# IBM Research Report

# Efficient Test Selection in Active Diagnosis via Entropy Approximation

**Alice X. Zheng**
Department of EECS
University of California at Berkeley
Berkeley, CA 94720-1776

**Irina Rish, Alina Beygelzimer**
IBM Research Division
Thomas J. Watson Research Center
P.O. Box 704
Yorktown Heights, NY 10598

**Research Division**
**Almaden - Austin - Beijing - Haifa - India - T. J. Watson - Tokyo - Zurich**

# Efficient Test Selection in Active Diagnosis via Entropy Approximation

**Alice X. Zheng**
Department of EECS
U. C. Berkeley
Berkeley, CA 94720-1776
alicez@eecs.berkeley.edu

**Irina Rish**
IBM T.J. Watson Research Center
19 Skyline Drive
Hawthorne, NY 10532
rish@us.ibm.com

**Alina Beygelzimer**
IBM T.J. Watson Research Center
19 Skyline Drive
Hawthorne, NY 10532
beygel@us.ibm.com

## Abstract

We address the problem of active diagnosis on a Bayesian network using most-informative test selection. Finding an optimal subset of tests in this setting is intractable in general. We show that it is difficult even to compute the next most-informative test using greedy test selection, as it involves several entropy terms whose exact computation is intractable. We propose an approximate approach that utilizes the loopy belief propagation infrastructure to simultaneously compute approximations of marginal and conditional entropies on multiple subsets of nodes. We apply our method to fault diagnosis in computer networks, and demonstrate promising empirical results on realistic Internet-like topologies.

## 1 Introduction

The problem of fault diagnosis appears in many places under various guises. Examples include medical diagnosis, computer system troubleshooting, decoding messages sent through a noisy channel, etc. In recent years, diagnosis has often been formulated as an inference problem on a Bayesian network, with the goal of assigning most likely states to unobserved nodes based on outcome of certain "test" nodes.

An important issue in diagnosis is the trade-off between the cost of performing tests and the achieved accuracy of diagnosis. It is often too expensive or even impossible to perform all tests. In this paper, we concentrate on the problem of *active* diagnosis, in which tests are selected sequentially to minimize the cost of testing. We use entropy as the cost function and select a set of tests providing maximum information, or minimum conditional entropy, about the unknown variables.

However, exact computation of conditional entropies in a general Bayesian network can be intractable. While much

existing research has addressed the problem of efficient and accurate probabilistic inference, other probabilistic quantities, such as conditional entropy and information gain, have not received nearly as much attention. There is a vast amount of literature on value of information and most-informative test selection [8, 3, 7, 15], but none of the previous work appears to focus on the computational complexity of most-informative test selection in a general Bayesian network setting.

We propose an approximation algorithm for computing marginal conditional entropy. The algorithm is based on loopy belief propagation, a successful approximate inference method. We illustrate the algorithm at work in the setting of fault diagnosis for distributed computer networks, and demonstrate promising empirical results. We also apply existing theoretical results on the optimality of certain greedy algorithms to our test selection problem, and analyze the effect of approximation error on the expected cost of active diagnosis. Our method is general enough to apply to other applications of Bayesian networks that require the computation of information gain and conditional entropies of subsets of nodes. In our application, it can efficiently compute the information gain for all candidate tests simultaneously.

The paper is structured as follows. Section 2 introduces necessary background and definitions. In section 3, we describe the general problem of active diagnosis and the emerging computational complexity issue. We propose a solution to this problem in section 4. Section 5 discusses an application of our approach in the context of distributed computer system diagnosis, while section 6 presents empirical results. We survey related work in section 7, and conclude in section 8.

## 2 Background and Definitions

Let $\mathbf{X} = \{X_1, X_2, \ldots, X_N\}$ denote a set of $N$ discrete random variables and $\mathbf{x}$ a possible realization of $\mathbf{X}$. A *Bayesian network* is a directed acyclic graph (DAG) $G$ with nodes corresponding to $X_1, X_2, \ldots, X_n$ and edges representing

direct dependencies [14]. The dependencies are quantified by associating each node $X_i$ with a local conditional probability distribution $P(x_i|\mathbf{pa}_i)$, where $\mathbf{pa}_i$ is an assignment to the parents of $X_i$ (nodes pointing to $X_i$ in the Bayesian network). The set of nodes $\{x_i, \mathbf{pa}_i\}$ is called a *family*. The joint probability distribution (PDF) over $\mathbf{X}$ is given as product

$$P(\mathbf{x}) = \prod_{i=1}^{n} P(x_i|\mathbf{pa}_i). \qquad (1)$$

We use $\mathbf{E} \subseteq \mathbf{X}$ to denote a possibly empty set of *evidence* nodes for which observation is available.

For convenience of presentation, we will also use the terminology of *factor graphs*[5], which unifies the directed and the undirected graphical representations of joint probability distributions. A factor graph is an undirected bipartite graph that contains factor nodes (usually shown as squares) and variable nodes (shown as circles). (See Fig. 1 for an example.) There is an edge between a variable node and a factor node if and only if the variable participates in the *potential function* for the corresponding factor. The joint distribution is assumed to be written in a factored form

$$P(\mathbf{x}) = \frac{1}{Z} \prod_{a} f_a(\mathbf{x}_a), \qquad (2)$$

where $Z$ is the normalization constant called the *partition function*, and the index $a$ ranges over all factors $f_a(\mathbf{x}_a)$, defined on the corresponding subsets $\mathbf{X}_a$ of $\mathbf{X}$.

The computation complexity of many probabilistic inference problems can be related to graphical properties. Exact inference algorithms require time and space exponential in the *treewidth*[14] of the graph, which is defined to be the size of the largest clique induced by inference, and can be as large as the size of the graph. Many common probabilistic inference problems are NP-complete [1]. This includes our problem of probabilistic diagnosis, which can be formulated as a *Maximum A posteriori Probability* (MAP) problem: given a set of observations (test outcomes), find the most likely states of unobserved variables.

Although probabilistic inference can be intractable in general, there exists a simple linear-time approximate inference algorithm known as *belief propagation (BP)* [14]. BP is provably correct on polytrees (i.e. Bayesian networks with no undirected cycles), and can be used as an approximation on general networks. Belief propagation passes probabilistic messages between the nodes and can be iterated until convergence. Convergence is guaranteed only for polytrees; otherwise BP is said to diverge.

Let $a$ denote a factor node and $i$ one of its variable nodes. Let $N(a)$ represent the neighbors of $a$, i.e., the set of variable nodes connected to that factor. Let $N(i)$ denote the neighbors of $i$, i.e., the set of factor nodes to which variable node $i$ belongs. The BP message from node $i$ to factor $a$ is defined as [10]:

$$n_{i \to a}(x_i) := \prod_{c \in N(i) \backslash a} m_{c \to i}(x_i), \qquad (3)$$

and the message from factor $a$ to node $i$ is define as

$$m_{a \to i}(x_i) := \sum_{\mathbf{x}_a \backslash x_i} f_a(\mathbf{x}_a) \prod_{j \in N(a) \backslash i} n_{j \to a}(x_j). \qquad (4)$$

Based on these messages, we can compute the beliefs about each node and about the probability potential for each factor, respectively:

$$b_i(x_i) \propto \prod_{a \in N(i)} m_{a \to i}(x_i) \qquad (5)$$

$$b_a(\mathbf{x}_a) \propto f_a(\mathbf{x}_a) \prod_{i \in N(a)} n_{i \to a}(x_i). \qquad (6)$$

Observations are incorporated into the process via $\delta$-functions as local potentials for each node in $\mathbf{E}$. When that is done, $b_i(x_i)$ becomes the approximation of the posterior probability $P(x_i|\mathbf{e})$.

## 3 The Active Test Selection Problem

In many diagnosis problems, the user has an opportunity to actively select tests in order to improve the accuracy of diagnosis. For example, in medical diagnosis, doctors face the *experiment design* problem of choosing which medical tests to perform next.

Let $\mathbf{S} = \{S_1, S_2, \ldots, S_N\}$ denote a set of unobserved random variables we wish to diagnose, and let $\mathbf{T} = \{T_1, T_2, \ldots, T_M\}$ denote the available set of tests. Our objective is to maximize diagnostic quality while minimizing the cost of testing. The diagnostic quality of a subset of tests $\mathbf{T}^*$ can be measured by the amount of uncertainty about $\mathbf{S}$ that remains after observing $\mathbf{T}^*$. From the information-theoretic perspective, this can be measured by the conditional entropy $H(\mathbf{S}|\mathbf{T}^*)$. Clearly, $H(\mathbf{S}|\mathbf{T}) \leq H(\mathbf{S}|\mathbf{T}^*)$ for all $\mathbf{T}^* \subseteq \mathbf{T}$. Thus the problem is to find $\mathbf{T}^* \subseteq \mathbf{T}$ which minimizes both $H(\mathbf{S}|\mathbf{T}^*)$ and the cost of testing. In the case of equally costly tests, this is equivalent to minimizing the number of tests. This problem is known to be NP-hard [?]. However, a simple greedy approach is to choose the next test to be $T^* = \arg\min_T H(\mathbf{S}|T, \mathbf{T}')$, where $\mathbf{T}'$ is the currently selected test set. We give a theoretical anlysis of performances bounds of this greedy approach in the Appendix section. Our current and previous empirical results [9][12] show that the approach also works well in practice.

We make a distinction between off-line test selection and online test selection. In the latter case, previous test outcomes are available when selecting the next test. We will

focus on the online approach, sometimes called *active diagnosis*, which is typically much more efficient in practice than its off-line version [**?**].

**Active Test Selection Problem:** given the observed outcome $\mathbf{t}'$ of previously selected sequence of tests $\mathbf{T}'$, select the the next test $T$ to be $\arg\min_T H(\mathbf{S}|T, \mathbf{t}')$ , where $H(\mathbf{S}|T, \mathbf{t}')$ is the conditional marginal entropy.

The joint entropy $H(\mathbf{X}, \mathbf{T})$ can be decomposed into sum of entropies over the families in a Bayesian networks and thus can be easily computed using only the input CPT specification. Conditional marginal entropies, on the other hand, do not generally have this property: although under certain independence conditions they decompose into functions over the families, computing those functions will require an inference. (See Appendix for proofs.)

**Lemma 1.** *Given a Bayesian network representing a joint PDF $P(\mathbf{X})$, the joint entropy $H(\mathbf{X})$ can be decomposed into the sum of entropies over the families: $H(\mathbf{X}) = \sum_{i=1}^{n} H(X_i|\mathbf{Pa_i})$.* □

**Lemma 2.** *Given a Bayesian network representing a joint PDF $P(\mathbf{S}, \mathbf{T})$, where $\forall i : pa(T_i) \subseteq \mathbf{S}$ (i.e. tests $T_i$ and $T_j$ are independent given a subset of $\mathbf{S}$), the observation $\mathbf{T}' = \mathbf{t}'$ of previously selected test set $\mathbf{T}'$, and a candidate test $T$, the conditional marginal entropy $H(\mathbf{S}|T, \mathbf{t}')$ can be written as*

$$H(\mathbf{S}|T, \mathbf{t}') = - \sum_{t, \mathbf{s}_{pa(T)}} P(\mathbf{s}_{pa(T)}, t|\mathbf{t}') \log P(t|\mathbf{s}_{pa(T)}) \quad (7)$$

$$+ \sum_{t} P(t|\mathbf{t}') \log P(t|\mathbf{t}') + const, \quad (8)$$

*where* const *is a constant expression.* □

As the proof of Lemma 3 demonstrates (see Appendix), the conditional test independence requirement above is indeed necessary for decomposing the conditional entropy.

Note that minimizing conditional entropy is a particular case of *value of information* analysis [7], where the next test $T$ is selected to minimize the expected value of certain *cost function* $c(\mathbf{s}, t, \mathbf{t}')$. The result of Lemma 2 can be generalized to this case if the cost function is decomposable over the families as follows:

**Lemma 3.** *Given a Bayesian network representing a joint PDF $P(\mathbf{S}, \mathbf{T})$, where $\forall i : pa(T_i) \subseteq \mathbf{S}$ (i.e. tests $T_i$ and $T_j$ are independent given a subset of $\mathbf{S}$), the observation $\mathbf{T}' = \mathbf{t}'$ of previously selected test set $\mathbf{T}'$, a candidate test $T$, and a cost function decomposable over the families, i.e. $c(t, \mathbf{s}|\mathbf{t}') = c(t, \mathbf{s}_{pa(t)})$, the expected cost of choosing test $t$ can be written as*

$$E_{P(s,t|t')}c(t, \mathbf{s}|\mathbf{t}') = \sum_{t, \mathbf{s}_{pa(T)}} P(\mathbf{s}_{pa(T)}, t|\mathbf{t}')c(t, \mathbf{s}_{pa(t)}). \quad (9)$$

The remaining part of this paper will focus on the conditional entropy cost. Let $A(T, \mathbf{S}_{pa(T)}|\mathbf{t}')$ denote the first term in Eqn. (8). This is the cross entropy between the posterior probability of $T$ and its parents, and the conditional probability of $T$ given its parents. The second term in Eqn. (8) is simply a negative conditional entropy, $-H(T|\mathbf{t}')$.

**Challenge:** since observations of test outcome correlate the parent nodes, the exact computation of the posterior probabilities in both entropy terms in Eqn. (8) is intractable. We can certainly use a existing approximation method to compute the marginal conditional probabilities $P(\mathbf{s}_{pa(T)}, t|\mathbf{t}')$ and $P(t|\mathbf{t}')$. But a more efficient approach is possible if we exploit the belief propagation infrastructure, as described in the next section.

## 4 Belief Propagation for Entropy Approximation (BPEA)

Let us consider the problem of computing conditional marginal entropies:

$$H(\mathbf{X}_a|\mathbf{e}) = - \sum_{\mathbf{x}_a} P(\mathbf{x}_a|\mathbf{e}) \log P(\mathbf{x}_a|\mathbf{e}) \quad (10)$$

$$\text{where} \quad P(\mathbf{x}_a|\mathbf{e}) = \sum_{\mathbf{x} \backslash \mathbf{x}_a} P(\mathbf{x}|\mathbf{e}),$$

where $\mathbf{x} \backslash \mathbf{x}_a$ are variable nodes not in $\mathbf{x}_a$. The trick is to replace the marginal posterior $P(\mathbf{x}_a|\mathbf{e})$ with its factorized BP approximation, and make use of the BP message passing mechanism to perform the summation over $\mathbf{x}_a$. We call this process Belief Propagation for Entropy Approximation (BPEA).

Pick any node $X_0$ from $\mathbf{X}_a$ and designate it as the root node. We modify the final message passed to $X_0$ as follows:

$$m'_{a \to 0}(x_0) := - \sum_{\mathbf{x}_a \backslash x_0} \tilde{b}_a(\mathbf{x}_a) \log \tilde{b}_a(\mathbf{x}_a). \quad (11)$$

Here, $\tilde{b}_a(\mathbf{x}_a)$ is the unnormalized belief of $X_a$, i.e., $\tilde{b}_a(\mathbf{x}_a) = \sigma b_a(\mathbf{x}_a)$, where $\sigma = \sum_{\mathbf{x}_a} \tilde{b}_a(\mathbf{x}_a)$ is the normalization constant that makes the belief sum to 1.

To get the marginal conditional entropy, we need to sum over the root node $X_0$ and normalize properly.

$$\tilde{h}(\mathbf{X}_a|\mathbf{e}) := \sum_{x_0} m'_{a \to 0}(x_0) \quad (12)$$

$$h(\mathbf{X}_a|\mathbf{e}) := \frac{\tilde{h}(\mathbf{X}_a|\mathbf{e})}{\sigma} + \log \sigma. \quad (13)$$

The proof of correctness is simple and is skipped due to space. It follows immediately that BPEA is exact whenever BP is exact.

The normalization constant $\sigma$ is already computed during normal BP iterations. The computation of $\tilde{b}_a(\cdot)$, $m'_{a \to i}$, and
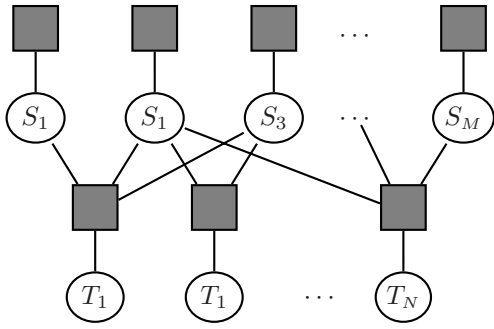
Figure 1: Factor graph of the fault diagnostic Bayes net.

$\tilde{h}(\cdot)$ can all be piggy-backed onto the same BP infrastructure, and therefore does not impact its overall complexity. Furthermore, due to the local and parallel message update procedure in BP, we can compute the marginal posterior entropies of multiple families in one single sweep. This is an important advantage in the active probing setup.

It is also easy to show (details omitted due to space restrictions) that the approach is extendable beyond the entropy computation, to an arbitrary cost function decomposable over families (see Lemma 3). The $c(t, \mathbf{s}_{pa(t)})$ cost function replaces the negative logarithm in Eqns 10 and 11.

## 5 Application: Fault Diagnosis in Computer Networks

Suppose we wish to monitor a system of networked computers. Let $\mathbf{S}$ represent the binary state of $N$ network elements. $S_i = 0$ indicates that the element is in normal operation mode, and $S_i = 1$ indicates that the element is faulty. We can take $S_i$ to be any system component whose state can be measured using a suite of tests. If the system is large, it is often impossible to test each individual component directly. A common solution is to test a subset of components with a single *test probe*. If all the test components are okay, the test would return a 0. Otherwise the test would return 1, but it does not reveal which components are faulty.

We assume there are machines designated as *probe stations*, which are instrumented to send out *probes* to test the response of the network elements represented by $\mathbf{S}$. Let $\mathbf{T}$ denote the available set of probes. A probe can be as simple as a *ping* request, which detects network availability. A more sophisticated probe might be an e-mail message or a webpage-access request. In the absence of noise a probe is a disjunctive test. More generally, it is a noisy-OR test [14]. The joint PDF of all tests and network nodes

forms the well-known QMR-DT model [11].

$$P(s_j) = (\alpha_j)^{s_j}(1 - \alpha_j)^{(1-s_j)}, \qquad (14)$$

$$P(t_i = 0|\mathbf{s}_{\mathbf{pa}_i}) = \rho_{i0} \prod_{j \in \mathbf{pa}_i} \rho_{ij}^{s_j} \qquad (15)$$

$$P(\mathbf{s}, \mathbf{t}) = \prod_i P(t_i|\mathbf{s}_{\mathbf{pa}_i}) \prod_j P(s_j). \qquad (16)$$

Here, $\alpha_j := P(s_j = 1)$ is the prior fault probability, $\rho_{ij}$ is the so-called inhibition probability, and $(1 - \rho_{i0})$ is the leak probability of an unaccounted-for faulty element. The inhibition probability is a measurement of the amount of noise in the network. Fig. 1 shows a factor graph representation of our model.

As discussed in Section 3, we adopt the *active probing* framework for fault diagnosis, sequentially selecting probes to minimize the conditional entropy. In previous work, a single-fault assumption was made, which effectively reduced $\mathbf{S}$ to one random variable with $N+1$ possible states. In general, however, multiple faults could exist in the system simultaneously, which requires the more complicated condition entropies given in Eqn. (8).

We deal with the two entropy terms separately. For $H(T|\mathbf{t}')$, we may use approximation methods such as BP or GBP to calculate the belief $b(t|\mathbf{t}')$, which can then be used to directly compute $H(T|\mathbf{t}')$. (Note that the summation over values of $T$ is simple since $T$ is binary-valued.) To calculate $A(T, \mathbf{S}_{\mathbf{pa}_T}|\mathbf{t}')$, we use the entropy approximation method (BPEA) as described in Section 4. Because BP message updates are done locally, we can compute $A(T, \mathbf{S}_{\mathbf{pa}_T}|\mathbf{t}')$ for all unobserved $T$ nodes during a single application of BP. Thus, picking the next probe requires only one run of the BPEA approximation algorithm.

For each candidate probe, we designate the probe node $T$ itself as the root node. The modified messages are:

$$\tilde{b}_t(t, \mathbf{s}_{\mathbf{pa}_T}) := P(t|\mathbf{s}_{\mathbf{pa}_T}) \prod_{j \in \mathbf{pa}_T} n_{j \to t}(s_j). \qquad (17)$$

$A(T, \mathbf{S}_{\mathbf{pa}_T}|\mathbf{t}')$ is a cross entropy term. Hence we do not take the log of $\tilde{b}$ during BPEA, but rather take the logarithm of the known probabilities $P(t|\mathbf{s}_{\mathbf{pa}_T})$. This simplifies the normalization step described in Eqn. (13) to:

$$A(T, \mathbf{S}_{\mathbf{pa}_T}|\mathbf{t}') = \frac{\tilde{A}(T, \mathbf{S}_{\mathbf{pa}_T}|\mathbf{t}')}{\sigma},$$

where $\sigma = \sum_{t, \mathbf{s}_{pa}(T)} \tilde{b}_t(t, \mathbf{s}_{\mathbf{pa}_T})$.

## 6 Empirical Results

We conduct our experiments on network topologies built by the INET generator[17], which simulates an Internet-like topology at the AS-level. Our dataset includes a set of

networks of 485 nodes, where the number of probe stations varies from 1 to 50.

The connection between probe nodes and network nodes are generated with two goals in mind: detection and diagnosis. A detection probe set needs to cover all network components, so that at least one probe has a positive probability of returning 0 when a component fails. A diagnosis probe set needs to not only cover all components, but also be able to distinguish between faulty componets. Optimal probe set design is NP-hard for either detection or diagnosis. For the datasets used here, we first use a greedy approach to obtain a probe set that covers all network components, then augment this set with additional probes in order to guarantee single-fault diagnosis. Interested readers may find a detailed discussions of probe set design for diagnosis Bayesian networks in [15, 16].

In our experiments, we measure the effects of prior fault probability $\alpha$ and inhibition probability $\rho$ on approximation and diagnostic quality. We compare the approximate entropy values and the quality of the selected probe set against the ground truth, which is obtained via the junction tree exact inference algorithm. In subsection 6.3, we also summarize how the type of network may effect computational efficiency. Since all measurements depend on the particular set of probe outcomes, we repeat all experiments on 10 different samples of the Bayes net.

We use the diagnostic quality of the probe set to determine when to stop the probe selection process: when the reduction in entropy (Eqn. (**??**)) for the past 5 iterations is no more than 0.00001, the selection process is deemed to converge. Otherwise we continue until all probes have been picked.

## 6.1 Approximation accuracy

First, we look at approximation accuracy. Recall that at each time step of the active probing process, we obtain a vector of approximate entropy values, one for each candidate probe $T$. We average the relative error between the approximate values and the exact values for all candidate probes, and further average over all time steps and samples. Let $M$ denote the total number of probes, $n$ the number of selected probes, $h_{ij}$ the approximate value for probe $j$ at the $i$th time step of probe selection, and $H_{ij}$ the corresponding exact values. We compute

$$R(h, H) := \frac{1}{n} \sum_{i=0}^{n-1} \frac{1}{M-i} \sum_{j=1}^{M-i} \frac{|h_{ij} - H_{ij}|}{|H_{ij}|}. \quad (18)$$

This experiment is conducted on the detection network with 10 probe stations augmented with single-node probes. Fig. 2(a-b) contains plot of the average, the minimum, and the maximum approximation errors, taken over 10 samples of probe outcomes. Relative error values are shown sepa-
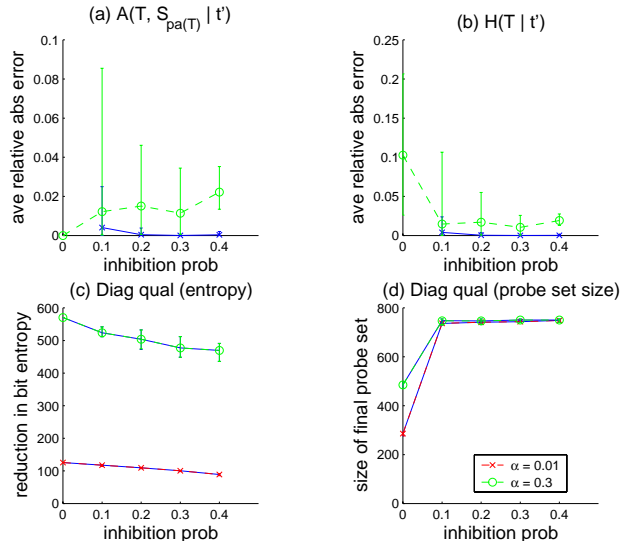


Figure 2: Approximation and diagnostic quality measurements on a augmented detection network, plotted against inhibition probability ($\rho$), and shown at different prior fault probabilities ($\alpha$).

rately for the two entropy terms in Eqn. (8). $A(T, \mathbf{S}_{\mathbf{pa}_T} | \mathbf{t}')$ is calculated using BPEA, whereas $H(T | \mathbf{t}')$ is obtained directly from the BP beliefs $b(t | \mathbf{t}')$. From the two plots, we can see that the approximation error is lower at lower levels of the prior fault probability. For both values of prior fault probability, and for all levels of inhibition probability, the errors do not exceed 2% on average. At the maximum, the approximation error does not exceed 10% for $A(T, \mathbf{S}_{\mathbf{pa}_T} | \mathbf{t}')$, and 20% for $H(T | \mathbf{t}')$.

On the other hand, there does not seem to be a certain relationship between approximation quality and inhibition probability. The BPEA approximation for $A(T, \mathbf{S}_{\mathbf{pa}_T} | \mathbf{t}')$ is slightly better at lower inhibition probabilities. But BP approximation seems to do better at higher inhibition.

## 6.2 Diagnostic quality

The quality of diagnosis is taken to be the reduction in conditional bit entropy of the state of the network elements. That is, if $\mathbf{t}'$ represents the observed outcomes of the final set of selected probes, we measure $H(\mathbf{S}) - H(\mathbf{S} | \mathbf{t}') = -\sum_{\mathbf{s}} P(\mathbf{s}) \log_2 P(\mathbf{s}) + + \sum_{\mathbf{s}} P(\mathbf{s} | \mathbf{t}') \log_2 P(\mathbf{s} | \mathbf{t}')$ .

Fig. 2(c) plots the diagnostic quality of approximate and exact algorithms obtained on the previously mentioned augmented detection network. Note that, for both levels of prior fault probability and at all levels of inhibition, the two algorithms are practically indistinguishable in terms of diagnostic quality. Fig. 2(d) looks at the size of the final probe set (i.e., the number of probes selected when the active probing process is deemed to converge); here again, the two algorithms have identical behavior.
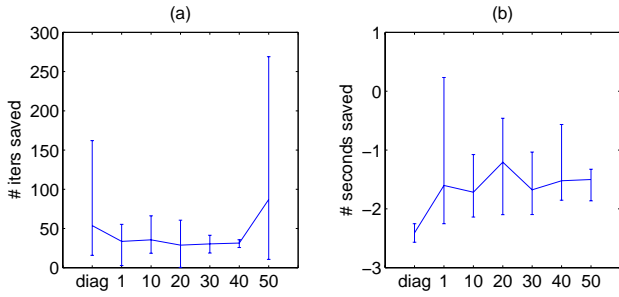
Figure 3: Efficiency of approximate method. (a) Average number of BP iterations saved by re-using messages; (b) Average speed-up compared to exact method.

## 6.3 Implementation and speed

We use the junction tree inference engine in Kevin Murphy's Bayes Net Toolbox [13] for Matlab to obtain exact singleton posterior probabilities. The approximate method is implemented on top of the belief propagation C++/mex code developed by Yair Weiss and Talya Meltzer. We speed up the approximation active probing process by re-using beliefs from previous iterations of BP.

Fig. 3(a) plots the average, maximum, and minimum number of BP iterations that we save through re-using messages. The x-axis denotes the type of network used. The label `diag` represents the diagnosis network with 1 probe station, and the rest are detection networks with various numbers of probe stations. Note that, on average, re-using messages shortens the BP convergence time by 40-50 iterations. This amounts to substantial savings in computation time over the entire active probing process.

Fig. 3(b) compares computation time of the approximate method to the exact method. On average, the approximate method turns out to be slower. Closer examination of the results show that, for most probe selection steps, BP converges under 10 iterations, which puts the approximate method ahead of the exact method. However, for a few of the probes, BP may take several hundred iterations to converge. Thus the average time requirement (per probe selection) of the approximate method is about 2 seconds longer than the exact method. However, keep in mind that, for networks with larger tree-width, the exact method is simply not feasible. Hence, in general, the approximate method is our only choice.

## 7 Related Work

The most-informative test selection was previously addressed in various work on diagnosis, decision analysis, feature selection in machine learning, and related areas. Given a cost function, a common decision-theoretic approach is to compute the *expected value-of-information* [8] of a candidate test, i.e. the expected cost of making a decision after observing the test outcome; using entropy as a cost function yields most-informative test selection. Value of information analysis (and particularly most-informative test selection) was considered in the context of model-based diagnosis [3], probabilistic diagnosis [14] and applied to many practical domains [15]. Previous research has addressed computational complexity of selecting a *set* of most-informative tests instead of a single test [7]. However, none of the previous approaches seem to address the efficiency of computing single-test information gain in a generic Bayesian network.

The most-informative test selection problem is quite similar to the optimal coding problem[2]. There is, however, an important difference. In the coding domain, one may separate source coding (compressing $\mathbf{S}$) from channel coding (adding redundance to improve decoding accuracy). Fault diagnosis, on the other hand, has to deal with a combination of the two, which manifests itself in the nature of available tests (described by the conditional probabilities $P(T_i|\mathbf{S}_{pa(i)})$). We may have no control over the source coding function, but we can still select the smallest, most informative subset of tests.

In the context of probing, i.e. disjunctive testing, optimal test selection is very similar to the *group testing* problem [4]. Given a set of Boolean variables representing objects that can be in two possible states (i.e. sick vs. healthy patients, failed vs. OK nodes), the objective of group testing is to find all 'failed' objects by using a sequence of disjunctive tests. Particularly, sequential test selection is known as *adaptive group testing* [4]. There is also a direct connection between adaptive group testing and Golomb codes [6]. Note that group testing assumes no constraints on the test selection (i.e., any subset of objects can be tested together), while in Bayesian networks the tests can be only selected from a fixed set. Even in a less restrictive case of probe selection we are still constrained by the network topology. Constrained group testing (and coding in general) appears to be more complicated, particularly for theoretical analysis, than its unconstrained version.

## 8 Conclusions

We propose an entropy approximation method based on loopy belief propagation, and examine its bahavior on the application of active probing for fault diagnosis in a networked computer system. The level of approximation error is found to vary with the level of noise. However, even with non-zero approximation errors, the diagnosis quality is practically identical to that obtained from the exact method. BPEA approximation takes slightly longer than the exact method on small networks. But it can handle much larger networks for which exact junction tree inference is infeasible. This highlights a promising direction for

active probing and fault diagnosis, and for entropy approximation on Bayesian networks in general.

## References

[1] G.F. Cooper. The computational complexity of probabilistic inference using Bayesian belief networks. *Artificial Intelligence*, 42(2–3):393–405, 1990.

[2] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley & Sons, 1991.

[3] J. de Kleer and B.C. Williams. Diagnosing Multiple Faults. *Artificial Intelligence*, 32(1), 1987.

[4] D-Z. Du and F.K. Hwang. *Combinatorial Group Testing and Its Applications (2nd edition)*. World Scientific, 2000.

[5] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Trans. Info. Theory*, pages 498–519, 2001.

[6] R. Gallager and D. Van Voorhis. Optimal source codes for geometrically distributed integer alphabets. *IEEE Trans. Information Theory*, IT-21:228–230, 1975.

[7] D. E. Heckerman, E. J. Horvitz, and B. Middleton. An approximate nonmyopic computation for value of information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15:292–298, 1993.

[8] R. Howard. Information value theory. *IEEE Trans Syst Sci Cybern*, 2(1):22–26, 1966.

[9] I. Rish, M. Brodie, N. Odintsova, S. Ma, G. Grabarnik. Real-time problem determination in distributed systems using active probing. In *NOMS*, 2004.

[10] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Constructing free energy approximations and generalized belief propagation algorithms. Technical Report TR-2004-040, MERL, May 2004.

[11] T. Jaakkola and M. Jordan. Variational probabilistic inference and the qmr-dt database. *Journal of Artificial Intelligence Research*, pages 291–322, 1999.

[12] M. Brodie, I. Rish, S. Ma, N. Odintsova. Active probing strategies for problem diagnosis in distributed systems. In *IJCAI*, 2003.

[13] K. Murphy. The bayes net toolbox for matlab. *Computing Science and Statistics*, 2001.

[14] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, San Mateo, California, 1988.

[15] I. Rish, M. Brodie, N. Odintsova, S. Ma, and G. Grabarnik. Real-time Problem Determination in Distributed Systems using Active Probing. In *Proceedings of 2004 IEEE/IFIP Network Operations and Management Symposium (NOMS 2004), Seoul, Korea*, 2004.

[16] Irina Rish, Mark Brodie, and Sheng Ma. Intelligent probing: a Cost-Efficient Approach to Fault Diagnosis in Computer Networks. *IBM Systems Journal*, 41(3):372–385, 2002.

[17] J. Winick and S. Jamin. Inet-3.0: Internet topology generator. Technical Report CSE-TR-456-02, University of Michigan, 2002.

## Appendix

**Lemma 1.** *Proof:* From Eqn. (1) we get $H(\mathbf{X}) = -\sum_{\mathbf{x}} P(\mathbf{x}) \sum_{i=1}^{n} \log P(x_i|\mathbf{pa}_i) = -\sum_{i=1}^{n} \sum_{x_i,\mathbf{pa}_i} P(x_i,\mathbf{pa}_i) \log P(x_i|\mathbf{pa}_i) = \sum_{i=1}^{n} H(X_i|\mathbf{Pa_i})$. $\qquad\square$

**Lemma 2.** *Proof:* $H(\mathbf{X}|T,\mathbf{t}') = H(\mathbf{X},T|\mathbf{t}') - H(T|\mathbf{t}')$ where $H(T|\mathbf{t}') = -\sum_t P(t|\mathbf{t}') \log P(t|\mathbf{t}')$, and

$$H(\mathbf{S},T|\mathbf{t}') = -\sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(\mathbf{s},t|\mathbf{t}')$$

$$= -\sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(\mathbf{s},t,\mathbf{t}') + \sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(\mathbf{t}')$$

$$= -\sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(\mathbf{s},t,\mathbf{t}') + \log P(\mathbf{t}'). \quad (19)$$

The last term in Eqn. (19) is independent of T and can be replaced by $const$. Since by lemma's condition, $pa(T_i) \subseteq \mathbf{S}$, the joint $P(\mathbf{s},t,\mathbf{t}')$ is factored as

$$P(\mathbf{s},t,\mathbf{t}') = P(t|\mathbf{s}_{pa(t)}) \prod_j P(t'_j|\mathbf{s}_{pa(j)}) P(\mathbf{x}).$$

Note that the above condition is essential since in general $P(\mathbf{s},t,\mathbf{t}')$ may not factorize, and no further simplifications would be possible. However, under this condition, the first term in Eqn. (19) can be written as

$$\sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(\mathbf{s},t,\mathbf{t}') = \sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') \log P(t|\mathbf{s}_{pa(t)})$$

$$+ \sum_{\mathbf{s}} P(\mathbf{s}|\mathbf{t}') \log P(\mathbf{x}) \prod_j P(t'_j|\mathbf{s}_{pa(j)}).$$

Again, the last term above does not involve $T$. The first term can be simplified as $\sum_{\mathbf{s}_{pa(t)},t} P(\mathbf{s}_{pa(t)},t|\mathbf{t}') \log P(t|\mathbf{s}_{pa(t)})$. Hence,

$$H(\mathbf{S}|T,\mathbf{t}') = -\sum_{t,\mathbf{s}_{pa(t)}} P(\mathbf{s}_{pa(t)},t|\mathbf{t}') \log P(t|\mathbf{s}_{pa(t)}) +$$

$$\sum_t P(t|\mathbf{t}') \log P(t|\mathbf{t}') + \text{ const. } \quad\square$$

**Lemma 3.** *Proof:* $E_{P(\mathbf{s},\mathbf{t}|\mathbf{t}')} c(t,\mathbf{s}|\mathbf{t}') = \sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') c(t,\mathbf{s}|\mathbf{t}') = \sum_{\mathbf{s},t} P(\mathbf{s},t|\mathbf{t}') c(t,\mathbf{s}_{pa(t)}) = \sum_{\mathbf{s}_{pa(t)},t} P(\mathbf{s}_{pa(t)},t|\mathbf{t}') c(t,\mathbf{s}_{pa(t)})$. $\qquad\square$

### Approximation Quality of the Greedy Strategy

Recall that we have a set $\{S_1, \ldots, S_N\}$ of binary variables representing the state of $N$ elements, and the goal is to decode this state using tests.

For simplicity, we will only consider the case when tests correspond to deterministic disjunctions, i.e., all inhibition probabilities and the leak probability are zero. Thus an outcome of a test splits the state space into the states consistent with the outcome and those that are not.

If any combination of the elements can be faulty, the system can, in principle, be in any one of $2^N$ possible states. However, the effective state space of a test $T$ involving $n$ elements contains $2^n$ "states," each corresponding to $2^{N-n}$ states in the original state space $\{0,1\}^N$. If the prior probability of fault is $\alpha$, then the probability of such partial assignment $\mathbf{s} \in \{0,1\}^n$ is $P(\mathbf{s}) = \alpha^{n_1}(1-\alpha)^{n-n_1}$, where $n_1$ is the number of faults (1's) in $\mathbf{s}$ (and assuming that faults are independent). Test $T$ splits this effective state space into two sets, corresponding to outcome 0 (with probability mass $(1-\alpha)^n$) and outcome 1 (with probability mass $1 - (1-\alpha)^n$) respectively. States within each set are indistinguishable by $T$.

A natural alternative to selecting the most informative test, is to pick the test that gives the most "balanced" partition of the current state space $\mathbf{S}^*$. Initially, when all states are indistinguishable, $\mathbf{S}^* = \{0,1\}^N$. Let $P(\mathbf{S}^* \mid T = 0)$ be the probability mass of states in $\mathbf{S}^*$ consistent with outcome 0 of test $T$. Similarly define $P(\mathbf{S}^* \mid T = 1)$ for outcome 1. At every step, the next test is taken to be $\text{argmin}_T |P(\mathbf{S}^* \mid T = 0) - P(\mathbf{S}^* \mid T = 1)|$. After the outcome of $T$ becomes known, we discard the states in $\mathbf{S}^*$ inconsistent with this outcome.

The same greedy strategy was used by Dasgupta [?] in the context of actively learning a concept by adaptive queries, and by Kosaraju et al. [?] in a general setting.
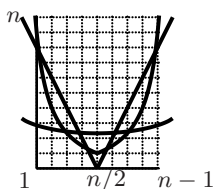
We show that this balance-based strategy is equivalent to the strategy based on choosing the most informative test. First we need to define the cost of a solution.

**Cost of diagnosis** Both greedy strategies produce a tree with leaf nodes corresponding to possible states, and non-leaf nodes corresponding to tests, assuming that tests are informative enough to distinguish among the states (otherwise leaves correspond to distributions over subsets of states). The cost $c(\mathbf{s})$ of diagnosing leaf $\mathbf{s}$ in the set of leaves $\mathbf{S}$ is the number of tests on the path from the root to $\mathbf{s}$. The cost of a tree $D$ is given by $c(D) = \sum_{\mathbf{s}\in\mathbf{S}} P(\mathbf{s})c(\mathbf{s})$, the expected cost needed to diagnose a leaf chosen according to $P$.

**Equivalence** Let $\mathbf{S}^*(\mathbf{t})$ be the set of states in $\mathbf{S}$ consistent with outcomes $\mathbf{t}$ of tests selected so far. For a test $T$, let $U$ be the set of leaves in $\mathbf{S}^*(\mathbf{t})$ consistent with outcome 0 of $T$. Also let $n_0 = |\overline{U}|$, $n_1 = |U|$, and $n = |\mathbf{S}^*(\mathbf{t})| = n_0 + n_1$.

To simplify the argument, assume that all weights in $\mathbf{S}^*(\mathbf{t})$ are equal; the argument readily holds for any set of weights.

The information gain is then maximized by $T$ minimizing $\frac{1}{n} \sum_{i\in\{0,1\}} n_i \log n_i$. The most balanced split in this case is the one minimizing $|n_0 - n_1|$. Both functions are shown below. The V-shaped curve is the (dis)balance of a split with the corresponding $n_0$. The smooth flat curve is the conditional entropy of the state given the result of the split (and all previous splits). The remaining curve shows the ratio $\max\{n_0, n_1\} / \min\{n_0, n_1\}$, another equivalent splitting criteria. It is clear that the functions achieve their minimum at the same point.



Both Dasgupta [?] and Kosaraju et al. [?] showed that the balance-based scheme results in a tree whose cost is within a factor of $\log(1/\min_{\mathbf{s}} P(\mathbf{s}))$ from optimal, where $\min_{\mathbf{s}} P(\mathbf{s})$ is the probability of the least probable state in $\mathbf{S}$. Furthermore, Kosaraju et al. [?] showed that a slight reweighting (needed if $P$ is exponentially unbalanced), results in a tree whose cost is within a factor of $O(\log N)$ from optimal, for any $P$. Thus a tree $D$ obtained by the greedy algorithm satisfies $C(D) \leq O(\log N)C(D^*)$, where $c(D^*)$ is the optimal cost.

As follows from the equivalence, the greedy algorithm that chooses the most informative split also results in a tree whose cost is within a factor of $O(\log n)$ from optimal.

Kosaraju et al. [?] also claimed that the guarantees hold for an algorithm that only *approximates* the most balanced partition. Assume without loss of generality that $P(\mathbf{S}^*|T = 0) \geq P(\mathbf{S}^*|T = 1)$, or in our case $n_0 \geq n_1$. The results says that if the best balance ratio $P(\mathbf{S}^*|T = 0)/P(\mathbf{S}^*|T = 1)$, or in our simple case, $n_0/n_1$, is approximated within a constant multiplicative factor, then the $O(\log N)$ approximation guarantee holds.

Let the most balanced split be $\{x^*, n - x^*\}$; without loss of generality, $x^* \geq n/2$. Assume that the approximate split is $\{x^* + a, n - x^* - a\}$ for some $0 < a \leq n - 1 - x^*$. We want to upper bound the ratio:

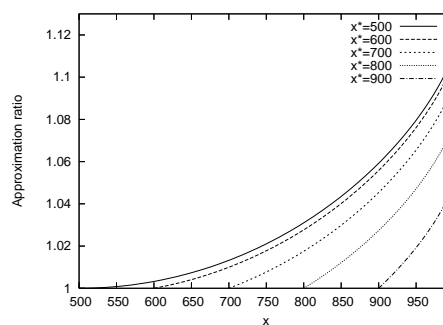$$\frac{(x^* + a)/(n - x^* - a)}{x^*/(n - x^*)} = 1 + \frac{an}{x^*(n - x^* - a)}.$$

Now $x^* \geq n/2$, and $n - x^* - a \geq 1$ since $a \leq n - 1 - x^*$. Thus the ratio is bounded by $1 + 2a$, which is constant if $a$ is constant.

It remains to show that good approximations of the conditional entropy result in small $a$'s. Note that *any* approximation of the conditional entropy is within a multiplicative factor of roughly $1 + (\log n)^{-1}$. Indeed, we have

$$\frac{x \log x + (n - x) \log(n - x)}{x^* \log x^* + (n - x^*) \log(n - x^*)},$$

which is maximized when $x$ is maximized and $x^*$ is minimized, or $x = n - 1$ and $x^* = n/2$, implying the upper bound.

The plot below shows the approximation ratio for the conditional entropy. Each curve corresponds to a particular exact split $x^*$ starting from $n/2$; here $n = 1000$. For every $x^*$, we plotted the approximation ratio as a function of approximate split $x = x^*+a$. Notice that the worst case ratio is roughly 1.1, as expected. Also notice that if the approximate value is sufficiently close to the exact value, the curves are well approximated by lines (with larger exact values being harder to approximate). Thus for good approximations, the approximation ratio for conditional entropy translates roughly "linearly" into $a$. Thus, according to the result of Kosaraju et al. [?], we should expect a $O(\log N)$ approximation ratio even if entropies are only approximate.



**State spaces** Note that if all combinations of faults are possible, a simple information-theoretic argument shows that we need at least $N$ tests to uniquely distinguish between these states, which is as inefficient as testing each element directly. Of course, we can stop the diagnosis when the conditional entropy is sufficiently low (i.e., when we have an almost deterministic distribution on some subset of states), and then output the most likely state. This way we can often approximate the state with significantly fewer tests.

Another common approach is to assume an upper bound on the number of faults. For example, if we have prior fault probabilities $\alpha_i = P(S_i = 1)$, the expected number of faults is $\sum_{i=1}^{N} \alpha_i$ (assuming that faults are independent); hence if $\alpha_i = \alpha$ for all $i$, this number is $\alpha N$. By Markov's inequality, the probability that the actual number of faults is more than $\alpha N c$ is at most $1/c$ for any $c \geq 1$, thus we can typically assume that the state space $\mathbf{S}$ is the set of all subsets of at most $\alpha N c$ elements, for appropriate $c$.