable to the broad range of graphical models captured by the unifying framework of Shenoy and Shafer. Future work will focus on such extensions and on empirical evaluations of the emerging

schemes.

The empirical analysis we provided is only preliminary. Yet it shows that  $mbe-m-opt$  scales even better than worst-case predict as a function of  $m$ . Comparison with other exact and approximation algorithms is left for future work.

## Acknowledgments

This work was partially supported by NSF grants IIS-0713118 and IIS-1065618 and NIH grant 5R01HG004175-03. Thanks to the reviewers for their helpful feedback.

## References

- [Darwiche, 2001] A. Darwiche. Decomposable negation normal form. *Journal of the ACM (JACM)*, 48(4):608–647, 2001.
- [Dechter and Rish, 2003] R. Dechter and I. Rish. Mini-buckets: A general scheme for bounded inference. *Journal of the ACM (JACM)*, 50(2):107–153, 2003.
- [Dechter, 1999] R. Dechter. Bucket elimination: A unifying framework for reasoning. *Artificial Intelligence*, 113(1):41–85, 1999.
- [Elliott, 2007] P.H. Elliott. Extracting the K Best Solutions from a Valued And-Or Acyclic Graph. Master’s thesis, Massachusetts Institute of Technology, 2007.
- [Fromer and Globerson, 2009] M. Fromer and A. Globerson. An LP View of the  $M$ -best MAP problem. *Advances in Neural Information Processing Systems*, 22:567–575, 2009.
- [Kask et al., 2005] Kalev Kask, Rina Dechter, Javier Larrosa, and Avi Dechter. Unifying cluster-tree decompositions for automated reasoning. *Artificial Intelligence Journal*, 2005.
- [Lawler, 1972] E.L. Lawler. A procedure for computing the  $k$  best solutions to discrete optimization problems and its application to the shortest path problem. *Management Science*, 18(7):401–405, 1972.
- [Nilsson, 1998] D. Nilsson. An efficient algorithm for finding the  $M$  most probable configurations in probabilistic expert systems. *Statistics and Computing*, 8(2):159–173, 1998.
- [Pearl, 1988] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988.
- [Seroussi and Golmard, 1994] B. Seroussi and JL Golmard. An algorithm directly finding the  $K$  most probable configurations in Bayesian networks. *International Journal of Approximate Reasoning*, 11(3):205–233, 1994.
- [Shenoy and Shafer, 1990] P. Shenoy and G. Shafer. Axioms for probability and belief-function propagation. *Uncertainty in Artificial Intelligence*, 4:169–198, 1990.
- [Yanover and Weiss, 2004] C. Yanover and Y. Weiss. Finding the  $M$  Most Probable Configurations Using Loopy Belief Propagation. In *Advances in Neural Information Processing Systems 16*. The MIT Press, 2004.