

# Most Frugal Explanations in Bayesian Networks

Johan Kwisthout\*

*Donders Institute for Brain, Cognition and Behaviour, Radboud University Nijmegen,  
Nijmegen, The Netherlands*

---

## Abstract

Inferring the most probable explanation to a set of variables, given a partial observation of the remaining variables, is one of the canonical computational problems in Bayesian networks, with widespread applications in AI and beyond. This problem, known as MAP, is computationally intractable (NP-hard) and remains so even when only an approximate solution is sought. We propose a heuristic formulation of the MAP problem, denoted as Inference to the Most Frugal Explanation (MFE), based on the observation that many intermediate variables (that are neither observed nor to be explained) are irrelevant with respect to the outcome of the explanatory process. An explanation based on few samples (often even a singleton sample) from these irrelevant variables is typically almost as good as an explanation based on (the computationally costly) marginalization over these variables. We show that while MFE is computationally intractable in general (as is MAP), it can be tractably approximated under plausible situational constraints, and its inferences are fairly robust with respect to which intermediate variables are considered to be relevant.

*Keywords:* Bayesian Abduction, Parameterized Complexity, Approximation, Heuristics, Computational Complexity

---

## 1. Introduction

Abduction or inference to the best explanation refers to the process of finding a suitable explanation (the *explanans*) of observed data or phenom-

---

\*Corresponding author: Donders Institute for Brain, Cognition and Behaviour, Radboud University Nijmegen, PO Box 9104, 6500HE Nijmegen, The Netherlands; e-mail j.kwisthout@donders.ru.nl; telephone +31 (0) 24 3616288

4 ena (the *explananda*). In the last decades, Bayesian notions of abduction  
5 have emerged due to the widespread popularity of Bayesian or probabilistic  
6 techniques for representing and reasoning with knowledge [5, 26, 30, 47, 52].  
7 They are used in decision support systems in a wide range of problem domains  
8 [e.g., 7, 11, 21, 23, 32, 45, 64] and as computational models of economic, so-  
9 cial, or cognitive processes [10, 25, 33, 48, 58, 60]. The natural interpretation  
10 of ‘best’ in such models is ‘most probable’: the explanation that is the most  
11 probable one given the evidence, i.e., that has maximum posterior proba-  
12 bility, is seen as the hypothesis that best explains the available evidence;  
13 this explanation is traditionally referred to as the MAP explanation and the  
14 computational problem of inferring this explanation as the MAP problem.<sup>1</sup>

15 However, computing or even approximating the MAP explanation is com-  
16 putationally costly (i.e., NP-hard), especially when there are many interme-  
17 diate (neither observed nor to be explained) variables that may influence  
18 the explanation [1, 4, 51, 56]. To compute the posterior probability distri-  
19 bution of the explanation variables, all these intermediate variables need to  
20 be marginalized over. One way of dealing with this intractability might be  
21 by assuming modularity of knowledge representations, i.e., by assuming that  
22 these representations are informationally encapsulated and do not have ac-  
23 cess to background knowledge. However, this is problematic as we cannot  
24 know beforehand which elements of background knowledge or observations  
25 may be relevant for determining the best explanation [17, 19].

26 Fortunately, even when a full Bayesian computation may not be feasible  
27 in large networks, we need not constrain inferences only to small or dis-  
28 connected knowledge structures. It is known that in general the posterior  
29 probability distribution of a (discrete) Bayesian network is skewed, i.e., a  
30 few joint value assignments cover most of the probability space [13], and  
31 that typically only few of the variables in a network are relevant for a par-

---

<sup>1</sup>Other relationships have been proposed that compete in providing ‘sufficiently ratio-  
nal’ relations between observed phenomena and their explanation that can be used to  
describe *why* we judge one explanation to be preferred over another [28, 44]. Examples  
include *maximum likelihood* [29], which does not take the prior probabilities of the hy-  
potheses into account, the *conservative Bayesian* approach [6], *generalized Bayes* factor  
[66], and various Bayesian formalisms of *coherence theory* [5, 15, 26, 49]. While the poste-  
rior probability of such explanations is not the deciding criterion to prefer one explanation  
over another, it is typically so that explanations we consider to be good for other reasons  
also have a high posterior probability compared to alternative explanations [27, 44].

32 ticular inference query [14]. We propose to utilize this property of Bayesian  
33 networks in order to make tractable (approximate) inferences to the best  
34 explanation over large and unencapsulated knowledge structures. We in-  
35 troduce a heuristic formulation of MAP, denoted as Inference to the Most  
36 Frugal Explanation (MFE), that is explicitly designed to reflect that only  
37 few intermediate variables are typically relevant in real-world situations. In  
38 this formulation we partition the set of intermediate variables in the network  
39 into a set of ‘relevant’ intermediate variables that are marginalized over, and  
40 a set of ‘irrelevant’ intermediate variables that we sample from in order to  
41 estimate an explanation.

42 Note that in the MFE formalism there is marginalization over *some* of the  
43 intermediate variables (the variables that are considered to be relevant), but  
44 not over those intermediate variables that are not considered to be relevant.  
45 Thus, MFE can be seen as a ‘compromise’ between computing the expla-  
46 nation with maximum posterior probability, where one marginalizes over all  
47 intermediate variables, and the previously proposed Most Simple Explana-  
48 tion (MSE) formalism [35] where there is no marginalization at all, i.e., all  
49 intermediate variables are seen as irrelevant. We want to emphasize that the  
50 notions ‘relevant’ and ‘irrelevant’ in the problem definition denote *subjective*  
51 partitions of the intermediate variables; we will revisit this issue in Section  
52 3.1.

53 We claim that this heuristic formalism of the MAP problem exhibits the  
54 following desirable properties:

- 55 1. The knowledge structures are *isotropic*, i.e., they are such that, po-  
56 tentially, everything can be relevant to the outcome of an inference  
57 process. They are also *Quinean*: candidate explanations are sensitive  
58 to the entire belief system [17, 18].
- 59 2. The inferences are provably computationally tractable (either to com-  
60 pute exactly or to approximate) under realistic assumptions with re-  
61 spect to situational constraints [43, 53].
- 62 3. The resulting explanations have an optimal or close-to-optimal poste-  
63 rior probability in many cases, i.e., MFE actually ‘tracks truth’ in the  
64 terms of Glass [28].

65 In the remainder of this paper, we will discuss some needed preliminaries  
66 in Section 2. In Section 3 we discuss MFE in more detail. We give a more

67 formal definition, including a formal definition of relevance in the context of  
68 Bayesian networks, and show how MFE can be tractably approximated under  
69 realistic assumptions despite computational intractability of the problem in  
70 general. In Section 4 we show that MFE typically gives an explanation  
71 that has a close-to-optimal posterior probability, even if only a subset of  
72 the relevant variables is considered. We discuss how MFE performs under  
73 various scenarios (e.g., when there are few or many relevant variables, when  
74 there are many hypotheses that are almost equally likely, or when there is  
75 a misalignment between the *actual* relevant variables and the variables that  
76 are mistakenly presumed to be relevant). We conclude our paper in Section  
77 5.

## 78 2. Preliminaries

79 In this section we will introduce some preliminaries from Bayesian net-  
80 works, in particular the MAP problem as standard formalization of Bayesian  
81 abduction. We will discuss the ALARM network which we will use as a  
82 running example throughout this paper. Lastly, we introduce some needed  
83 concepts from complexity theory, in particular the complexity class PP, ora-  
84 cles, and fixed parameter tractability.

### 85 2.1. Bayesian networks and Bayesian abduction

86 A Bayesian or probabilistic network  $\mathcal{B}$  is a graphical structure that mod-  
87 els a set of stochastic variables, the conditional independences among these  
88 variables, and a joint probability distribution over these variables [52].  $\mathcal{B}$   
89 includes a directed acyclic graph  $\mathbf{G}_{\mathcal{B}} = (\mathbf{V}, \mathbf{A})$ , modeling the variables and  
90 conditional independences in the network, and a set of parameter probabili-  
91 ties  $\text{Pr}$  in the form of conditional probability tables (CPTs), capturing the  
92 strengths of the relationships between the variables. The network models a  
93 joint probability distribution  $\text{Pr}(\mathbf{V}) = \prod_{i=1}^n \text{Pr}(V_i \mid \pi(V_i))$  over its variables,  
94 where  $\pi(V_i)$  denotes the parents of  $V_i$  in  $\mathbf{G}_{\mathcal{B}}$ . We will use upper case letters  
95 to denote individual nodes in the network, upper case bold letters to denote  
96 sets of nodes, lower case letters to denote value assignments to nodes, and  
97 lower case bold letters to denote joint value assignments to sets of nodes. We  
98 will sometimes write  $\text{Pr}(\mathbf{x} \mid \mathbf{y})$  as a shorthand for  $\text{Pr}(\mathbf{X} = \mathbf{x} \mid \mathbf{Y} = \mathbf{y})$  if no  
99 ambiguity can occur.

100 In a Bayesian abduction task there are three types of variables: the *evid-*  
101 *dence* variables, the *explanation* variables, and a set of variables called *inter-*

102 *mediate* variables that are neither evidence nor explanation variables. The  
 103 evidence variables are instantiated, i.e., have been assigned a value; the joint  
 104 value assignment constitutes the explananda (what is to be explained, viz.,  
 105 the observations, data, or evidence). The explanation variables together form  
 106 the *hypothesis space*: a set of possible explanations for the observations; a  
 107 particular joint value assignment to these variables constitutes an explanans  
 108 (the actual explanation of the observations). When determining what is the  
 109 *best* explanation, typically we also need to consider other variables that are  
 110 not directly observed, nor are to be explained: the intermediate variables. By  
 111 convention, we will use  $\mathbf{E}$ ,  $\mathbf{H}$ , and  $\mathbf{I}$ , to denote the sets of evidence variables,  
 112 explanation variables, and intermediate variables, respectively. We will use  
 113  $\mathbf{e}$  to denote the evidence, viz., the (observed) joint value assignment to the  
 114 evidence variables.

115 The problem of inferring the *most probable* explanation, i.e., the joint  
 116 value assignment for the explanation set that has maximum posterior prob-  
 117 ability given the evidence, is defined as MAP, or also PARTIAL MAP or  
 118 MARGINAL MAP to emphasize that the probability of any such joint value  
 119 assignment is computed by marginalization over the intermediate variables.  
 120 MAP is formally defined as follows.

121 MAXIMUM A POSTERIORI PROBABILITY (MAP)

122 **Instance:** A Bayesian network  $\mathcal{B} = (\mathbf{G}_{\mathcal{B}}, \text{Pr})$ , where  $\mathbf{V}$  is partitioned into  
 123 evidence variables  $\mathbf{E}$  with joint value assignment  $\mathbf{e}$ , explanation variables  
 124  $\mathbf{H}$ , and intermediate variables  $\mathbf{I}$ .

125 **Output:** The joint value assignment  $\mathbf{h}$  to the nodes in  $\mathbf{H}$  that has  
 126 maximum posterior probability given the evidence  $\mathbf{e}$ .

## 127 2.2. The ALARM network

128 The ALARM network (Figure 1) will be used throughout this paper as a  
 129 running example. This network is constructed as a part of the ALARM moni-  
 130 toring system, providing users with text messages denoting possible problems  
 131 in anesthesia monitoring [2]. It consists of thirty-seven discrete random vari-  
 132 ables. Eight of these variables are designed as diagnostic variables that are to  
 133 be explained, indicating problems like pulmonary embolism or a kinked tube;  
 134 another sixteen variables indicate measurable or observable findings. The re-  
 135 maining thirteen variables are intermediate variables, i.e., they are neither  
 136 diagnostic variables, nor can be observed (in principle or in practice). Apart  
 137 from its practical use in the system described above, the ALARM network

138 is one of the most prominent benchmark networks in the Bayesian network  
 139 community.<sup>2</sup>

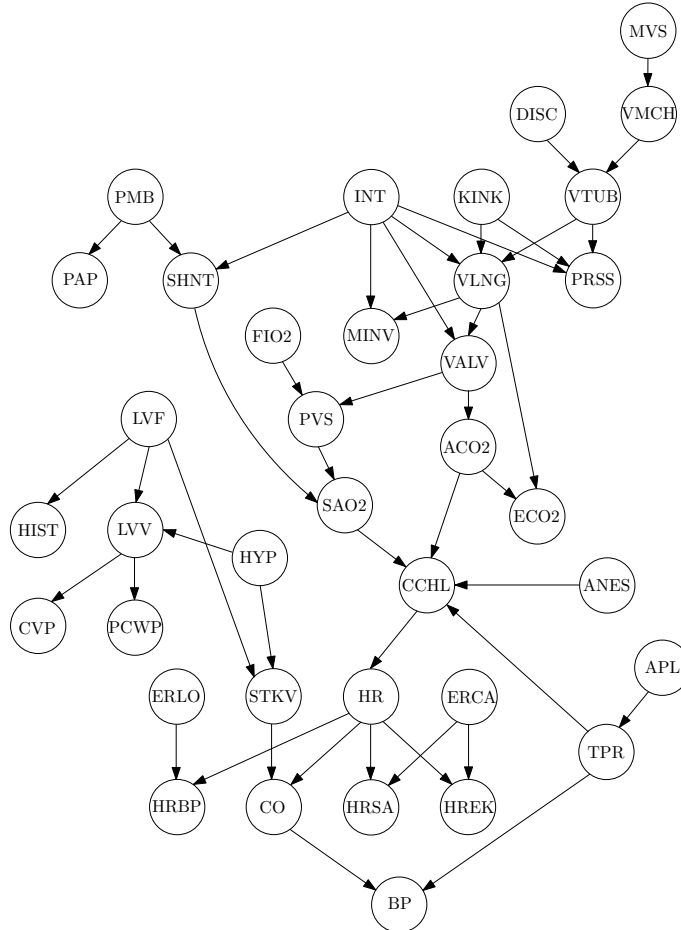


Figure 1: The ALARM network [2]

140 As an example, consider that a high breathing pressure was detected  
 141 (PRSS = high) and that minute ventilation was low (MINV = low); all  
 142 other observable variables take their default (i.e., non-alarming) value. From  
 143 these findings a probability of 0.92 for the diagnosis ‘kinked tube’ (KINK =

---

<sup>2</sup>See, e.g., <http://www.cs.huji.ac.il/site/labs/compbio/Repository/>

144 true) can be computed. Likewise, we can compute that the most probable  
145 joint explanation for the diagnostic variables, given that PCWP (pulmonary  
146 capillary wedge pressure) and BP (blood pressure) are high, is that HYP  
147 = true (hypovolemia, viz., loss of blood volume) and all other diagnostic  
148 variables are negative. This joint value assignment has probability 0.58. The  
149 second-best explanation (all diagnostic variables are negative, despite the  
150 two alarming conditions) has probability 0.11.

### 151 2.3. Complexity theory

152 In the remainder, we assume that the reader is familiar with basic con-  
153 cepts of computational complexity theory, such as Turing Machines, the com-  
154 plexity classes P and NP, and intractability proofs. For more background we  
155 refer to classical textbooks like [22] and [50]. In addition to these basic con-  
156 cepts we will introduce concepts that are in particular relevant to Bayesian  
157 computations, in particular Probabilistic Turing Machines, Oracle Turing  
158 Machines, the complexity class PP and the Counting Hierarchy; the inter-  
159 ested reader will find more background in [34] or [8]. Finally, we will briefly,  
160 and somewhat informally, introduce parameterized complexity theory. A  
161 more thorough introduction can be found in [12] or [16].

162 A Probabilistic Turing Machine (PTM) augments the more traditional  
163 Non-deterministic Turing Machine (NTM) with a probability distribution  
164 associated with each state transition. Without loss of generality we may  
165 assume that state transitions are binary and that the probability distribution  
166 at each transition is uniform. A PTM accepts a language  $L$  if the probability  
167 of ending in an accepting state when given some input  $x$  is strictly larger than  
168  $1/2$  if and only if  $x \in L$ . Given uniformly distributed binary state transitions  
169 this is exactly the case if the majority of computation paths accepts. The  
170 complexity class PP is defined as the class of languages accepted by some  
171 PTM in polynomial time. Observe that  $\text{NP} \subseteq \text{PP}$ ; the inclusion is thought to  
172 be strict. PP contains complete problems, the canonical one being MAJSAT:  
173 given a Boolean formula  $\phi$ , does the majority of truth assignments to the  
174 variables satisfy it?

175 An Oracle Turing Machine (OTM) is a Turing Machine enhanced with  
176 a so-called *oracle tape* and an oracle  $O$  for deciding membership queries  
177 for a particular language  $L_O$ . Apart from its usual operations, the OTM  
178 can write a string  $y$  on the oracle tape and ‘summon the oracle’. In the  
179 next state, the OTM will have either replaced the string with 1 if  $y \in L_O$ ,  
180 or 0 if  $y \notin L_O$ . The oracle can thus be seen as a ‘black box’ that answers

181 membership queries in constant time. Note that both accepting and rejecting  
182 answers of the oracle can be used. Various complexity classes are defined  
183 using oracles; for example, the class  $\text{NP}^{\text{PP}}$  includes exactly those languages  
184 that can be decided on an NTM with an oracle for PP-complete languages.  
185 Using the class PP and hierarchies of oracles the *Counting Hierarchy* [61] can  
186 be defined as a generalization of the Polynomial Hierarchy [59], including  
187 classes as  $\text{NP}^{\text{PP}}$ ,  $\text{P}^{\text{NP}^{\text{PP}}}$ , or  $\text{NP}^{\text{PP}^{\text{PP}}}$ . Canonical complete problems for such  
188 classes include various SATISFIABILITY variants, using the quantifiers  $\forall$ ,  $\exists$ ,  
189 and MAJ to bind subsets of variables [61, 63].

190 Sometimes problems are intractable (i.e., NP-hard) in general, but be-  
191 come tractable if some *parameters* of the problem can be assumed to be  
192 small. Informally, a problem is called fixed-parameter tractable for a pa-  
193 rameter  $k$  (or a set of parameters  $\{k_1, \dots, k_m\}$ ) if it can be solved in time,  
194 exponential (or even worse) *only* in  $k$  and polynomial in the input size  $|x|$ .  
195 In practice, this means that problem instances can be solved efficiently, even  
196 when the problem is NP-hard in general, if  $k$  is known to be small. If an  
197 NP-hard problem  $\Pi$  is fixed-parameter tractable for a particular parameter  
198 set  $k$  then  $k$  is denoted a *source of complexity* [53] of  $\Pi$ : bounding  $k$  renders  
199 the problem tractable, whereas leaving  $k$  unbounded ensures intractability  
200 under usual complexity-theoretic assumptions like  $\text{P} \neq \text{NP}$ . On the other  
201 hand, if  $\Pi$  remains NP-hard independent of the value of parameter  $k$ , then  $\Pi$   
202 is para-NP-hard with respect to  $k$ : bounding  $k$  does not render the problem  
203 tractable. The notion of fixed-parameter tractability can be extended to deal  
204 with *rational*, rather than integer, parameters [36]. Informally, if a problem  
205 is fixed-rational tractable for a (rational) parameter  $k$ , then the problem can  
206 be solved tractably if  $k$  is close to 0 (or, depending on the definition, to 1).  
207 For readability, we will liberally mix integer and rational parameters in the  
208 remainder.

### 209 3. Most Frugal Explanations

210 In real-world applications there are many intermediate variables that are  
211 neither observed nor to be explained, yet may influence the explanation.  
212 Some of these variables can considerably affect the outcome of the abduction  
213 process. Most of these variables, however, are irrelevant as they are not  
214 expected to influence the outcome of the abduction process in all but maybe  
215 the very rarest of cases [14]. To compute the most probable explanation of  
216 the evidence, however, one needs to marginalize over all these variables, that



217 is, take their prior or conditional probability distribution into account. This  
218 seems like a waste of computing resources in cases where we might as well  
219 have assigned an arbitrary value to these variables and still arrive at the  
220 same explanation.

221 One way of ensuring tractability of inference may be by ‘weeding out’  
222 the irrelevant aspects in the knowledge structure prior to inference, reducing  
223 the network to a simplified version. For example, one might try to iden-  
224 tify intermediate variables in the network that are conditionally independent  
225 of the explanation variables, given the evidence. While this can be done  
226 tractably in principle [24], it may still leave us with many variables that are  
227 conditionally dependent, yet do not influence the most probable explanation  
228 of the evidence. These variables are still in a sense redundant for finding  
229 explanations, as illustrated in the following example.

230 **Example 1.** Consider in the ALARM network the observations that PCWP  
231 and BP are high and the other observable variables take their non-alarming  
232 states. The actual value of ACO2 does not influence the most probable value  
233 of the observable variables in the network, i.e.,  $\operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{h}, \mathbf{e}, \mathbf{i}, \text{ACO2} =$   
234  $\text{high}) = \operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{h}, \mathbf{e}, \mathbf{i}, \text{ACO2} = \text{mid}) = \operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{h}, \mathbf{e}, \mathbf{i}, \text{ACO2} =$   
235  $\text{low})$  for every joint value assignment  $\mathbf{i}$  to the intermediate variables other  
236 than ACO2. However, ACO2 is not conditionally independent of (e.g.) KINK  
237 given the observed evidence variables.

238 An alternative to only selecting those intermediate variables that are con-  
239 ditionally dependent on the explanation variables is to apply a stronger cri-  
240 terion for relevance, e.g., selecting only those variables whose value may  
241 potentially change the most probable explanation. However, finding these  
242 variables itself would require potentially intractable computations as we will  
243 illustrate in Section 3.1 and formally prove in the Appendix. Furthermore,  
244 we might want to even constrain the number of variables to select even more  
245 by demanding not only that their value *might* change the most probable ex-  
246 planation (e.g., in some extraordinary combination of values for the other  
247 variables), but in fact actually *does* change the most probable explanation  
248 in a non-trivial number of situations. In addition, it is preferable to have a  
249 means of trading off the quality of a solution and the time needed to obtain  
250 a solution.

251 **Example 2** (Adapted from [35]). Mr. Jones typically comes to work by  
252 train. Today Mr. Jones is late while he has been seen to leave his house at

253 the usual time. One explanation can be that the train is delayed. However, it  
254 might also be the case that Mr. Jones was the unlucky individual who walked  
255 through 11th Street at 8.03 AM and was shot during an armed bank robbery,  
256 while mistakenly taken for a police constable. When trying to explain why  
257 Mr. Jones is not at his desk on 8.30 AM, there are a number of variables  
258 we might take into account, for example whether he has to change trains.  
259 A whole lot of variables are typically not taken into account because they  
260 are normally not relevant in most of the cases, for example the color of Mr.  
261 Jones’s coat, or whether walked on the left or right pavement in 11th Street.  
262 Only in the awkward coincidence that Mr. Jones was in the wrong place at  
263 the wrong time they become relevant to explain why he is not at work.

264 Our approach is not to reduce the network to only include those interme-  
265 diate variables we consider to be relevant and do inference on the resulting  
266 pruned network. In contrast, we propose that (the computationally costly)  
267 marginalization is done only on a subset of the intermediate variables (the  
268 variables that are considered to be relevant), and that a sampling strategy  
269 is used for the remaining intermediate variables that are not considered to  
270 be relevant. Such a sampling strategy may be very simple (‘decide using a  
271 singleton sample’) or more complex (‘compute the best explanation on  $N$   
272 samples and take a majority vote’). This allows for a trade-off between time  
273 to compute a solution and the quality of the result obtained, by having both  
274 a degree of freedom on which variables to include in the set of relevant inter-  
275 mediate variables and a degree of freedom on how many samples to take on  
276 the remaining intermediate variables. In Section 4 we will discuss the effects  
277 of such choices using computer simulations on random networks.

278 We now formally define the Most Frugal Explanation problem as follows<sup>3</sup>:

279 **MOST FRUGAL EXPLANATION (MFE)**

280 **Instance:** A Bayesian network  $\mathcal{B}$ , partitioned into a set of observed  
281 evidence variables  $\mathbf{E}$ , a set of explanation variables  $\mathbf{H}$ , a set of ‘relevant’

---

<sup>3</sup>To improve readability, this formulation does not explicate how to deal with the fol-  
lowing borderline cases: a) for any given joint value assignment to the irrelevant interme-  
diate variables, multiple hypotheses have the same posterior probability; and b) multiple  
hypotheses are most probable for the same maximum number of (possibly distinct) hy-  
potheses. The implementation of the algorithm described in Section 3.3 dealt with both  
these borderline cases by randomly selecting one of the competing hypotheses in case of a  
tie.

282 intermediate variables  $\mathbf{I}^+$  that are marginalized over, and a set of  
283 ‘irrelevant’ intermediate variables  $\mathbf{I}^-$  that are not marginalized over.  
284 **Output:** The joint value assignment to the variables in the explanation set  
285 that is most probable for the maximum number of joint value assignments  
286 to the irrelevant intermediate variables.

287 The approach sketched above guarantees that, as in the MAP problem,  
288 the knowledge structures remain both isotropic and Quinean, i.e., everything  
289 still can be relevant to the outcome of the inference process and the candi-  
290 date explanations remain sensitive to the entire belief system, as claimed in  
291 Section 1. For example, when new evidence arises (say, that we learn of a  
292 bank robbery where an innocent passerby was shot), the color of Mr. Jones’s  
293 coat suddenly may become relevant to explaining his absence. We will elab-  
294 orate on the *tractability* claim in Section 3.2 and on the *tracking truth* claim  
295 in Section 4.2.

296 **Example 3.** As in the previous example, we assume that in the ALARM  
297 network PCWP and BP have been observed to be high and the other ob-  
298 servable variables take their non-alarming states. Furthermore, let us assume  
299 that we consider VTUB, SHNT, VLNG, VALV and LVV to be relevant in-  
300 termediate variables, and VMCH, PVS, ACO2, CCHL, ERLO, STKV, HR,  
301 and ERCA to be irrelevant variables. The most *frugal* joint explanation for  
302 the diagnostic variables is still that HYP = true while all other diagnostic  
303 variables are negative: in 31% of the joint value assignments to these irrele-  
304 vant intermediate variables, this is the most probable explanation. In 16% of  
305 the assignments ‘all negative’ is the most probable explanation, and in 24%  
306 of the assignments HYP = true and INT = one sided (one sided intubation,  
307 rather than normal) is the most probable explanation of the observations.  
308 If, in addition, we also consider VMCH, PVS, and STKV to be relevant,  
309 then every joint value assignment to ACO2, CCHL, ERLO, HR, and ERCA  
310 will have HYP = true as the most probable explanation for the observations.  
311 In other words, rather than marginalizing over these variables, we might  
312 have assigned just an arbitrary joint value assignment to these variables, de-  
313 creasing the computational burden. If we had considered less intermediate  
314 variables to be relevant, this strategy may still often work, but has a chance  
315 of error, if we pick a sample for which a different explanation is the most  
316 probable one. We can decrease this error by taking more samples and take  
317 a majority vote.

318 Note that MFE is not *guaranteed* to give the MAP explanation, unless we  
319 marginalize over all intermediate variables (i.e., consider all variables to be  
320 relevant). When the set of irrelevant variables is non-empty, the most frugal  
321 explanation may differ from the MAP explanation, even when using a voting  
322 strategy based on *all* joint value assignments to the irrelevant intermediate  
323 variables, since both explanations are computed differently. Take for example  
324 the small network with two ternary variables  $H$  with values  $\{h_1, h_2, h_3\}$  and  
325  $I$  with values  $\{i_1, i_2, i_3\}$ , with  $I$  uniformly distributed and  $H$  conditioned on  
326  $I$  as follows:

$$\begin{aligned} \Pr(h_1 \mid i_1) &= 0.4 & \Pr(h_2 \mid I = i_1) &= 0.3 & \Pr(h_3 \mid i_1) &= 0.3 \\ \Pr(h_1 \mid i_2) &= 0.4 & \Pr(h_2 \mid I = i_2) &= 0.3 & \Pr(h_3 \mid i_2) &= 0.3 \\ \Pr(h_1 \mid i_3) &= 0.1 & \Pr(h_2 \mid I = i_3) &= 0.6 & \Pr(h_3 \mid i_3) &= 0.3 \end{aligned}$$

327 Now, the most *probable* explanation of  $H$ , marginalized on  $I$ , would be  $H =$   
328  $h_2$ , but the most *frugal* explanation of  $H$  with irrelevant variable  $I$  would be  
329  $H = h_1$  as this is the most probable explanation for two out of three value  
330 assignments to  $I$ . Only in borderline cases MAP and MFE are guaranteed  
331 to give the same results independent of the number of samples taken and  
332 the partition in relevant and irrelevant intermediate variables. This will, for  
333 example, be the case when the MAP explanation has a probability of 1 and  
334 all the intermediate variables are uniformly distributed. In this case, every  
335 joint value assignment to any subset of the intermediate variables gives the  
336 MAP explanation as most frugal explanation.<sup>4</sup>

### 337 3.1. Relevance

338 Until now, we have quite liberally used the notion ‘relevance’. It is im-  
339 portant here to note that we consider the relevance of *intermediate* variables.  
340 This is in contrast with Shimony’s well-known account [55] where relevance  
341 is a property of *explanation* variables, i.e., the relevance criterion partitions  
342 the non-observed variables in MAP variables—that are to be explained—and  
343 intermediate variables that do not need to be assigned a value in the expla-  
344 nation. In this paper we assume that the partition between the explanation  
345 variables  $\mathbf{H}$  and the intermediate variables  $\mathbf{I}$  is already made. However, in  
346 our model we again partition the intermediate variables  $\mathbf{I}$  and perform full  
347 inference only on the *relevant* intermediate variables  $\mathbf{I}^+$ .

---

<sup>4</sup>We thank one of the anonymous reviewers for this observation.

348 It will be clear that the formal notion of (conditional) independence is  
 349 too strong to capture relevance as we understand it: even if an intermediate  
 350 variable is formally not independent of all the explanation variables, condi-  
 351 tioned on the observed evidence variables, its influence may still be too small  
 352 to have an impact on which explanation to select as the most probable as we  
 353 saw in the previous sub-section. In contrast, we define relevance as a statis-  
 354 tical property of an intermediate variable that is partly based on Druzdzel  
 355 and Suermondt’s [14] definition of relevance of variables in a Bayesian model,  
 356 and partly on Wilson and Sperber’s [65] relevance theory, and is related to  
 357 the definition in [37]. According to Druzdzel and Suermondt a variable in  
 358 a Bayesian model is relevant for a set  $\mathbf{T}$  of variables, given an observation  
 359  $\mathbf{E}$ , if it is “needed to reason about the impact of observing  $\mathbf{E}$  on  $\mathbf{T}$ ” [14,  
 360 p.60]. Our operationalization of “needed to reason” is inspired by Wilson  
 361 and Sperber, who state that “an input is relevant to an individual when its  
 362 processing in a context of available assumptions yields (...) a worthwhile  
 363 difference to the individual’s representation of the world” [65, p.608]. The  
 364 term ‘worthwhile difference’ in this quote refers to the balance between the  
 365 actual effects of processing that particular input and the effort required to  
 366 do so. We therefore define the relevance of an intermediate variable as a  
 367 *measure*, indicating how sensitive explanations are to changes in its value  
 368 assignment. Informally, an intermediate variable  $I$  has a low relevance when  
 369 there are only few possible worlds in which the most probable explanation  
 370 changes when the value of  $I$  changes.<sup>5</sup>

371 **Definition 4.** Let  $\mathcal{B} = (\mathbf{G}_{\mathcal{B}}, \text{Pr})$  be a Bayesian network partitioned into  
 372 evidence nodes  $\mathbf{E}$  with joint value assignment  $\mathbf{e}$ , intermediate nodes  $\mathbf{I}$ , and  
 373 an explanation set  $\mathbf{H}$ . Let  $I \in \mathbf{I}$ , and let  $\Omega(\mathbf{I} \setminus \{I\})$  denote the set of joint  
 374 value assignments to the intermediate variables other than  $I$ . The *relevance*  
 375 of  $I$ , denoted as  $\mathcal{R}(I)$ , is the fraction of joint value assignments  $\mathbf{i}$  in  $\Omega(\mathbf{I} \setminus \{I\})$   
 376 for which  $\text{argmax}_{\mathbf{h}} \text{Pr}(\mathbf{h}, \mathbf{e}, \mathbf{i}, i)$  is not identical for all  $i \in \Omega(I)$ .

377 As computing the relevance of a variable  $I$  is NP-hard, i.e., intractable  
 378 in general (see the Appendix for a formal proof), we introduce the notion  
 379 *estimated relevance of  $I$*  as a subjective assessment of  $\mathcal{R}(I)$  which may or may  
 380 not correspond to the actual value. Such a subjective assessment might be

---

<sup>5</sup>Note that the *size of the effect* on the probability distribution of  $\mathbf{H}$  is not taken into  
 account here, only that the distribution alters sufficiently enough for the most probable  
 joint value assignment to ‘flip over’ to a different value.

381 based on heuristics, previous knowledge, or by approximating the relevance,  
 382 e.g., by sampling a few instances of  $\Omega(\mathbf{I} \setminus \{I\})$ . Where confusion may arise, we  
 383 will use the term *intrinsic relevance* to refer to the actual statistical property  
 384 ‘relevance’ of  $I$ , in contrast to the subjective assessment thereof. Note that  
 385 both intrinsic and estimated relevance of a variable are relative to a particular  
 386 set of candidate explanations  $\mathbf{H}$ , and conditional on a particular observation,  
 387 i.e., a value assignment  $\mathbf{e}$  to the evidence nodes  $\mathbf{E}$ .

388 **Example 5.** Let, in the ALARM network, pulmonary capillary wedge pres-  
 389 sure and blood pressure be high, and let all other observable variables take  
 390 their non-alarming default values. The intrinsic relevance of the intermediate  
 391 variables for the diagnosis is given in Figure 2.

392 When solving an MFE problem, we marginalize over the ‘relevant inter-  
 393 mediate variables’. This assumes some (subjective) threshold on the (esti-  
 394 mated or intrinsic) relevance of the intermediate variables that determine  
 395 which variables are considered to be relevant and which are considered to  
 396 be irrelevant. For example, if the threshold would be 0.85 then only SHNT  
 397 and LVV would be relevant intermediate variables in the ALARM network,  
 398 but if the threshold would be 0.40 then also VMCH, VTUB, VLNG, VALV,  
 399 and STKV would be relevant variables. That influences the results, as the  
 400 distribution of MFE explanations tends to be flatter when less variables are  
 401 marginalized over. With a threshold of 0.85 there are 24 explanations that  
 402 are sometimes the most probable explanation, with the actual MAP expla-  
 403 nation occurring most often (26%). With a threshold of 0.40 there are just  
 404 three such explanations, with the MAP explanation occurring in 75% of the  
 405 cases. Thus, the distribution of MFE explanations is typically more ‘skewed’  
 406 towards one explanation when more variables are considered to be relevant.

### 407 3.2. Complexity Analysis

408 To assess the computational complexity of MFE, we first define a decision  
 409 variant.

410 MOST FRUGAL EXPLANATION (MFE)

411 **Instance:** A Bayesian network  $\mathcal{B} = (\mathbf{G}_{\mathcal{B}}, \text{Pr})$ , where  $\mathbf{V}$  is partitioned into a  
 412 set of evidence nodes  $\mathbf{E}$  with a joint value assignment  $\mathbf{e}$ , an explanation set  
 413  $\mathbf{H}$ , a set of *relevant* intermediate variables  $\mathbf{I}^+$ , and a set of *irrelevant*  
 414 intermediate variables  $\mathbf{I}^-$ ; a rational number  $0 \leq q < 1$  and an integer  
 415  $0 \leq k < |\Omega(\mathbf{I}^-)|$ .

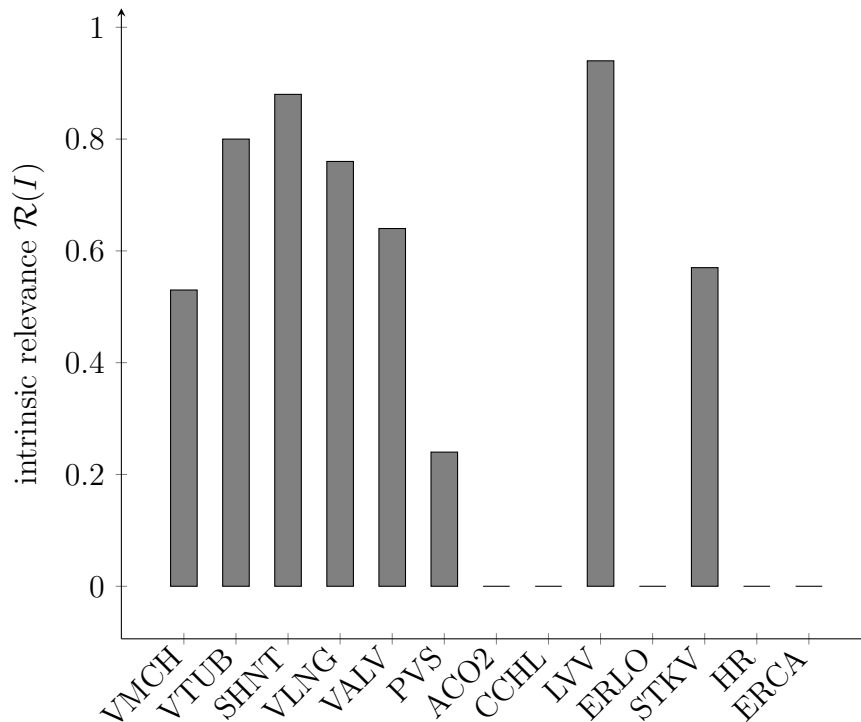


Figure 2: The intrinsic relevance of the intermediate variables of the ALARM network for the diagnostic variables given  $PCWP = \text{TRUE}$  and  $BP = \text{TRUE}$ . Note that the left ventricular end-diastolic blood volume (LVV) is highly relevant for the diagnosis, while the amount of catecholamines in the blood (CCHL) is irrelevant given these observations

---

416 **Question:** Is there a joint value assignment  $\mathbf{h}$  to the nodes in  $\mathbf{H}$  such that  
417 for more than  $k$  disjoint joint value assignments  $\mathbf{i}$  to  $\mathbf{I}^-$ ,  $\Pr(\mathbf{h}, \mathbf{i}, \mathbf{e}) > q$ ?

418 It will be immediately clear that MFE is intractable, as it has the  $\text{NP}^{\text{PP}}$ -  
419 complete MAP [51] and MSE [35] problems as special cases for  $\mathbf{I}^- = \emptyset$ ,  
420 respectively  $\mathbf{I}^+ = \emptyset$ . In this section we show that MFE happens to be even  
421 harder, viz., that it is  $\text{NP}^{\text{PPPP}}$ -complete, making it one of few real world-  
422 problems that are complete for that class<sup>6</sup>. The canonical SATISFIABILITY-

---

<sup>6</sup>Informally, one could imagine that for solving MFE one needs to counter *three* sources of complexity: selecting a joint value assignment out of potentially exponentially many

423 variant that is complete for this class is E-MAJMAJSAT, defined as follows  
 424 [61].

425 E-MAJMAJSAT

426 **Instance:** A Boolean formula  $\phi$  whose  $n$  variables  $x_1 \dots x_n$  are partitioned  
 427 into three sets  $\mathbf{E} = x_1 \dots x_k$ ,  $\mathbf{M}_1 = x_{k+1} \dots x_l$ , and  $\mathbf{M}_2 = x_{l+1} \dots x_n$  for  
 428 some numbers  $k, l$  with  $1 \leq k \leq l \leq n$ .

429 **Question:** Is there a truth assignment to the variables in  $\mathbf{E}$  such that for  
 430 the majority of truth assignments to the variables in  $\mathbf{M}_1$  it holds, that the  
 431 majority of truth assignments to the variables in  $\mathbf{M}_2$  yield a satisfying  
 432 truth instantiation to  $\mathbf{E} \cup \mathbf{M}_1 \cup \mathbf{M}_2$ ?

433 As an example, consider the formula  $\phi_{\text{ex}} = x_1 \wedge (x_2 \vee x_3) \wedge (x_4 \vee x_5)$  with  
 434  $\mathbf{E} = \{x_1\}$ ,  $\mathbf{M}_1 = \{x_2, x_3\}$  and  $\mathbf{M}_2 = \{x_4, x_5\}$ . This is a *yes* example of E-  
 435 MAJMAJSAT: for  $x_1 = \text{TRUE}$ , three out of four truth assignments to  $\{x_2, x_3\}$   
 436 (all but  $x_2 = x_3 = \text{FALSE}$ ) are such that the majority of truth assignments  
 437 to  $\{x_4, x_5\}$  satisfy  $\phi_{\text{ex}}$ .

438 To prove  $\text{NP}^{\text{PPP}}$ -completeness of the MFE problem, we construct a  
 439 Bayesian network  $\mathcal{B}_\phi$  from an E-MAJMAJSAT instance  $(\phi, \mathbf{E}, \mathbf{M}_1, \mathbf{M}_2)$ . For  
 440 each propositional variable  $x_i$  in  $\phi$ , a binary stochastic variable  $X_i$  is added  
 441 to  $\mathcal{B}_\phi$ , with uniformly distributed values TRUE and FALSE. These stochastic  
 442 variables in  $\mathcal{B}_\phi$  are three-partitioned into sets  $\mathbf{X}_\mathbf{E}$ ,  $\mathbf{X}_{\mathbf{M}_1}$ , and  $\mathbf{X}_{\mathbf{M}_2}$  according  
 443 to the partition of  $\phi$ . For each logical operator in  $\phi$  an additional binary  
 444 variable in  $\mathcal{B}_\phi$  is introduced, whose parents are the variables that correspond  
 445 to the input of the operator, and whose conditional probability table is equal  
 446 to the truth table of that operator. The variable associated with the top-  
 447 level operator in  $\phi$  is denoted as  $V_\phi$ , the set of variables associated with the  
 448 remaining operators is denoted as  $\text{Op}_\phi$ . Figure 3 shows the graphical struc-  
 449 ture of the Bayesian network constructed for the example E-MAJMAJSAT  
 450 instance given above.

451 **Theorem 6.** MFE is  $\text{NP}^{\text{PPP}}$ -complete.

---

candidate assignments to the explanation set; solving an inference problem over the vari-  
 ables in the set  $\mathbf{I}^+$ , and deciding upon a threshold of the joint value assignments to the  
 set  $\mathbf{I}^-$ . While the ‘selecting’ aspect is typically associated with problems in NP, ‘infer-  
 ence’ and ‘threshold testing’ are typically associated with problems in PP. Hence, as these  
 three sub-problems work on top of each other, the complexity class that corresponds to  
 this problem is  $\text{NP}^{\text{PPP}}$ .



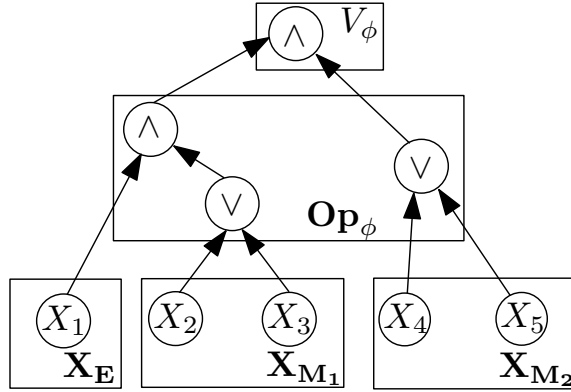


Figure 3: Example of the construction of  $\mathcal{B}_{\phi_{\text{ex}}}$  for the Boolean formula  $\phi_{\text{ex}} = x_1 \wedge (x_2 \vee x_3) \wedge (x_4 \vee x_5)$

452 *Proof.* Membership in  $\text{NP}^{\text{PPP}}$  follows from the following algorithm: non-  
 453 deterministically guess a value assignment  $\mathbf{h}$ , and test whether there are at  
 454 least  $k$  joint value assignments  $\mathbf{i}^-$  to  $\mathbf{I}^-$  such that  $\Pr(\mathbf{h}, \mathbf{i}^-, \mathbf{e}) > q$ . This  
 455 inference problem can be decided (for given value assignment  $\mathbf{h}$  and  $\mathbf{i}^-$ ) us-  
 456 ing a PTM capable of deciding INFERENCE (marginalizing over the variables  
 457 in  $\mathbf{I}^+$ ). We can decide whether there are at least  $k$  such joint value assign-  
 458 ments  $\mathbf{i}^-$  using an PTM capable of threshold counting. Thus, as both decid-  
 459 ing INFERENCE and threshold counting are PP-complete problems, we can  
 460 solve this problem by augmenting an NTM with an oracle for PPP-complete  
 461 problems; note that we cannot ‘merge’ both oracles as the ‘threshold’ oracle  
 462 machine must accept inputs for which the INFERENCE oracle answers ‘no’ as  
 463 well as inputs for which the oracle answers ‘yes’.

464 To prove  $\text{NP}^{\text{PPP}}$ -hardness, we reduce MFE from E-MAJMAJSAT. We  
 465 fix  $q = 1/2$  and  $k = |\Omega(\mathbf{I}^-)|/2$ . Let  $(\phi, \mathbf{E}, \mathbf{M}_1, \mathbf{M}_2)$  be an instance of E-  
 466 MAJMAJSAT and let  $\mathcal{B}_\phi$  be the network constructed from that instance as  
 467 shown above. We claim the following: If and only if there exists a satisfying  
 468 solution to  $(\phi, \mathbf{E}, \mathbf{M}_1, \mathbf{M}_2)$ , there is a joint value assignment to  $\mathbf{x}_\mathbf{E}$  such that  
 469  $\Pr(V_\phi = \text{TRUE}, \mathbf{x}_\mathbf{E}, \mathbf{x}_{\mathbf{M}_2}) > 1/2$  for the majority of joint value assignments  
 470  $\mathbf{x}_{\mathbf{M}_2}$  to  $\mathbf{X}_{\mathbf{M}_2}$ .

471  $\Rightarrow$  Let  $(\phi, \mathbf{E}, \mathbf{M}_1, \mathbf{M}_2)$  denote the satisfiable E-MAJMAJSAT instance. Note  
 472 that in  $\mathcal{B}_\phi$  any particular joint value assignment  $\mathbf{x}_\mathbf{E} \cup \mathbf{x}_{\mathbf{M}_1} \cup \mathbf{x}_{\mathbf{M}_2}$  to

473  $\mathbf{X}_{\mathbf{E}} \cup \mathbf{X}_{\mathbf{M}_1} \cup \mathbf{X}_{\mathbf{M}_2}$  yields  $\Pr(V_\phi = \text{TRUE}, \mathbf{x}_{\mathbf{E}}, \mathbf{x}_{\mathbf{M}_1}, \mathbf{x}_{\mathbf{M}_2}) = 1$ , if and only  
474 if the corresponding truth assignment to  $\mathbf{E} \cup \mathbf{M}_1 \cup \mathbf{M}_2$  satisfies  $\phi$ , and  
475 0 otherwise. When marginalizing over  $\mathbf{x}_{\mathbf{M}_1}$  (and  $\mathbf{Op}_\phi$ ) we thus have  
476 that a joint value assignment  $\mathbf{x}_{\mathbf{E}} \cup \mathbf{x}_{\mathbf{M}_2}$  to  $\mathbf{X}_{\mathbf{E}} \cup \mathbf{X}_{\mathbf{M}_2}$  yields  $\Pr(V_\phi =$   
477  $\text{TRUE}, \mathbf{x}_{\mathbf{E}}, \mathbf{x}_{\mathbf{M}_2}) > 1/2$  if and only if the majority of truth assignments  
478 to  $\mathbf{M}_1$ , together with the given truth assignment to  $\mathbf{E} \cup \mathbf{M}_2$ , satisfy  $\phi$ .  
479 Thus, given that this is the case for the majority of truth assignments  
480 to  $\mathbf{M}_2$ , we have that  $\Pr(V_\phi = \text{TRUE}, \mathbf{x}_{\mathbf{E}}, \mathbf{x}_{\mathbf{M}_2}) > 1/2$  for the majority  
481 of joint value assignments  $\mathbf{x}_{\mathbf{M}_2}$  to  $\mathbf{X}_{\mathbf{M}_2}$ . We conclude that the corre-  
482 sponding instance  $(\mathcal{B}_\phi, V_\phi = \text{TRUE}, \mathbf{X}_{\mathbf{E}}, \mathbf{X}_{\mathbf{M}_1} \cup \mathbf{Op}_\phi, \mathbf{X}_{\mathbf{M}_2}, 1/2, |\Omega(\mathbf{X}_{\mathbf{M}_2})|/2)$   
483 of MFE is satisfiable.

484  $\Leftarrow$  Let  $(\mathcal{B}_\phi, V_\phi = \text{TRUE}, \mathbf{X}_{\mathbf{E}}, \mathbf{X}_{\mathbf{M}_1} \cup \mathbf{Op}_\phi, \mathbf{X}_{\mathbf{M}_2}, 1/2, |\Omega(\mathbf{X}_{\mathbf{M}_2})|/2)$  be a satisfiable  
485 instance of MFE, i.e., there exists a joint value assignment  $\mathbf{x}_{\mathbf{E}}$  to  $\mathbf{X}_{\mathbf{E}}$   
486 such that for the majority of joint value assignments  $\mathbf{x}_{\mathbf{M}_2}$  to  $\mathbf{X}_{\mathbf{M}_2}$ ,  
487  $\Pr(V_\phi = \text{TRUE}, \mathbf{x}_{\mathbf{E}}, \mathbf{x}_{\mathbf{M}_2}) > 1/2$ . For each of these assignments  $\mathbf{x}_{\mathbf{M}_2}$  to  
488  $\mathbf{X}_{\mathbf{M}_2}$ ,  $\Pr(V_\phi = \text{TRUE}, \mathbf{x}_{\mathbf{E}}, \mathbf{x}_{\mathbf{M}_2}) > 1/2$  if and only if the majority of joint  
489 value assignments  $\mathbf{x}_{\mathbf{M}_1}$  to  $\mathbf{X}_{\mathbf{M}_1}$  satisfy  $\phi$ .

490 Since the reduction can be done in polynomial time, this proves that MFE  
491 is  $\text{NP}^{\text{PPP}}$ -complete.  $\square$

492 Given the intractability of MFE for unconstrained domains, it may not be  
493 clear how MFE as a heuristic mechanism for Bayesian abduction can scale  
494 up to task situations of real-world complexity. One approach may be to  
495 seek to approximate MFE, rather than to compute it exactly. Unfortunately,  
496 *approximating* MFE is NP-hard as well. Given that MFE has both MAP and  
497 MSE as special cases (for  $\mathbf{I}^- = \emptyset$ , respectively  $\mathbf{I}^+ = \emptyset$ ), it is intractable to  
498 infer an explanation that has a probability that is close to optimal [51], that  
499 is similar to the most frugal explanation [40], or that is likely to be the most  
500 frugal explanation with a bounded margin of error [42]. By and of itself,  
501 for unconstrained domains, approximation of MFE does not buy tractability  
502 [43].

### 503 3.3. Parameterized Complexity

504 An alternative approach to ensure computational tractability is to study  
505 how the complexity of MFE depends on situational constraints. This ap-  
506 proach has firm roots in the theory of parameterized complexity as described  
507 in Section 2. Building on known fixed parameter tractability results for MAP

508 [36] and MSE [42], we will consider the parameters *treewidth* and *cardinality*  
 509 of the Bayesian network, the *size* of  $\mathbf{I}^+$ , and a *decisiveness* measure on the  
 510 probability distribution. An overview is given in Table 1.

Parameter	Description
Treewidth ( $t$ )	A measure on the network topology [see, e.g., 3].
Cardinality ( $c$ )	The maximum number of values any variable can take.
#Relevants ( $ \mathbf{I}^+ $ )	The number of relevant intermediate variables that we marginalize over.
Decisiveness ( $d$ )	A measure on the probability distribution [42], denoting the probability that for a given evidence set $\mathbf{E}$ with evidence $\mathbf{e}$ and explanation set $\mathbf{H}$ , two random joint value assignments $\mathbf{i}_1$ and $\mathbf{i}_2$ to the irrelevant variables $\mathbf{I}^-$ would yield the same most probable explanations. Decisiveness is high if a robust majority of the joint value assignments to $\mathbf{I}^-$ yields a particular most probable explanation.

Table 1: Overview of parameters for MFE.

511 For  $\mathbf{I}^+ = \emptyset$ , MAP can be solved in  $O(c^t \cdot n)$  for a network with  $n$  variables,  
 512 and since  $\Pr(X = x) = \sum_{y \in \Omega(Y)} \Pr(X = x, Y = y)$ , we have that MAP can  
 513 be solved in  $O(c^t \cdot c^{|\mathbf{I}^+|} \cdot n)$ . Note that even when we can tractably decide upon  
 514 the most probable explanation for a given joint value assignment  $\mathbf{i}$  to  $\mathbf{I}^-$  (i.e.,  
 515 when  $c$ ,  $t$ , and  $|\mathbf{I}^+|$  are bounded) we still need to test at least  $\lfloor c^{|\mathbf{I}^-|}/2 \rfloor + 1$  joint  
 516 value assignments to  $|\mathbf{I}^-|$  to decide MFE exactly, even when  $d = 1$ . However,  
 517 in that case we can tractably find an explanation that is *very likely* to be the  
 518 MFE if  $d$  is close to 1. Consider the following algorithm for MFE (adapted  
 519 from [35]):

---

**Algorithm 1** Compute the Most Frugal Explanation

---

Sampled-MFE( $\mathcal{B}, \mathbf{H}, \mathbf{I}^+, \mathbf{I}^-, \mathbf{e}, N$ )

- 1: **for**  $n = 1$  to  $N$  **do**
  - 2:   Choose  $\mathbf{i} \in \mathbf{I}^-$  at random
  - 3:   Determine  $\mathbf{h} = \operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{H} = \mathbf{h}, \mathbf{i}, \mathbf{e})$
  - 4:   Collate the joint value assignments  $\mathbf{h}$
  - 5: **end for**
  - 6: Decide upon the joint value assignment  $\mathbf{h}_{\text{maj}}$  that was picked most often
  - 7: **return**  $\mathbf{h}_{\text{maj}}$
- 

520     This randomized algorithm repeatedly picks a joint value assignment  
521  $\mathbf{i} \in \mathbf{I}^-$  at random, determines the most probable explanation, and at the end  
522 decides upon which explanation was found most often. Due to its stochastic  
523 nature, this algorithm is not guaranteed to give correct answers all the  
524 time. However, the error margin  $\epsilon$  can be made sufficiently low by choosing  
525  $N$  large enough. If there are only two competing most probable explanations,  
526 the threshold value of  $N$  can be computed using the *Chernoff bound*:  
527  $N \geq \frac{1}{(p-1/2)^2} \ln 1/\sqrt{\epsilon}$  (more sophisticated methods are to be used to compute  
528 or approximate  $N$  when there are more than two competing explanations).  
529 Assume we require an error margin of less than 0.1, then the number of re-  
530 peats depends on the probability  $p$  of picking a joint value assignment  $\mathbf{i}$  for  
531 which  $\mathbf{h}_{\text{maj}}$  is the most probable explanation. This probability corresponds  
532 to the *decisiveness* parameter  $d$  that was introduced in Table 1. If decisiveness  
533 is high (say  $d = 0.85$ ), then  $N$  can be fairly low ( $N \geq 10$ ), however, if  
534 the distribution of explanations is very flat, and consequently, decisiveness is  
535 low, then an exponential number of repetitions is needed.

536     If  $d$  is bounded (i.e., larger than a particular fixed threshold) we thus need  
537 only polynomially many repetitions to obtain any constant error rate. When  
538 in addition determining the most probable explanation is easy—in particular,  
539 when the treewidth and cardinality of  $\mathcal{B}$  are low and there are few relevant  
540 variables in the set  $\mathbf{I}^+$ —the algorithm thus runs in polynomial time, and thus  
541 MFE can be decided in polynomial time, with a small possibility of error.

542 *3.4. Discussion*

543     In the previous subsections we showed that MFE is intractable in general,  
544 both to compute exactly and to approximate, yet can be tractably approxi-  
545 mated (with a so-called expectation-approximation [42]) when the treewidth

546 of the network is low, the cardinality of the variables is small, the number of  
547 relevant intermediate variables is low, *and* the probability distribution for a  
548 given explanation set  $\mathbf{H}$ , evidence  $\mathbf{e}$  and relevant intermediate variables  $\mathbf{I}^+$   
549 is fairly decisive, i.e., skewed towards a single MFE explanation. We also  
550 know that MAP can be tractably computed exactly<sup>7</sup> when the treewidth of  
551 the network is low, the cardinality of the variables is small, and either the  
552 MAP explanation has a high probability, or the total number of intermediate  
553 variables is low [36]. How do these constraints compare to each other?

554 For MAP, the constraint on the total number of intermediate variables  
555 seems implausible. In real-world knowledge structures there are many inter-  
556 mediate variables, and while only some of them may contribute to the MAP  
557 explanation, we still need to marginalize over all of them to compute MAP.  
558 Likewise, when there are many candidate hypotheses, it is not obvious that  
559 the most probable one has a high (i.e., close to 1) probability. Note that the  
560 actual fixed-parameter tractable algorithm [4, 36] has a running time with  
561  $\frac{\log p}{\log 1-p}$  in the exponent, where  $p$  denotes the probability of the MAP explana-  
562 tion. This exponent quickly grows with decreasing  $p$ . Furthermore, treewidth  
563 and cardinality actually refer to the treewidth of the *reduced* junction tree,  
564 where observed variables are absorbed in the cliques. Given that we sample  
565 over the set  $\mathbf{I}^-$  in MFE, but not in MAP, both parameters (treewidth and  
566 cardinality) will typically have much lower values in MFE as compared to  
567 MAP. That is, it is more plausible that these constraints are met in MFE  
568 than that they are met in MAP.

569 Given the theoretical considerations in [14] it seems plausible that the  
570 *decisiveness* constraint is met in many practical situations. Surely, one could  
571 argue that the fixed parameter tractability of MFE is misguided, as the set  
572 of candidate hypotheses and the observations are given in the input of the  
573 formal problem, and it is known beforehand what the relevant variables are.  
574 Thus, the problem of finding candidate hypotheses, the problem of deciding  
575 what counts as evidence, and the problem of deciding which variables are  
576 relevant are left out of the problem definition. We acknowledge that this  
577 is indeed the case, and that the problem of non-demonstrative inference is  
578 much broader than ‘merely’ inferring the best explanation out of a set of

---

<sup>7</sup>There are to the best of our knowledge no stronger (or even *different*) fixed parameter tractable results for *approximate* MAP than the results listed above for exact computations.

579 candidate explanations [39]; yet, this is no different for MAP, at least when  
580 it comes to deciding upon the candidate hypotheses and the observations.  
581 With respect to the partition between irrelevant and relevant intermediate  
582 variables we will show in Section 4 that MFE is fairly robust: including even  
583 a few variables with a high intrinsic relevance will suffice to find relatively  
584 good MFE explanations.

## 585 4. Simulations

586 In Section 3 we illustrated, using the ALARM example, that computing  
587 MFE can give similar results as when MAP is computed, while requiring  
588 less variables to be marginalized over. In this section, we will simulate MFE  
589 on random graphs to obtain empirical results to support that claim. We  
590 will also illustrate that, in order to obtain a good explanation using only  
591 few samples, the decisiveness of the probability distribution indeed must be  
592 high. Finally we show how MFE behaves under various scenarios where  
593 the intrinsic and estimated relevance of the intermediate variables (i.e., the  
594 actual relevance and the subjective assessment thereof) do not match. As  
595 the goal of these simulations is to explore how MFE behaves under scenarios  
596 that can be considered either natural (occurring in real-world networks) or  
597 artificial, we use randomly generated networks, rather than an existing set  
598 of benchmark networks, like the ALARM network, in our simulations.

### 599 4.1. Method

600 We generated 100 random Bayesian networks, each consisting of 40 vari-  
601 ables, using the (second) method described in [51]. Each variable had ei-  
602 ther two, three, or four possible values, and the in-degree of the nodes was  
603 limited to five. With each variable, a random conditional probability dis-  
604 tribution was associated. We randomly selected five explanation variables  
605 and five evidence variables, and set a random joint value assignment to the  
606 evidence variables. Given the variation on the cardinality of the variables,  
607 the number of candidate joint value assignments to the explanation variables  
608 could vary from  $2^5$  to  $4^5$ ; in practice, it ranged from 48 to 576 (mean 208.5,  
609 standard deviation 107.4). See also the on-line supplementary materials:  
610 <http://www.dcc.ru.nl/~johank/MFE/>.

611 Using the Bayes Net Toolbox for MATLAB [46] we computed, for each  
612 network, the posterior distribution over the explanation variables, approx-  
613 imated the relevance of each intermediate variable, and approximated the

614 MFE distribution under various conditions. The results presented below are  
 615 based on 91 random networks. The MATLAB software was unable to com-  
 616 pute the MAP of seven networks due to memory limitations, and the results  
 617 of two networks were lost due to hardware failure. In Figures 4 and 5 some  
 618 typical results are given for illustrative purposes.

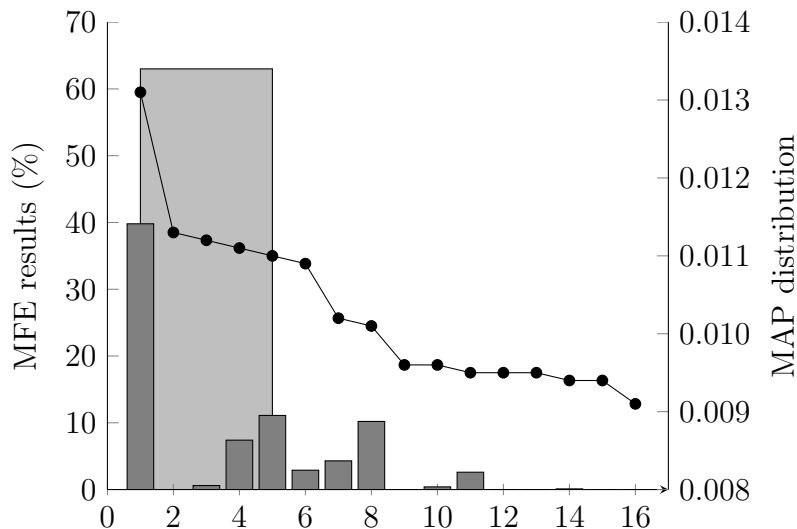


Figure 4: MAP distribution and MFE results for the 16 most probable joint value assignments of one of the random networks (#99) for a particular set of relevant intermediate variables, using 1000 samples. The light gray bar denotes the cumulative MFE result of the five most probable joint value assignments. Note that the most probable joint value assignment (which has a probability of 0.0131) is also the most frugal explanation, as it is the MAP for about 40% of the joint value assignments to the irrelevant intermediate variables. The ‘second-best MAP’, while it has a relative high posterior probability, is *always* ‘second-best’: there are no joint value assignments to the irrelevant intermediate variables in which this particular explanation has the highest probability. There *are* other explanations, with a lower posterior probability, that become the most probable explanation for some particular value assignments to these irrelevant intermediate variables. Note that in this situation there is no error as the most probable and most frugal explanation are identical.

---

619 *4.2. Tracking Truth*

620 We compared the MAP explanation with the MFE explanation using 100  
 621 samples of the irrelevant variables, varying the  $I^+/I^-$  partition. In particular

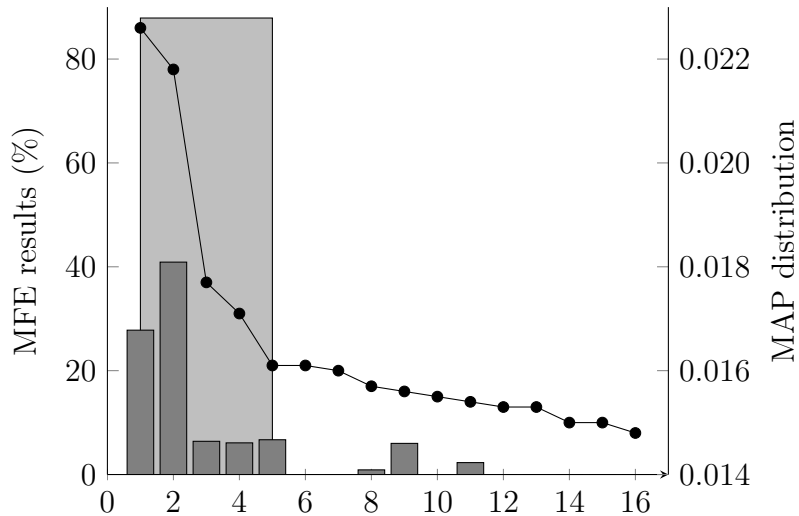


Figure 5: A similar plot as in Figure 4, but in this random network (#68) the most frugal explanation is the second most probable explanation, yielding a difference between the ‘marginalizing’ and the ‘sampling’ approach. Note, however, that both explanations are almost as good: they differ in a single variable, and the probability ratio is 0.965, meaning that the probability of the most frugal explanation is only slightly lower than the probability of the most probable explanation.

---

622 we compared the explanations where all variables are deemed irrelevant ( $I^+ =$   
623  $\emptyset$ ), where  $I^+$  consisted of the five intermediate variables with the highest  
624 relevance, and where  $I^+$  consisted of the intermediate variables that have a  
625 relevance of more than 0.00, 0.05, 0.10, 0.25, respectively 0.50. To assess  
626 how similar the most frugal explanations are to the MAP results, we used  
627 three different error measures: (1) the structural deviation from MAP (how  
628 many variables have different values, i.e., the Hamming distance between  
629 the MFE and MAP explanations), (2) the rank  $k$  of the MFE explanation,  
630 indicating that the MFE explanation is the  $k$ -th MAP instead of the most  
631 probable explanation, and (3) the ratio of the MFE probability and the MAP  
632 probability, indicating the proportion of probability mass that was allocated  
633 to the MFE explanation.

634 Furthermore, we estimated how often the MFE was picked relative to  
635 other explanations, indicating how likely it is that a singleton sample over  
636 the irrelevant variables would yield this particular explanation. This yields a



637 measure on how many samples are needed to arrive at a confident decision.  
 638 Lastly, we estimated the likelihood of picking the MAP explanation and one  
 639 of the five most probable explanations using a single sample. This indicates  
 640 how likely it is that an arbitrary singleton sample will yield an explanation  
 641 with the maximum, respectively a relatively high, posterior probability.

642 The results are summarized in Table 2 and Figure 6. The scatter plots  
 643 in Figure 6 illustrate the spread of the errors along different networks. In  
 644 general one can conclude that MFE explanations are reasonably close to the  
 645 MAP explanations, when there is marginalization over those variables that  
 646 are ‘sufficiently relevant’. From the results it follows that including the five  
 647 most relevant variables gives fairly good results, and that including variables  
 648 that have a relevance of less than 0.25 does not significantly improve the  
 649 average MFE results. Including no relevant variables at all (i.e., computing  
 650 the Most Simple Explanation [35]) gives considerably worse results, however.

Cond.	$I^+$ size	ratio	rank	dist.	% MFE	% MAP	% 5-MAP
None	0	0.66	25.90	2.05	0.08	0.03	0.14
Best 5	5	0.82	10.73	1.30	0.13	0.08	0.27
> 0.50	11.32	0.87	5.36	0.87	0.25	0.17	0.46
> 0.25	14.93	0.91	4.59	0.79	0.38	0.25	0.58
> 0.10	15.79	0.91	5.56	0.81	0.39	0.25	0.60
> 0.05	15.99	0.91	6.09	0.75	0.41	0.27	0.60
> 0.00	16.35	0.92	4.12	0.75	0.41	0.26	0.61

Table 2: Overview of simulation results. In this simulation the partition between relevant and irrelevant variables was varied and ranged from ‘none’ (all variables are irrelevant), ‘best 5’ (the five variables with the highest relevance are deemed relevant, to a relevance threshold between 0.50 and 0.00, yielding an average  $I^+$  size between 11.32 and 16.35).

### 651 4.3. Number of Samples

652 As shown in Section 3.3, approximating the MFE (i.e., finding the ex-  
 653 planation which is very likely the MFE) can be done by sampling, where  
 654 the number of samples needed to guarantee a particular confidence level is  
 655 related to the decisiveness of the network. When decisiveness is low, and con-  
 656 sequently the MFE distribution is flat (many competing explanations, none  
 657 of which has a high probability of being the most probable explanation for a  
 658 random joint value assignment to the irrelevant intermediate variables), we

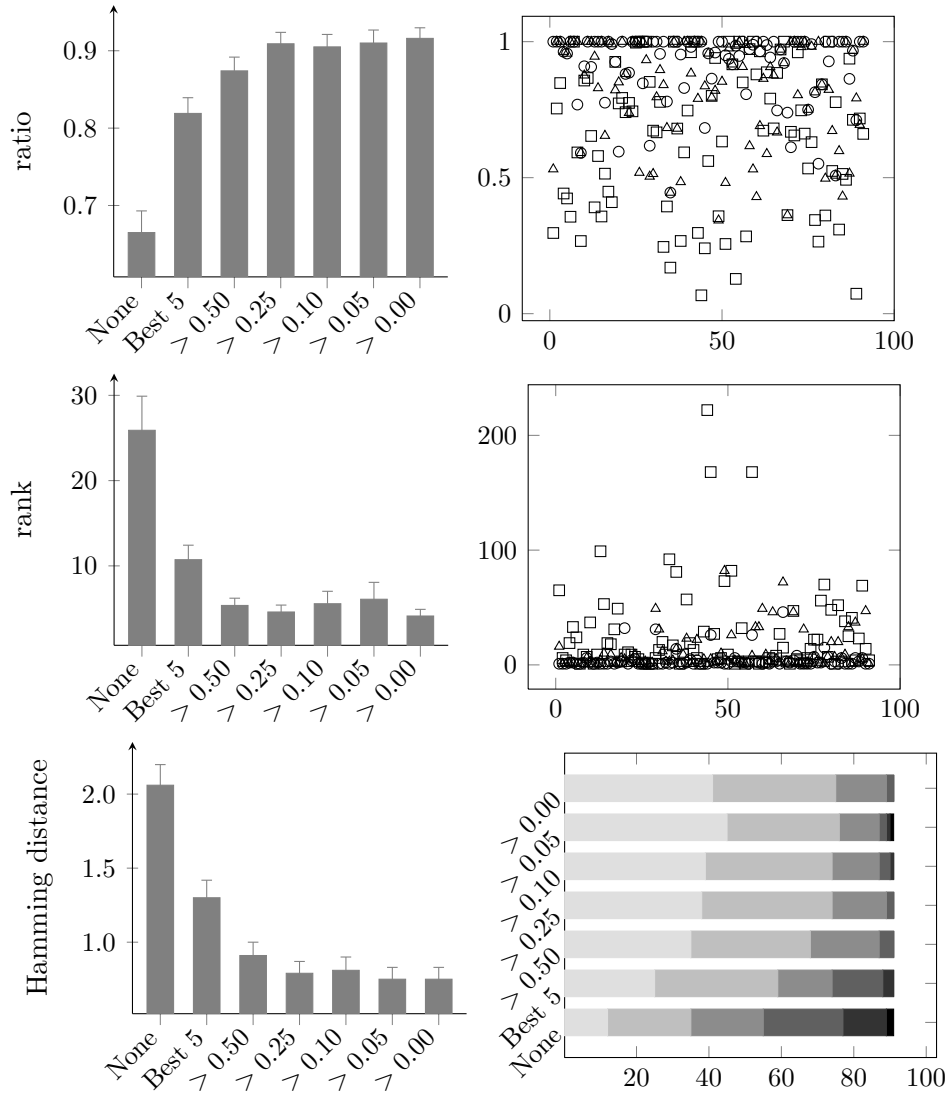


Figure 6: On the left: Three error indicators of MFE versus MAP explanations: The ratio between their probabilities, rank of the MFE explanation, and Hamming distance between MFE and MAP for various  $I^+/I^-$  settings. On the right: Scatter plots of ratio and rank, and stacked box plot for Hamming distance. In the scatter plots, results of all random networks are shown, for the conditions where all variables are irrelevant ('None', square), the five variables with the highest relevancy were deemed relevant ('best 5', triangle) and where all variables with non-zero relevancy were relevant ('> 0.00', circle). The stacked box plot illustrates the distribution of the Hamming distance between MFE and MAP explanation, where darker colors indicate a higher Hamming distance. Error bars indicate standard error of the mean.

659 need much more samples to make confident decisions. This is illustrated by  
 660 the following figures. In Figure 7 we see a typical result for a random network  
 661 which is highly skewed towards a singleton explanation, and in Figure 8 the  
 662 results of a random network with a low decisiveness are shown.

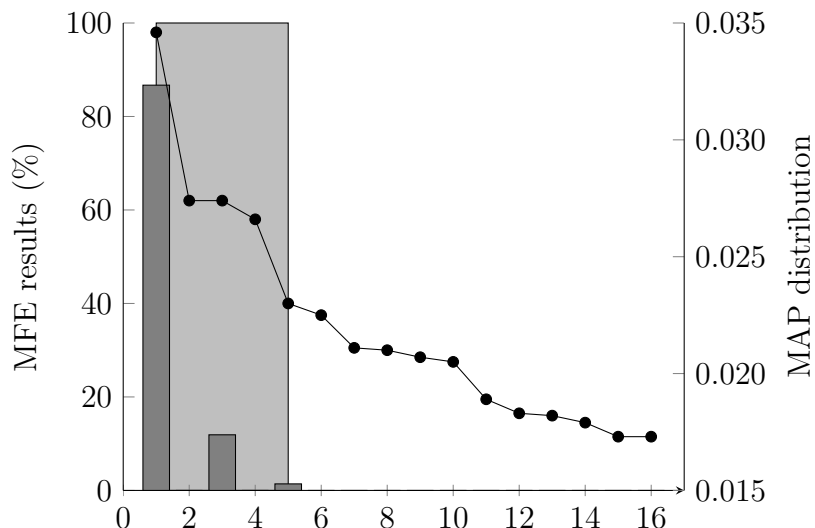


Figure 7: This plot shows part of the MAP distribution and MFE results using 1000 samples for a random network (#93) with a very steep distribution of the MFE explanations. This network is strongly skewed towards the most probable explanation which is picked in 83% of the samples, so that an arbitrary singleton sample is quite likely to be the MFE; we can be guaranteed to obtain the most frugal explanation with 95% confidence by generating thirteen samples and decide which explanation is most often picked. Even a single sample is guaranteed to correspond to one of the five most probable examples.

---

663 However, even when there is no explanation which stands out, the sam-  
 664 pling algorithm can still give good results. In Figure 9 we show a typical  
 665 result when there are *a few* competing explanations that all have a relatively  
 666 high probability. While it may take many samples to decide on which of  
 667 them is the MFE, we still can be quite sure that a singleton sample of the  
 668 irrelevant intermediate variables would yield *one of them* as the most prob-  
 669 able explanation; here, sampling seems like a reasonable strategy to obtain  
 670 an explanation that is likely to have a reasonably high probability.

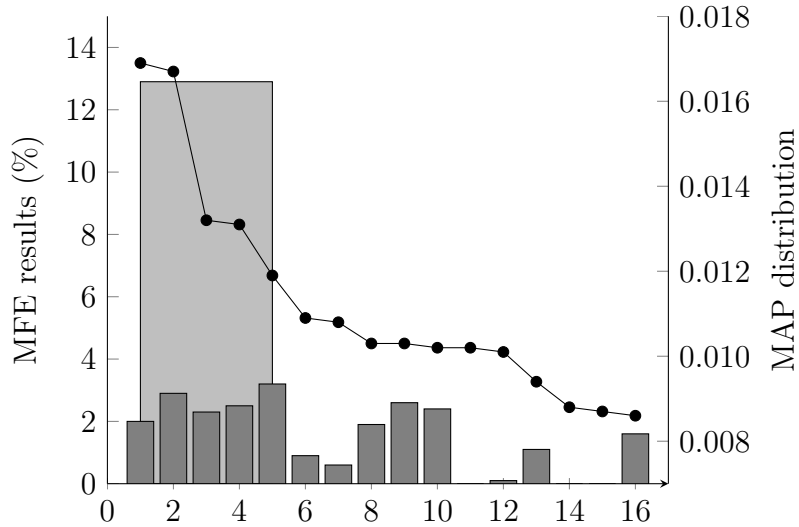


Figure 8: This plot shows part of the MAP distribution and MFE results using 1000 samples for a random network (#89) with a very flat distribution of the MFE explanations. No explanation really stands out; the most frugal explanation being picked in just over 3% of the samples. In this network, that is not at all skewed towards any particular explanation, an arbitrary sample can have a low posterior probability, and we will need a massive number of samples to decide with reasonable confidence about which explanation is the MFE.

671 *4.4. Other parameters*

672 Obviously, the  $\mathbf{I}^+/\mathbf{I}^-$  partition influences the quality of the MFE solution  
673 in terms of the three error measures introduced in Section 4.2. We also in-  
674 vestigated whether the size of the hypothesis space, the number of relevant  
675 variables, or the probability of the most probable explanation influences this  
676 quality. First we observe that these parameters are not independent. There  
677 is a strong negative correlation (-.65) between the size of the explanation set  
678 and the probability of the most probable explanation. This can be explained  
679 by the random nature of the networks and the probability distribution they  
680 capture: on average, if there are more candidate explanations in the explana-  
681 tion set, the average probability of each of them is lower, and so it is expected  
682 that the average probability of the most probable explanation is also lower.  
683 The results of the correlation analysis are shown in Table 3, and can be  
684 summarized as follows. Neither explanation set size, intrinsic relevance, or

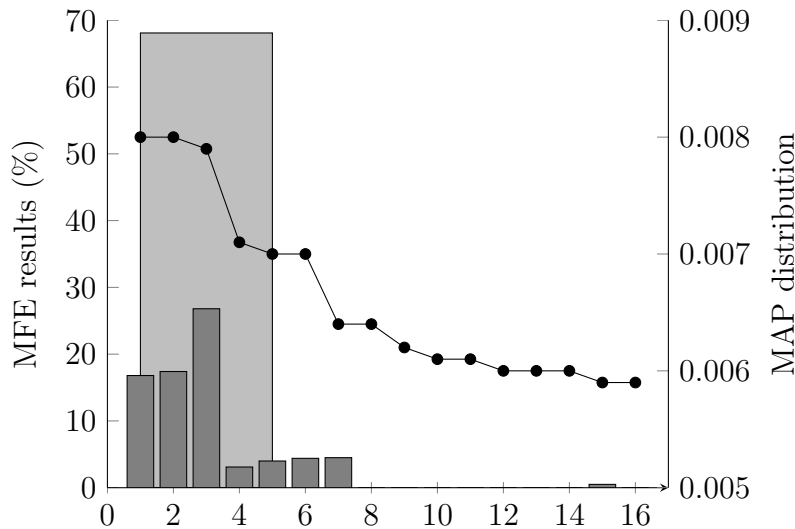


Figure 9: This plot shows part of the MAP distribution and MFE results using 1000 samples for a random network (#70) where three explanations are often picked as the most probable, and quite some samples are needed to decide on the most frugal explanation with confidence. However, since one of these three (almost equally probable) most probable explanations is picked in 61% of the samples, we can expect that few samples, possibly just a singleton sample, may return a quite good explanation.

---

685 probability of the most probable explanation (MPE) correlates with the ratio  
 686 between probability of MPE and probability of MFE. There is a weak corre-  
 687 lation between explanation set size and rank, and a weak negative correlation  
 688 between probability of MPE and rank: the bigger the explanation size, the  
 689 larger the average rank  $k$ . Neither explanation set size, intrinsic relevance, or  
 690 probability of MPE correlates (or correlates only very weakly) with distance  
 691 errors.

#### 692 4.5. Wrong judgments

693 Obviously, taking more intermediate variables into account (i.e., consider-  
 694 ing more variables to be relevant) helps to obtain better results; still, we can  
 695 make reasonable good inferences using only the five most relevant interme-  
 696 diate variables. But what if ones subjective assessment of what is relevant  
 697 does not match the intrinsic relevance of these variables? Figure 10 illus-  
 698 trates what typically happens when there is a mismatch between intrinsic

Cond.	explanation set size			intrinsic relevance			probability of MPE		
	ratio	rank	dist.	ratio	rank	dist.	ratio	rank	dist.
MSE	-.01	.15	.15	-.09	.02	.18	-.11	-.23*	-.20
Best 5	-.16	.22*	.27*	.13	.18	.07	-.15	-.35**	-.40**
> 0.50	.08	.12	.02	-.11	.01	.18	-.04	-.17	-.16
> 0.25	-.09	.24*	.12	-.11	.05	-.02	.06	-.22*	-.12
> 0.10	-.10	.26*	.21*	-.08	.01	-.06	.05	-.17	-.18
> 0.05	-.08	.22*	.10	-.08	.01	-.02	.02	-.13	-.11
> 0.00	.06	.17	.03	-.20	.11	.01	-.01	-.19	-.03

Table 3: Overview of correlations (Pearson’s  $r$ ) with significance levels. \* indicates significance at the  $p < .05$  level, \*\* indicates significance at the  $p < .01$  level

699 and estimated relevance. Here we plotted the results of the  $> 0.00$  (top left)  
700 and Best 5 (bottom right) conditions, as well as some conditions in which  
701 there is a mismatch between intrinsic and expected relevance. In particu-  
702 lar, we omitted the two (top right), five (middle left), ten (middle right),  
703 respectively fifteen (bottom left) most relevant variables.

704 This example illustrates a graceful degradation of the results, especially  
705 when the cumulative results of the five most probable joint value assignments  
706 are compared. Observe that including the twenty-five *least* relevant variables  
707 leads to comparable results as including the five *most* relevant variables.  
708 Clearly, it helps to know what is relevant, yet there is margin for error.

#### 709 4.6. Discussion

710 The simulation results, as illustrated by Table 2 and Figure 6, clearly  
711 show that MFE ‘tracks truth’ quite well, even when only part of the relevant  
712 intermediate variables are taken into account. However, when more interme-  
713 diate variables are marginalized over, we can be more confident of the results.  
714 In these cases the distribution of explanations typically is narrower and it is  
715 more likely that a majority vote using few samples, or even a singleton sam-  
716 ple, results in an explanation that is close to the most probable explanation.  
717 The simulations also indicate that indeed the probability distribution must  
718 be skewed towards one or a few explanations for obtaining acceptable results  
719 with few samples - and that indeed many distributions *are* skewed, even in  
720 completely random networks.

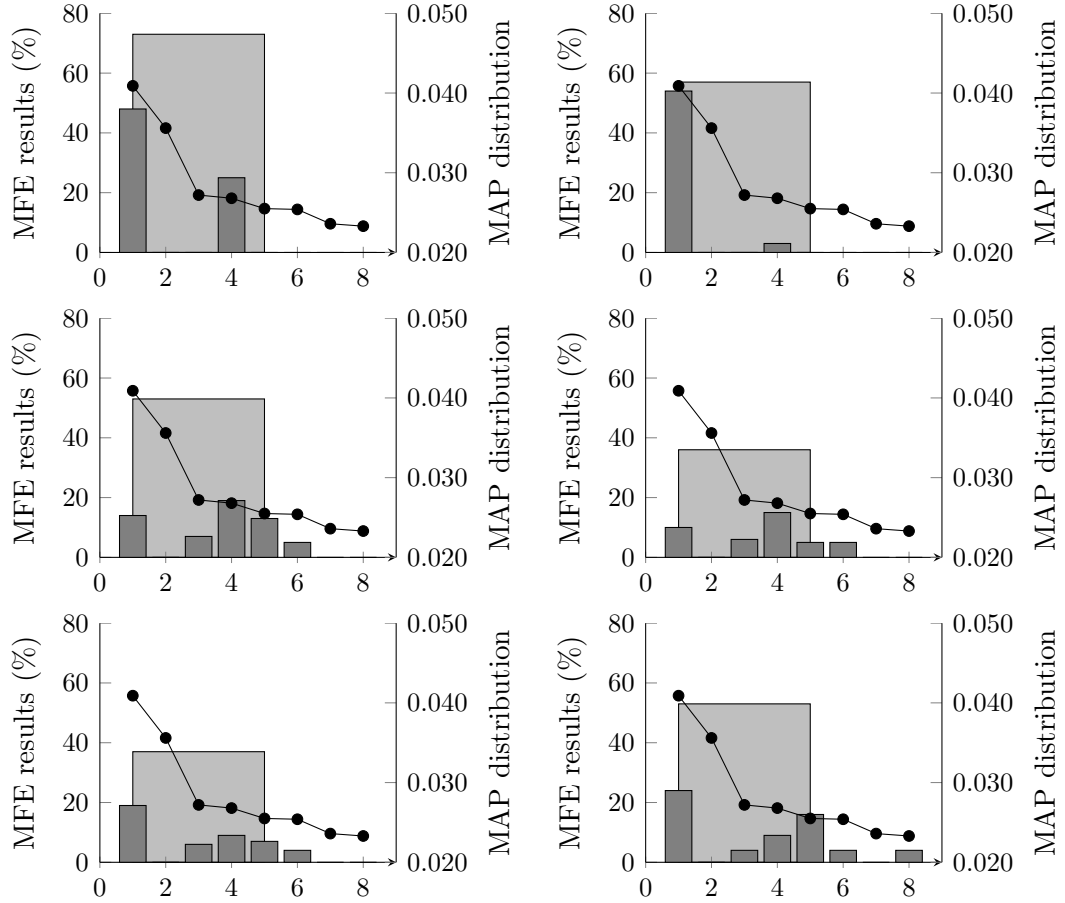


Figure 10: This plot shows part of the MAP distribution and MFE results of a random network (#78) with different partitions of the intermediate variables, where the subjective assessment that yields the partition may not match the actual relevance of the variables. Shown are the results when all variables with non-zero relevancy are deemed relevant (top left, 19 variables in  $\mathbf{I}^+$ ), all *but* the two most relevant variables (top right, 28 variables in  $\mathbf{I}^+$ ), all *but* the five most relevant variables (middle left, 25 variables in  $\mathbf{I}^+$ ), all *but* the ten most relevant variables (middle right, 20 variables in  $\mathbf{I}^+$ ), only the fifteen *least* relevant variables (bottom left, 15 variables in  $\mathbf{I}^+$ ), and only the five *most* relevant variables (bottom right, 5 variables in  $\mathbf{I}^+$ ).

721 **5. Conclusion**

722 In this paper we proposed Most Frugal Explanation (MFE) as a tractable  
723 heuristic alternative to (approximate) MAP for deciding upon the best ex-  
724 planation in Bayesian networks. While the MFE problem is intractable in  
725 general—its decision variant is  $\text{NP}^{\text{PPP}}$ -complete, and thus even harder than  
726 the  $\text{NP}^{\text{PP}}$ -complete MAP problem [51], the  $\text{PP}^{\text{PP}}$ -complete Same-Decision  
727 Probability problem [9], or the  $\text{P}^{\text{PPP}}$ -complete  $k$ -th MAP problem [41]—it  
728 can be tractably approximated under situational constraints that are ar-  
729 guably more realistic in large real-world applications than the constraints  
730 that are needed to render MAP (fixed-parameter) tractable. Notably, the  
731  $\{c, tw, 1 - p\}$ -fixed-parameter tractable algorithm for MAP [4] has a running  
732 time with  $\frac{\log p}{\log 1-p}$  in the exponent. In the random networks used in the simu-  
733 lations, the *average* probability of the most probable explanation was 0.0245,  
734 which would yield an unpractical exponent of  $\frac{\log 0.0245}{\log 0.9755} \approx 150$ . In contrast,  
735 even when only half of the total set of intermediate variables are considered  
736 as relevant, for an arbitrary sample over the rest of the intermediate variables  
737 we will find the MFE in about 40% of the cases, and an explanation that is  
738 one of the five best in about 60% of the cases.

739 In future work we wish to investigate the possible explanatory power  
740 of MFE in cognitive science. In recent years it has been proposed that  
741 human cognizers make decisions using (Bayesian) sampling [31, 57, 62] and  
742 approximate Bayesian inferences using exemplars [54]; studies show that we  
743 have a hard time solving problems with many relevant aspects [20]. The  
744 parameterized complexity results of the MFE framework may theoretically  
745 explain why such approaches work fine in practice and under what conditions  
746 the limits of our cognitive capacities are reached.

747 **6. Acknowledgments**

748 The author wishes to thank the members of the *Theoretical Cognitive*  
749 *Science* group at the Donders Center for Cognition for useful discussions and  
750 comments on earlier versions of this paper, and the anonymous reviewers  
751 that gave valuable suggestions for improvement. He is in particular indebted  
752 to Todd Wareham for suggesting the term “Most Frugal Explanations” to  
753 denote the problem of finding an explanation for observations without taking  
754 care of everything that is only marginally relevant. A previous shorter version  
755 of this paper appeared in the Benelux Conference on AI [38].



756 **7. Vitae**

757 Johan Kwisthout is a postdoctoral researcher at the Donders Institute for  
758 Brain, Cognition and Behaviour, and lecturer in the Department of Artificial  
759 Intelligence, at the Raboud University Nijmegen. His research interests are  
760 reasoning under uncertainty, computational complexity, and the application  
761 of both fields of study in (theoretical) cognitive (neuro-)science. In particular,  
762 he is interested in topics related to (Bayesian) abduction and relevance, both  
763 from a philosophical, cognitive, and computational perspective.

764 **Appendix: Computing relevance is NP-hard**

765 In Definition 4 we formally defined the intrinsic relevance of an interme-  
766 diate variable as a measure indicating the sensitivity of explanations to its  
767 value. We here show that computing the intrinsic relevance of such a variable  
768 is NP-hard. The decision problem used in this proof is defined as follows:

769 **INTRINSIC RELEVANCE**

770 **Instance:** A Bayesian network  $\mathcal{B} = (\mathbf{G}_{\mathcal{B}}, \text{Pr})$ , where  $\mathbf{V}$  is partitioned into  
771 evidence variables  $\mathbf{E}$  with joint value assignment  $\mathbf{e}$ , explanation variables  
772  $\mathbf{H}$ , and intermediate variables  $\mathbf{I}$ , and a designated variable  $I \in \mathbf{I}$ .

773 **Question:** Is the *intrinsic relevance*  $\mathcal{R}(I) > 0$ ?

774 We reduce from the following NP-complete decision problem [37]:

775 **ISA-RELEVANT VARIABLE**

776 **Instance:** A Boolean formula  $\phi$  with  $n$  variables, describing the  
777 characteristic function  $\mathbf{1}_{\phi} : \{\text{FALSE}, \text{TRUE}\}^n \rightarrow \{1, 0\}$ , designated variable  
778  $x_r \in \phi$ .

779 **Question:** Is  $x_r$  a relevant variable in  $\phi$ , that is, is

780  $\mathbf{1}_{\phi}(x_r = \text{TRUE}) \neq \mathbf{1}_{\phi}(x_r = \text{FALSE})$ ?

781 Here, the characteristic function  $\mathbf{1}_{\phi}$  of a Boolean formula  $\phi$  maps truth  
782 assignments to  $\phi$  to  $\{0, 1\}$ , such that  $\mathbf{1}_{\phi}(x) = 1$  if and only if  $x$  denotes a  
783 satisfying truth assignment to  $\phi$ , and 0 otherwise. We will use the formula  
784  $\phi_{\text{ex}} = \neg(x_1 \vee x_2) \wedge x_3$  as a running example, where  $x_3$  is the variable of  
785 interest. Note that  $x_3$  is relevant, since for  $x_1 = x_2 = \text{FALSE}$ ,  $\mathbf{1}_{\phi}(x_3 =$   
786  $\text{TRUE}) \neq \mathbf{1}_{\phi}(x_3 = \text{FALSE})$ .

787 We construct a Bayesian network  $\mathcal{B}_\phi$  from  $\phi$  as follows. For each propo-  
788 sitional variable  $x_i \in \phi$  we add a binary stochastic variable  $X_i \in \mathcal{B}_\phi$  with  
789 uniformly distributed values TRUE and FALSE. We add an additional binary  
790 variable  $X_r^T$  with observed value TRUE. For each logical operator  $o_j$  in  $\phi$ ,  
791 we add *two* binary stochastic variables  $O_j$  and  $O_j^T$  in  $\mathcal{B}_\phi$ . The parents of  
792 the variables  $O_j$  are the variables  $O_k$  that represent the sub-formulas bound  
793 by  $O_j$ ; in case such a sub-formula is a literal, the corresponding parent is a  
794 variable  $X_i$ . In contrast, the parents of the variables  $O_j^T$  are the variables  
795  $O_k^T$  (for sub-formula),  $X_i$  (for literals *except*  $x_r$ ), respectively  $X_r^T$  (for the  
796 literal  $x_r$ ). The variables corresponding with the top-level operator in  $\phi$  are  
797 denoted with  $V_\phi$ , respectively  $V_\phi^T$ .

798 Furthermore, an additional binary variable  $C$  is introduced in  $\mathcal{B}_\phi$ , acting  
799 as ‘comparator’ variable.  $C$  has both  $V_\phi$  and  $V_\phi^T$  as parents and condi-  
800 tional probability  $\Pr(C = \text{TRUE} \mid V_\phi, V_\phi^T) = 1$  if  $V_\phi \neq V_\phi^T$  and  $\Pr(C =$   
801  $\text{TRUE} \mid V_\phi, V_\phi^T) = 0$  if  $V_\phi = V_\phi^T$ . An example of this construction is given in  
802 Figure 11 for the formula  $\phi_{\text{ex}}$ . We set  $\mathbf{H} = C$ ,  $\mathbf{E} = X_r^T$ , and  $I = X_r$ .

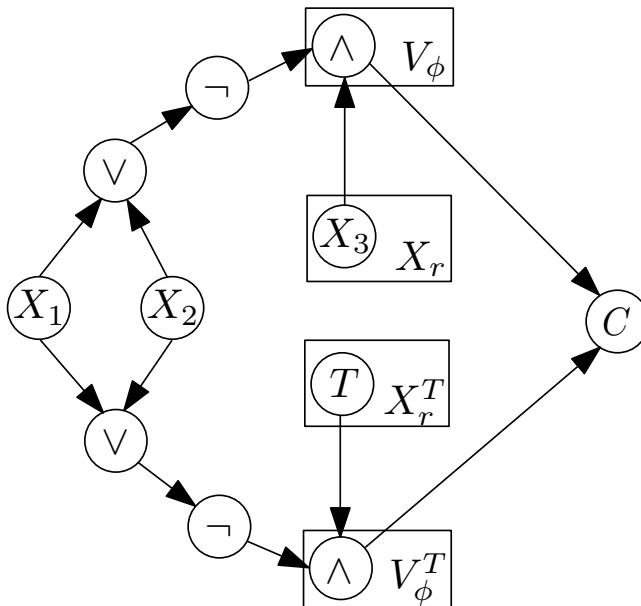


Figure 11: Example of the construction of  $\mathcal{B}_{\phi_{\text{ex}}}$  for the formula  $\phi_{\text{ex}} = \neg(x_1 \vee x_2) \wedge x_3$

803 **Theorem 7.** INTRINSIC RELEVANCE *is NP-complete.*

804 *Proof.* Membership in NP follows from the following polynomial-time verify-  
 805 ing algorithm for *yes*-instances: given a suitable joint value assignment  $\mathbf{i}$  to  $\mathbf{I} \setminus$   
 806  $\{I\}$  and assignments  $i_1, i_2$  to  $I$ , we can easily check that  $\operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{h}, \mathbf{e}, \mathbf{i}, I =$   
 807  $i_1) \neq \operatorname{argmax}_{\mathbf{h}} \Pr(\mathbf{h}, \mathbf{e}, \mathbf{i}, I = i_2)$ , and thus that  $\mathcal{R}(I) > 0$ .

808 To prove NP-hardness, we reduce ISA-RELEVANT VARIABLE to INTRIN-  
 809 SIC RELEVANCE. Let  $(\phi, x_r)$  be an instance of ISA-RELEVANT VARIABLE.  
 810 From  $(\phi, x_r)$ , we construct  $(\mathcal{B}_\phi, I)$  as described above. If  $(\phi, x_r)$  is a *yes*-  
 811 instance of ISA-RELEVANT VARIABLE, then the characteristic function  $\mathbf{1}_\phi$   
 812 is not identical for  $x_r = \text{FALSE}$  and  $x_r = \text{TRUE}$ . In other words, there  
 813 is at least one truth assignment  $\mathbf{t}$  to the variables in  $\phi \setminus \{x_r\}$  such that  
 814 either  $\mathbf{t} \cup \{x_r = \text{TRUE}\}$  is satisfying  $\phi$  and  $\mathbf{t} \cup \{x_r = \text{FALSE}\}$  is not sat-  
 815 isfying  $\phi$ , or vice versa. Let  $\mathbf{i}$  be the joint value assignment to  $\mathbf{I} \setminus \{X_r\}$   
 816 that corresponds to the truth assignment  $\mathbf{t}$ , and in addition fixes the val-  
 817 ues of the operator variables  $O_j^T$  and  $O_j$  according to their (determinis-  
 818 tic) conditional probability tables. Now, we have that for the truth as-  
 819 signment  $X_r = \text{TRUE}$ ,  $\Pr(C = \text{TRUE} \mid \mathbf{i}, X_r^T = \text{TRUE}) = 1$  and thus  
 820  $\operatorname{argmax}_c \Pr(C = c, \mathbf{i}, X_r = \text{FALSE}) = \text{TRUE}$ . By definition, we have that for  
 821 the truth assignment  $X_r = \text{FALSE}$  that  $\Pr(C = \text{TRUE} \mid \mathbf{i}, X_r^T = \text{FALSE}) = 0$   
 822 and thus  $\operatorname{argmax}_c \Pr(C = c, \mathbf{i}, X_r = \text{FALSE}) = \text{FALSE}$ . Hence, the intrinsic  
 823 relevance  $\mathcal{R}(X_r) > 0$  and thus  $(\mathcal{B}_\phi, I)$  is a *yes*-instance of INTRINSIC RELE-  
 824 VANCE.

825 Now we assume that  $\mathcal{R}(I) > 0$ , implying that there is at least one  
 826 truth assignment  $\mathbf{i}$  to  $\mathbf{I} \setminus \{X_r\}$  such that  $\Pr(C = \text{TRUE} \mid \mathbf{i}, X_r^T = \text{FALSE}) \neq$   
 827  $\operatorname{argmax}_c \Pr(C = c, \mathbf{i}, X_r = \text{FALSE})$  where the joint value assignment to the  
 828 operator variables  $O_j^T$  and  $O_j$  matches the deterministic conditional prob-  
 829 abilities imposed by the joint value assignment to the variables  $X_i$ . This  
 830 implies that the characteristic function  $\mathbf{1}_\phi$  is not identical for  $x_r = \text{FALSE}$   
 831 and  $x_r = \text{TRUE}$ , hence, that  $(\phi, x_r)$  is a *yes*-instance of ISA-RELEVANT  
 832 VARIABLE.

833 As the reduction can be done in polynomial time, this proves that IN-  
 834 TRINSIC RELEVANCE is NP-complete.  $\square$

## 835 References

- 836 [1] Abdelbar, A.M., Hedetniemi, S.M., 1998. Approximating MAPs for  
 837 belief networks is NP-hard and other theorems. Artificial Intelligence  
 838 102, 21–38.

- 839 [2] Beinlich, I., Suermondt, G., Chavez, R., Cooper, G., 1989. The ALARM  
840 monitoring system: A case study with two probabilistic inference tech-  
841 niques for belief networks, in: Hunter, J., Cookson, J., Wyatt, J. (Eds.),  
842 Proceedings of the Second European Conference on AI and Medicine,  
843 Springer-Verlag. pp. 247–256.
- 844 [3] Bodlaender, H.L., 2006. Treewidth: characterizations, applications, and  
845 computations, in: Proceedings of the 32nd International Workshop on  
846 Graph-Theoretic Concepts in Computer Science, pp. 1–14.
- 847 [4] Bodlaender, H.L., van den Eijkhof, F., van der Gaag, L.C., 2002. On  
848 the complexity of the MPA problem in probabilistic networks, in: van  
849 Harmelen, F. (Ed.), Proceedings of the Fifteenth European Conference  
850 on Artificial Intelligence, IOS Press, Amsterdam. pp. 675–679.
- 851 [5] Bovens, L., Olsson, E.J., 2000. Coherentism, reliability and Bayesian  
852 networks. *Mind* 109, 686–719.
- 853 [6] Chajewska, U., Halpern, J., 1997. Defining explanation in probabilistic  
854 systems, in: Geiger, D., Shenoy, P. (Eds.), Proceedings of the Thirteenth  
855 Conference on Uncertainty in Artificial Intelligence, Morgan Kaufmann,  
856 San Francisco, CA. pp. 62–71.
- 857 [7] Cofiño, A.S., Cano, R., Sordo, C., Gutiérrez, J.M., 2002. Bayesian net-  
858 works for probabilistic weather prediction, in: van Harmelen, F. (Ed.),  
859 Proceedings of the Fifteenth European Conference on Artificial Intelli-  
860 gence, IOS Press, Amsterdam. pp. 695–699.
- 861 [8] Darwiche, A., 2009. Modeling and Reasoning with Bayesian Networks.  
862 CU Press, Cambridge, UK.
- 863 [9] Darwiche, A., Choi, A., 2010. Same-decision probability: a confidence  
864 measure for threshold-based decisions under noisy sensors, in: Myllymki,  
865 P., Roos, T., Jaakkola, T. (Eds.), Proceedings of the Fifth European  
866 Workshop on Probabilistic Graphical Models, p. 113120.
- 867 [10] Demirer, R., Mau, R., Shenoy, C., 2006. Bayesian Networks: A decision  
868 tool to improve portfolio risk analysis. *Journal of Applied Finance* 16,  
869 106–119.

- 870 [11] Dey, S., Stori, J.A., 2005. A Bayesian network approach to root cause  
871 diagnosis of process variations. *International Journal of Machine Tools*  
872 *and Manufacture* 45, 75–91.
- 873 [12] Downey, R.G., Fellows, M.R., 1999. *Parameterized Complexity*. Springer  
874 Verlag, Berlin.
- 875 [13] Druzdzal, M., 1994. Some properties of joint probability distributions,  
876 in: de Mantaras, R.L., Poole, D. (Eds.), *Proceedings of the Tenth Con-*  
877 *ference on Uncertainty in Artificial Intelligence*, San Francisco, Ca: Mor-  
878 gan Kaufmann Publishers. pp. 187–194.
- 879 [14] Druzdzal, M.J., Suermondt, H.J., 1994. Relevance in probabilistic mod-  
880 els: “backyards” in a “small world”, in: *Working notes of the AAAI–*  
881 *1994 Fall Symposium Series: Relevance*, pp. 60–63.
- 882 [15] Fitelson, B., 2003. A probabilistic theory of coherence. *Analysis* 63,  
883 194–199.
- 884 [16] Flum, G., Grohe, M., 2006. *Parameterized Complexity Theory*.  
885 Springer, Berlin.
- 886 [17] Fodor, J.A., 1983. *The Modularity of Mind*. Cambridge, MA: MIT  
887 Press.
- 888 [18] Fodor, J.A., 1987. Modules, frames, fridgeons, sleeping dogs, and the  
889 music of the spheres, in: Pylyshyn, Z.W. (Ed.), *The Robot’s Dilemma:*  
890 *The Frame Problem in Artificial Intelligence*. Ablex Publishing, pp. 139–  
891 150.
- 892 [19] Fodor, J.A., Lepore, E., 1992. *Holism: A shopper’s guide*. volume 16.  
893 Blackwell Oxford.
- 894 [20] Funke, J., 1991. Solving complex problems: Exploration and control of  
895 complex social systems, in: Sternberg, R.J., Frensch, P.A. (Eds.), *Com-*  
896 *plex Problem Solving: Principles and Mechanisms*. Lawrence Erlbaum  
897 Associates, pp. 185–222.
- 898 [21] van der Gaag, L.C., Renooij, S., Witteman, C.L.M., Aleman, B.M.P.,  
899 Taal, B.G., 2002. Probabilities for a probabilistic network: a case study  
900 in oesophageal cancer. *Artificial Intelligence in Medicine* 25, 123–148.

- 901 [22] Garey, M.R., Johnson, D.S., 1979. Computers and Intractability. A  
902 Guide to the Theory of NP-Completeness. W. H. Freeman and Co., San  
903 Francisco, CA.
- 904 [23] Geenen, P.L., Elbers, A.R.W., van der Gaag, L.C., van der Loeffen,  
905 W.L.A., 2006. Development of a probabilistic network for clinical detec-  
906 tion of classical swine fever, in: Proceedings of the Eleventh Symposium  
907 of the International Society for Veterinary Epidemiology and Economics,  
908 pp. 667–669.
- 909 [24] Geiger, D., Verma, T., Pearl, J., 1990. Identifying independence in  
910 Bayesian networks. *Networks* 20, 507–534.
- 911 [25] Gemela, J., 2001. Financial analysis using Bayesian networks. *Applied  
912 Stochastic Models in Business and Industry* 17, 57–67.
- 913 [26] Glass, D.H., 2007. Coherence measures and inference to the best expla-  
914 nation. *Synthese* 157, 275–296.
- 915 [27] Glass, D.H., 2009. Inference to the best explanation: a comparison of  
916 approaches, in: Bishop, M. (Ed.), *Proceedings of the Second Symposium  
917 on Computing and Philosophy*, The Society for the Study of Artificial  
918 Intelligence and the Simulation of Behaviour. pp. 22–27.
- 919 [28] Glass, D.H., 2012. Inference to the best explanation: does it track truth?  
920 *Synthese* 185, 411–427.
- 921 [29] Hempel, C.G., 1965. *Aspects of Scientific Explanation*. Free Press, New  
922 York.
- 923 [30] Jaynes, E., 2003. *Probability Theory: The Logic of Science*. Cambridge  
924 University Press.
- 925 [31] Johnson-Laird, P.N., Legrenzi, P., Girotto, V., Legrenzi, M.S., Caverni,  
926 J., 1999. Naive probability: A mental model theory of extensional rea-  
927 soning. *Psychological Review* 106, 62–88.
- 928 [32] Kennett, R.J., Korb, K.B., Nicholson, A.E., 2001. Seabreeze prediction  
929 using Bayesian networks, in: Cheung, D.W.L., Williams, G., Li, Q.  
930 (Eds.), *Proceedings of the Fifth Pacific-Asia Conference on Advances  
931 in Knowledge Discovery and Data Mining*, Springer Verlag, Berlin. pp.  
932 148–153.

- 933 [33] Kragt, M.E., Newhama, L.T.H., Jakemana, A.J., 2009. A Bayesian  
934 network approach to integrating economic and biophysical modelling,  
935 in: Anderssen, R., Braddock, R., Newham, L. (Eds.), Proceedings of  
936 the Eighteenth World IMACS / MODSIM Congress on Modelling and  
937 Simulation, pp. 2377–2383.
- 938 [34] Kwisthout, J., 2009. The Computational Complexity of Probabilistic  
939 Networks. Ph.D. thesis. Faculty of Science, Utrecht University, The  
940 Netherlands.
- 941 [35] Kwisthout, J., 2010. Two new notions of abduction in Bayesian net-  
942 works, in: et al., P.B. (Ed.), Proceedings of the 22nd Benelux Conference  
943 on Artificial Intelligence (BNAIC’10), pp. 82–89.
- 944 [36] Kwisthout, J., 2011. Most probable explanations in Bayesian networks:  
945 Complexity and tractability. *International Journal of Approximate Rea-*  
946 *soning* 52, 1452 – 1469.
- 947 [37] Kwisthout, J., 2012. Relevancy in problem solving: A computational  
948 framework. *The Journal of Problem Solving* 5, 17 – 32.
- 949 [38] Kwisthout, J., 2013a. Most frugal explanations: Occam’s razor applied  
950 to Bayesian abduction, in: Hindriks, K., de Weerd, M., van Riemsdijk,  
951 B., Warnier, M. (Eds.), Proceedings of the 25th Benelux Conference on  
952 AI (BNAIC’13), pp. 96–103.
- 953 [39] Kwisthout, J., 2013b. Most inforbable explanations: Finding expla-  
954 nations in Bayesian networks that are both probable and informative,  
955 in: van der Gaag, L. (Ed.), Proceedings of the Twelfth European Con-  
956 ference on Symbolic and Quantitative Approaches to Reasoning with  
957 Uncertainty, pp. 328–339.
- 958 [40] Kwisthout, J., 2013c. Structure approximation of most probable expla-  
959 nations in Bayesian networks, in: van der Gaag, L. (Ed.), Proceedings  
960 of the Twelfth European Conference on Symbolic and Quantitative Ap-  
961 proaches to Reasoning with Uncertainty, pp. 340–351.
- 962 [41] Kwisthout, J., Bodlaender, H.L., van der Gaag, L.C., 2011a. The com-  
963 plexity of finding  $k$ th most probable explanations in probabilistic net-  
964 works, in: Cerná, I., Gyimóthy, T., Hromkovic, J., Jefferey, K., Královic,

- 965 R., Vukolic, M., Wolf, S. (Eds.), Proceedings of the 37th International  
966 Conference on Current Trends in Theory and Practice of Computer Sci-  
967 ence, pp. 356–367.
- 968 [42] Kwisthout, J., van Rooij, I., 2013. Bridging the gap between theory and  
969 practice of approximate Bayesian inference. *Cognitive Systems Research*  
970 24, 2–8.
- 971 [43] Kwisthout, J., Wareham, T., van Rooij, I., 2011b. Bayesian intractabil-  
972 ity is not an ailment approximation can cure. *Cognitive Science* 35,  
973 779–1007.
- 974 [44] Lipton, P., 2004. *Inference to the Best Explanation*. London, UK:  
975 Routledge. 2nd edition.
- 976 [45] Lucas, P.J.F., de Bruijn, N., Schurink, K., Hoepelman, A., 2000. A prob-  
977 abilistic and decision-theoretic approach to the management of infectious  
978 disease at the ICU. *Artificial Intelligence in Medicine* 3, 251–279.
- 979 [46] Murphy, K., 2001. *The Bayes Net Toolbox for MATLAB*. *Computing*  
980 *Science and Statistics* 33, 2001.
- 981 [47] Neapolitan, R.E., 1990. *Probabilistic Reasoning in Expert Systems.*  
982 *Theory and Algorithms*. Wiley/Interscience, New York, NY.
- 983 [48] Nedeveschi, S., Sandhu, J.S., Pal, J., Fonseca, R., Toyama, K., 2006.  
984 Bayesian networks: an exploratory tool for understanding ICT adoption,  
985 in: Toyama, K. (Ed.), *Proceedings of the International Conference on*  
986 *Information and Communication Technologies and Development*, pp.  
987 277–284.
- 988 [49] Olsson, E.J., 2002. What is the problem of coherence and truth? *Journal*  
989 *of Philosophy* 99, 246–272.
- 990 [50] Papadimitriou, C.H., 1994. *Computational Complexity*. Addison-  
991 Wesley.
- 992 [51] Park, J.D., Darwiche, A., 2004. Complexity results and approxima-  
993 tion settings for MAP explanations. *Journal of Artificial Intelligence*  
994 *Research* 21, 101–133.



- 995 [52] Pearl, J., 1988. Probabilistic Reasoning in Intelligent Systems: Networks  
996 of Plausible Inference. Morgan Kaufmann, Palo Alto, CA.
- 997 [53] van Rooij, I., 2008. The tractable cognition thesis. *Cognitive Science*  
998 32, 939–984.
- 999 [54] Shi, L., Feldman, N., Griffiths, T., 2008. Performing Bayesian inference  
1000 with exemplar models, in: Sloutsky, V., Love, B., McRae, K. (Eds.),  
1001 Proceedings of the 30th annual conference of the Cognitive Science So-  
1002 ciety, pp. 745–750.
- 1003 [55] Shimony, S., 1993. The role of relevance in explanation I: Irrelevance as  
1004 statistical independence. *International Journal of Approximate Reason-*  
1005 *ing* 8, 281–324.
- 1006 [56] Shimony, S.E., 1994. Finding MAPs for belief networks is NP-hard.  
1007 *Artificial Intelligence* 68, 399–410.
- 1008 [57] Stewart, N., Chater, N., Brown, G.D.A., 2006. Decision by sampling.  
1009 *Cognitive Psychology* 53, 1–26.
- 1010 [58] Sticha, P.J., Buede, D.M., Rees, R.L., 2006. Bayesian model of the effect  
1011 of personality in predicting decisionmaker behavior, in: van der Gaag,  
1012 L. (Ed.), Proceedings of the Fourth Bayesian Modelling Applications  
1013 Workshop.
- 1014 [59] Stockmeyer, L., 1977. The polynomial-time hierarchy. *Theoretical Com-*  
1015 *puter Science* 3, 1–22.
- 1016 [60] Tenenbaum, J.B., 2011. How to grow a mind: Statistics, structure, and  
1017 abstraction. *Science* 331, 1279–1285.
- 1018 [61] Torán, J., 1991. Complexity classes defined by counting quantifiers.  
1019 *Journal of the ACM* 38, 752–773.
- 1020 [62] Vul, E., Goodman, N.D., Griffiths, T.L., Tenenbaum, J.B., 2009. One  
1021 and done? Optimal decisions from very few samples, in: Taatgen, N.,  
1022 van Rijn, H., Schomaker, L., Nerbonne, J. (Eds.), Proceedings of the  
1023 31st Annual Meeting of the Cognitive Science Society, pp. 66–72.
- 1024 [63] Wagner, K.W., 1986. The complexity of combinatorial problems with  
1025 succinct input representation. *Acta Informatica* 23, 325–356.

- 1026 [64] Wasyluk, H., Onisko, A., Druzdzal, M.J., 2001. Support of diagnosis  
1027 of liver disorders based on a causal Bayesian network model. *Medical*  
1028 *Science Monitor* 7, 327–332.
- 1029 [65] Wilson, D., Sperber, D., 2004. Relevance theory, in: R., H.L., Ward, G.  
1030 (Eds.), *Handbook of Pragmatics*. Blackwell, Oxford, UK, pp. 607–632.
- 1031 [66] Yuan, C., Lim, H., Lu, T., 2011. Most relevant explanations in Bayesian  
1032 networks. *Journal of Artificial Intelligence Research* 42, 309–352.