

Empirical Hardness of Finding Optimal Bayesian Network Structures

Algorithm Selection and Runtime Prediction

Brandon Malone · Kustaa Kangas ·
Matti Järvisalo · Mikko Koivisto ·
Petri Myllymäki

Received: date / Accepted: date

Abstract Various algorithms have been proposed for finding a Bayesian network structure that is guaranteed to maximize a given scoring function. Implementations of state-of-the-art algorithms, *solvers*, for this Bayesian network structure learning problem rely on adaptive search strategies, such as branch-and-bound and integer linear programming techniques. Thus, the time requirements of the solvers are not well characterized by simple functions of the instance size. Furthermore, no single solver dominates the others in speed. Given a problem instance, it is thus *a priori* unclear which solver will perform best and how fast it will solve the instance.

We show that for a given solver the hardness of a problem instance can be efficiently predicted based on a collection of non-trivial features which go beyond the basic parameters of instance size. Specifically, we train and test statistical models on empirical data, based on the largest evaluation of state-of-the-art exact solvers to date. We demonstrate that we can predict the runtimes to a reasonable degree of accuracy. These predictions enable effective selection of solvers that perform well in terms of runtimes on a particular instance. Thus, this work contributes a highly efficient portfolio solver that makes use of several individual solvers.

This work is supported by Academy of Finland, grants #125637, #251170 (COIN Centre of Excellence in Computational Inference Research), #255675, #276412, and #284591; Finnish Funding Agency for Technology and Innovation (project D2I); and Research Funds of the University of Helsinki.

Brandon Malone

Section of Bioinformatics and Systems Cardiology, Department of Internal Medicine III and Klaus Tschira Institute for Integrative Computational Cardiology, University of Heidelberg, Germany, DZHK (German Centre for Cardiovascular Research), Partner site Heidelberg/Mannheim, Germany E-mail: brandon.malone@uni-heidelberg.de

Kustaa Kangas · Matti Järvisalo · Mikko Koivisto · Petri Myllymäki

Helsinki Institute for Information Technology HIIT, Department of Computer Science, University of Helsinki, Finland

1 Introduction

Since the formalization and popularization of Bayesian networks [55] for modeling and reasoning with multiple variables, much research has been devoted to learning them from data [28]. One of the main challenges has been to learn the model structure, represented by a directed acyclic graph (DAG) on the variables. Cast as a problem of finding a DAG that is a global optimum of a score function for given data, the *Bayesian network structure learning* problem (BNSL) is notoriously NP-hard; the hardness is chiefly due to the acyclicity constraint imposed on the DAG to be learned [14]. To cope with the computational hardness, early work on structure learning resorted to local search algorithms. While local search algorithms oftentimes perform well, they are unfortunately unable to guarantee global optimality. The uncertainty about the quality of the found network hampers the use of the network [50] in probabilistic inference and causal discovery.

The last decade has raised hopes of solving larger problem instances to optimality. The first algorithms guaranteed to find an optimum adopted a dynamic programming approach to avoid exhaustive search in the space of DAGs [53, 37, 63, 62]. Later algorithms have expedited the dynamic programming approaches using the A* search algorithm with various admissible heuristics [72], or have employed quite different approaches, such as branch and bound in the space of (cyclic) directed graphs [11], integer linear programming (ILP) [35, 16, 17], and constraint programming (CP) [5]. In this work, we focus on such *complete solvers* for BNSL, which we call simply *solvers*. Our interest is in unsupervised learning of a joint structure over the variables, only noting in passing that alternative methods have been developed for supervised learning of the relationship between a designated response variable and the other predictor variables (see, e.g., a recent survey [7] and references therein).

Due to the intrinsic differences between the algorithmic approaches underlying BNSL solvers, it is not surprising that their relative efficiency varies greatly on a per-instance basis. To exemplify this, a comparison of the runtimes of three current state-of-the-art solvers, based on A*, ILP, and CP, is illustrated in Figure 1 using typical benchmark datasets. Evidently, no single one of these three solvers dominates the other two.

Figure 1 suggests that, to improve over the existing solvers, an alternative to developing yet another solver is to design *algorithm portfolios* which select a solver to run on a per-instance basis, ideally combining the best-case performance of the different solvers. Indeed, in this work we do not focus on developing or improving an individual algorithmic approach. Instead, we aim to characterize how the performance of different algorithmic approaches depends on the problem instance, which is the key to the design of efficient algorithm portfolios. The underlying motivation for developing such techniques is the aim of improving the efficiency of state of the art in complete solvers in solving hard BNSL instances.

In this quest, it is vital to discover a collection of *features* that are efficient to compute and yet informative about the hardness of an instance for a solver. Prior work has identified two simple features, namely the number of variables and the number of so-called *candidate parent sets*, denoted by n and m , respectively. To explain the observed orthogonal performance characteristics shown in Figure 1, it has been suggested, roughly, that typical instances can be solved to optimality by A*, if n is at most 40 (no matter how large m), and by ILP if m is moderate, say, at

most some tens of thousands (no matter how large n) [17,72]; for the more recent CP approach, we are not aware of any comparable description. Beyond this rough characterization, the practical time complexity of the best-performing solvers is currently poorly understood. This stems from the sophisticated search heuristics employed by the solvers, which tend to be sensitive to small variations in the instances, thus resulting in somewhat chaotic-looking behavior of runtimes. Furthermore, the gap between the analytic worst-case and best-case runtime bounds, in terms of n and m , is huge, and typical instances fall somewhere in between the two extremes.

The starting point of our work is the following basic open question:

Q1 For *determining the fastest* of the available solvers on a given instance, do the simple features, the number of variables and the number of candidate parent sets, suffice?

We answer this question in the affirmative. Our result is empirical in that it relies on training and testing a statistical model with a large set of problem instances collected from various sources. We show that a simple set of features yields a model which accurately predicts the fastest solver for a given instance based on the parameters n and m only. Furthermore, we show how this yields an algorithm portfolio that almost always runs as fast as the fastest solver, thus significantly outperforming any fixed solver on a large collection of instances.

However, a closer inspection reveals that the predicted runtimes of the model based on the simple features often differ from the actual runtimes by one to two orders of magnitude. The large deviations suggest that, if the interest is in accurate estimation of the runtimes, then the simple feature set may not suffice. This observation motivates our second question:

Q2 For *predicting the runtime of a solver* on a given instance, can the accuracy be significantly improved by including additional efficiently computable features?

Also to this question our answer is affirmative. We introduce and study several additional features that capture the hardness of the problem more accurately for a given solver. We focus on what are currently the three top-performing solver

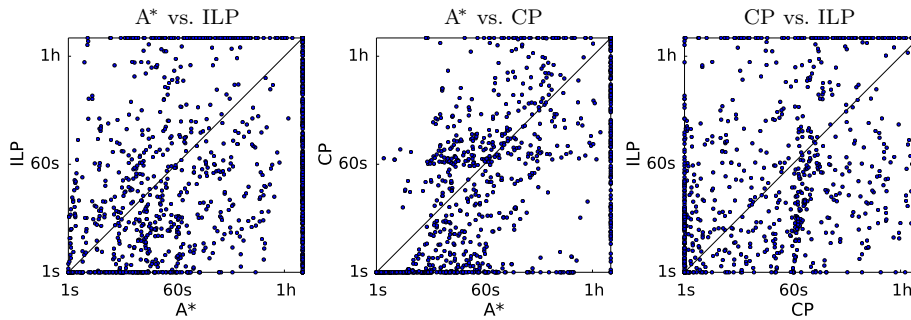


Fig. 1 Comparison of three state-of-the-art algorithms for finding an optimal Bayesian network. Runtimes below 1 or above 7,200 seconds are rounded to 1 and 7,200, respectively. The specific solver parameterizations are A*-comp (A*), ilp-162 (ILP), and cpbayes (CP); see Section 5 for descriptions of the solvers and the datasets.

families based on A* [72], ILP [17], and CP [5], which clearly dominate earlier approaches based on dynamic programming and branch-and-bound [51]. Specifically, we show that models with a wider variety of features yield at times significant improvements in prediction accuracy.

Besides the aforementioned contributions, the empirical work associated with this paper also provides the most elaborate evaluation of state-of-the-art solvers to date, significant in its own right.

The present work extends and revises substantially our preliminary study reported at the AAAI-14 conference [51]. Here we have thoroughly revised the methodology and analysis presented throughout the paper. We have expanded the portfolio itself to include the very recent CP-based solver [5]. At the same time, we have updated the runtime results to the most recent versions of the A*-based and ILP-based solvers. Furthermore, we provide a more fine-grained analysis by categorizing datasets based on their origin. Our results show that the origin of the dataset significantly affects the relative runtime performance of solvers. To this end, we have also increased the number of synthetic data sets considerably, from a few dozens to several hundred. Finally, we provide a more extensive discussion of the characteristics of the learned models, such as preprocessing strategies.

1.1 Related Work

Due to the wide range of potential applications, the general research area of algorithm selection, with tight connections to machine learning and algorithm portfolio design, is very diverse. Instead of aiming at a full review of the relevant literature, here we aim at a brief overview of the research area by providing references to some of the key early works on the topic and some of the more recent works most closely related to ours. For an expanded discussion of the literature on algorithm selection and runtime prediction, we refer the reader to two recent surveys on the topic with further pointers to related work [34, 39].

Research on algorithm selection for various types of important computational problems has its roots in [58], where the algorithm selection problem was introduced, and feature-based modeling was proposed to facilitate the selection of the best-performing algorithm for a given problem instance, considering various example problems. Later works, including [12, 21, 48], demonstrated the efficacy of applying machine learning techniques, such as Bayesian approaches [31], to learn models from empirical performance data.

More recently, empirical hardness models [45, 46] have been applied in the construction of solver portfolios [26] for various NP-hard search problems [40], including Boolean satisfiability (SAT) (e.g., in [70]), constraint programming (e.g., in [24, 32]), quantified Boolean formula satisfiability (e.g., in [57]), answer set programming (e.g., in [30]), as well as for the traveling salesperson problem (e.g., in [41]). To the best of our knowledge, for the important problem of Bayesian network structure learning, the present work is the first to adopt the approach.

In terms of terminology, we investigate algorithm selection in the context of learning Bayesian networks, which is an *unsupervised* learning task. Nevertheless, this work is well-situated in the context of meta-learning [25], which most often considers supervised settings. The BNSL features we propose in Section 3.1 are exactly a set of *meta-features* for this particular domain. The regression models

we learn (Section 3.2) capture meta-knowledge about the state-of-the-art BNSL solvers.

Previous work [42, 44] has suggested that in many cases, a small set of features can lead to accurate predictions; indeed, in Section 5.2 we show that a very small number of features leads to near-optimal algorithm selection performance. Furthermore, while that work relied on qualitative visual analysis, in Section 6.4 we quantify the utility of each feature using the Gini importance [9].

Recently, a simple “Best in Sample” approach [59] was shown to be very effective for algorithm (classifier) selection in the supervised setting. Briefly, this approach trains each classifier in the portfolio using a very small subset of the data; it then selects the classifier to use based on performance on the subset. “Probing” features—a central form of features in, for example, SAT portfolios [70]—we apply in the context of BNSL (see Section 3.1) are similar in spirit to this approach, though adapted to the unsupervised learning setting. In terms of evaluation, our virtual best solver comparisons in Section 5 are quite similar to Loss Curves [43], which have previously been used in the context of meta-learning.

1.2 Organization

The remainder of this paper is organized as follows. We begin in Section 2 by describing the problem of structure learning in Bayesian networks and by giving an overview of the algorithmic techniques underlying the state-of-the-art solvers. Section 3 presents the building blocks of our empirical hardness model: we introduce several BNSL features; we choose an appropriate statistical learning framework; and we describe the methods we use for training and evaluating the models. In Section 4, we present the experimental setting, namely technical details of the investigated solvers and characteristics of the collected problem instances. Results on learning solver portfolios and on predicting runtimes of individual solvers are reported in Sections 5 and 6, respectively. Finally, we discuss some questions that are left open and directions for future research in Section 7.

2 Learning Bayesian Networks

A Bayesian network (G, P) consists of a directed acyclic graph (DAG) G on a set of random variables X_1, \dots, X_n and a joint distribution P of the variables such that P factorizes into a product of the conditional distributions $P(X_i | G_i)$. Here G_i denotes the set of parents of X_i in G ; we call a variable X_j a *parent* of X_i , and X_i a *child* of X_j , if G contains an arc from X_j to X_i .

2.1 The Structure Learning Problem

In its simplest form, structure learning in Bayesian networks concerns finding a DAG that best fits some observed data on the variables.¹ Throughout this work, we only deal with this optimization formulation, here only mentioning that there

¹ Strictly speaking, the data are assumed to consist of N independent and identically distributed tuples (X_1^t, \dots, X_n^t) , $t = 1, \dots, N$, so the dimension of the data is $N \times n$.

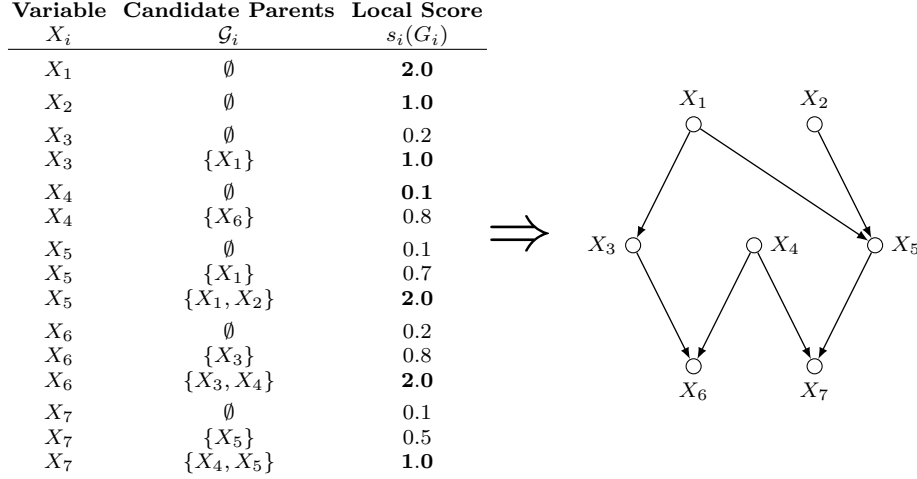


Fig. 2 An optimal DAG (on the right) for a given scoring function s (on the left). There are $n = 7$ variables and $m = 15$ candidate parent sets in total. The optimal score, 8.1, is the sum of the local scores shown in bold face. Observe that choosing $G_4 = \{X_6\}$ would have increased the score but violated the acyclicity constraint.

are also other popular formulations based on frequentist (multiple) hypothesis testing [65, 13] and Bayesian model averaging [49, 23, 37].

The goodness of fit is typically measured by a real-valued scoring function s , which associates a DAG G with a real-valued score $s(G)$.² Frequently used scoring functions are based on (penalized) maximum likelihood, minimum description length, or Bayesian principles (e.g., BDeu and other forms of marginal likelihood). Additionally, they *decompose* [28] into a sum of local scores $s_i(G_i)$ for each variable X_i and its set of parents G_i . In principle, for each i the local scores are defined for all the 2^{n-1} possible parent sets. However, in practice this number is greatly reduced by enforcing a small upper bound for the size of the parent sets G_i or by pruning, as preprocessing, parent sets that provably cannot belong to an optimal DAG [67, 11]. Applying one or both of these reductions results in a collection of *candidate parent sets*, which we will denote by \mathcal{G}_i .

This motivates the following formulation of the *Bayesian network structure learning* problem (BNSL).

INPUT: Local scores $s_i(G_i)$ for a collection of candidate parent sets $G_i \in \mathcal{G}_i$ for $i = 1, \dots, n$.
 TASK: Find a DAG G that maximizes the score $s(G) = \sum_i s_i(G_i)$.

Along with the number of variables n , another key parameter describing the input size is the total number of candidate parent sets $m = \sum_i |\mathcal{G}_i|$. See Figure 2 for an example instance of the BNSL problem.

² The score does not depend on the parameters of the unspecified distribution P , which are treated as nuisance parameters and absorbed by the scoring function (e.g., estimated or integrated away).

2.2 Overview of Complete Solvers for BNSL

We call an algorithm that is guaranteed to find a global optimum and prove its optimality for the BNSL problem a *complete solver* for BNSL, or simply a *solver*. In the next paragraphs we review some state-of-the-art solvers that fit the scope of our study. We omit algorithms that assume significant additional constraints given as input [56] or massive parallel processing [66, 54].

Several works [53, 37, 62] have proposed dynamic programming algorithms to solve BNSL. The solvers are based on the early observation [10, 15] that for any fixed ordering of the n variables, the decomposability of the score enables efficient optimization over all DAGs compatible with the ordering. The algorithms proceed by adding one variable at a time, only tabulating partial solutions for the explored *subsets* of the variables. Thus the runtime scales roughly as 2^n .

Yuan and Malone [72] formulated BNSL as a state-space search through the dynamic programming lattice and applied the A* search algorithm. Unlike the other sophisticated solvers, A* maintains the meaningful worst-case time bound of dynamic programming. To this end, they developed several admissible heuristics which relax the acyclicity constraint; these allow the algorithm to prune suboptimal paths during search, thus typically avoiding visiting all the variable subsets.

The branch-and-bound style algorithm by de Campos and Ji [11] searches in a relaxed space of directed graphs that may contain cycles. It begins by allowing all variables to choose their optimal parents, which typically results in some number of cycles. Then, any found cyclic solutions are iteratively ruled out: it finds a cycle and breaks it by removing one arc in it, branching over the possible choices of the arc. It examines graphs in a best-first order, so the first acyclic graph it finds is an optimal DAG. In this way, the algorithm ignores many cyclic graphs.

Integer linear programming (ILP) algorithms by Jaakkola et al. [35] and by Bartlett and Cussens [16, 17] search in a geometric space, in which DAGs appear as vertices of an embedded polytope, corresponding to integral solutions to a linear program (LP). A series of LP relaxations are solved, and the solution to each relaxation is checked for integrality; an integral solution corresponds to an optimal DAG. The search space is effectively pruned by employing domain-specific cutting planes.

A very recent development in solvers for BNSL is the constraint programming (CP) based approach by van Beek and Hoffmann [5], constituting a constraint-based depth-first branch-and-bound approach to BNSL. As a key ingredient, the approach uses an improved constraint model with problem-specific dominance, symmetry breaking, and acyclicity constraints and propagators. It also employs cost-based pruning rules applied during search, together with domain-specific search heuristics. The approach combines some of the ideas applied in A*, specifically pattern databases, for obtaining bounds on the scoring function.

3 Empirical Hardness Models

In this work, we focus on the *hardness* of a BNSL instance, relative to a particular solver. We define the hardness of instance I for solver S simply as the runtime

$T_S(I)$ of the solver S on the instance I .³ Due to the sophisticated heuristics underlying the state-of-the-art BNSL solvers, evaluating the empirical hardness is presumably (that is, under standard complexity-theoretic assumptions) computationally intractable; indeed, the fastest method we are aware of for evaluating $T_S(I)$ is actually running S on I .

Rather than exactly evaluating the function T_S , we take a machine learning approach to approximate it: from a large collection of example instances for which the actual runtimes are known (computed), we learn a *model* which is efficient to evaluate at any given instance. Underlying this approach is the hypothesis that an accurate *empirical hardness model* [45] can be built based on a set of efficiently computable *features* of BNSL instances; by a feature we refer to a mapping from the instances to the real numbers. This approach naturally gives rise to the following supervised machine learning problem, for a fixed solver S .

- INPUT: A training set of BNSL instances (represented as collections of feature values) and the respective runtimes of the solver S .
- TASK: Learn a function \hat{T}_S which minimizes the average prediction error on an unseen set of BNSL instances.

We next introduce several categories of efficiently-computable features of BNSL instances. Most of these features have not previously been used for characterizing the hardness of BNSL. We then explain our training and testing strategies.

3.1 Features for BNSL

We use four different sets of features based on complementary strategies to characterize BNSL instances: **Basic**, **Basic extended**, **Upper bounding**, and **Probing**. Table 1 lists the features in each set. Further, we define the set **All** as the union of **Basic**, **Basic extended**, **Upper bounding**, and **Probing**.

The **Basic** features are the number of variables n and the mean number of candidate parent sets per variable, m/n , which can be viewed as a natural measure of the “density” of an instance. The features in **Basic extended** are other simple features that summarize the size distribution of the collections \mathcal{G}_i and the candidate parent sets G_i in each \mathcal{G}_i . During training, we take the logarithm of the features related to the number of candidate parent sets (Features 2–5).

In the **Upper bounding** set, the features are characteristics of a directed graph that is an optimal solution to a relaxation of the original BNSL problem. Notice here especially the features based on strongly connected components (SCCs), which can be seen as a proxy for cyclicity.⁴ In the **Simple UB** subset, a graph is obtained by letting each variable select its best parent set according to the scores. The resulting graph may contain cycles, and the associated score is a guaranteed upper bound on the score of an optimal DAG. Many of the reviewed state-of-the-art solvers either implicitly or explicitly use this upper bounding technique; however, they do not use this information to estimate the difficulty

³ While, in principle, the function T_S also depends on external factors such as the specific hardware on which the solver is run, we do not consider those factors in this work.

⁴ Note that counting the number of cycles in a given graph is, in terms of computational complexity, presumably highly intractable, whereas SCC computation is achieved fast with well-known polynomial-time algorithms.

of a given instance. The features summarize structural properties of the graph: in- and out-degree distribution over the variables, and the number and size of non-trivial strongly connected components. In the **Pattern database UB** subset, the features are the same but the graph is obtained by solving a more sophisticated relaxation of the BNSL problem using the *pattern databases* technique [71]. Briefly, this strategy optimally breaks cycles among some subsets of variables but allows cycles among larger groups; it is a strictly tighter relaxation than the **Simple UB**. Both A* and CP explicitly make use of the pattern database relaxation.

Probing refers to running a solver for a fixed number of seconds and collecting statistics about its behavior during the run. Probing has previously been shown to be an important form of features, for example, in the context of Boolean satisfia-

Table 1 BNSL features

Basic

1. **Number of variables**
2. **Mean number of CPSs** (candidate parent sets)

Basic extended

- 3–5. **Number of CPSs** max, sum, sd (standard deviation)
- 6–8. **CPS cardinalities** max, mean, sd

Upper bounding

Simple UB

- 9–11. **Node in-degree** max, mean, sd
- 12–14. **Node out-degree** max, mean, sd
- 15–17. **Node degree** max, mean, sd
18. **Number of root nodes** (no parents)
19. **Number of leaf nodes** (no children)
20. **Number of non-trivial SCCs** (strongly connected components)
- 21–23. **Size of non-trivial SCCs** max, mean, sd

Pattern database UB

- 24–38. The same features as for *Simple UB* but calculated on the graph derived from the pattern databases

Probing

Greedy probing

- 39–41. **Node in-degree** max, mean, sd
- 42–44. **Node out-degree** max, mean, sd
- 45–47. **Node degree** max, mean, sd
48. **Number of root nodes**
49. **Number of leaf nodes**
50. **Error bound**, derived from the score of the graph and the pattern database upper bound

A probing*

- 51–62. The same features as for *Greedy probing* but calculated on the graph learned with A* probing

ILP probing

- 63–74. The same features as for *Greedy probing* but calculated on the graph learned with ILP probing

CP probing

- 75–86. The same features as for *Greedy probing* but calculated on the graph learned with CP probing

bility within the SATzilla portfolio approach [70]. Hutter *et al.* [34] survey the use of probing features in other domains. Here in the context of BNSL we consider four probing strategies: greedy hill climbing with a TABU list and random restarts, an anytime variant of A* [52], and the default versions of ILP [17] and CP [5]. All of these algorithms have anytime characteristics, so they can be stopped at any time and output the best DAG found so far. Furthermore, the A*, ILP, and CP implementations give guaranteed error bounds on the quality of the found DAGs in terms of the BNSL objective function; an error bound can also be calculated for the DAG found using greedy hill climbing by using the upper bounding techniques discussed above. Probing is implemented in practice by running each algorithm for 5 seconds and then collecting several features, including in- and out-degree statistics and the error bound. We refer to these feature subsets of **Probing** as **Greedy probing**, **A* probing**, **ILP probing**, and **CP probing**, respectively.

3.2 Model Training and Evaluation

In this work, we use the AUTO-SKLEARN system [20] to learn an explicit empirical hardness model \hat{T}_S for each solver S ; Briefly, AUTO-SKLEARN uses a Bayesian optimization strategy for learning good model classes and hyperparameters for those model classes for a given training set; additionally, preprocessing strategies, such as polynomial expansion or feature selection, and associated hyperparameters are included in this optimization. Importantly, this approach avoids the difficult step of manually choosing hyperparameters in an *ad hoc* fashion. We refer the reader to the original publication [20] for more details.

In total, AUTO-SKLEARN selects from amongst eleven preprocessing strategies, including higher dimensional projection techniques like polynomial expansion and feature selection strategies based on, for example, mutual information. The default learning strategy for AUTO-SKLEARN includes twelve model classes for regression and selects an ensemble of up to 50 regressors with optimized hyperparameters. In order to learn interpretable models and avoid potential overfitting, we restricted the use of AUTO-SKLEARN to learn the hyperparameters for a single preprocessor and random forest.⁵ As described in detail in Section 4.2, this study includes three types of BNSL instances: REAL, SAMPLED and SYNTHETIC. For model training, we used all of the three types of datasets.

The portfolios and prediction accuracy are evaluated using an “outer” 10-fold cross-validation scheme. In other words, the data is partitioned into 10 non-overlapping subsets. For each fold, nine of the subsets are used to train the model. As a first step in training, we normalize each feature so that it has zero mean and unit variance; the same mean and variance are later used to scale the test data. We then use AUTO-SKLEARN to learn the respective models. Internally, AUTO-SKLEARN further splits the training data in an “inner” cross-validation approach to avoid overfitting. We give 5 hours for training time for each fold. The remaining subset is used for testing, which only takes a few seconds; each subset is used as the testing set once. Importantly, the subset used for testing is not at all seen by AUTO-SKLEARN during training.

⁵ The choice of preprocessor was not restricted.

For testing, we predict the runtime of each testing instance using the appropriate model for each solver. For the algorithm selection analysis in Section 5.2, we then select the solver with the lowest predicted runtime. In order to accurately reflect the entire cost of algorithm selection, we report the runtime of a portfolio on a given instance as *the sum of the runtimes* of (i) feature computation for all feature sets used in the respective models and (ii) the selected solver.

4 Experimental Setup

We continue with a detailed description of our experimental setup, including descriptions of the solver parameterizations used, the data sets used in the experiments, as well as the computing infrastructure used.

4.1 Solvers

We begin by describing the exact parameterizations of complete BNSL solvers used in the experiments. Specifically, we evaluate three complete approaches: Integer-Linear Programming (ILP), A*-based state-space search (A*), and a constraint programming based approach (CP). Importantly, these approaches constitute the current state-of-the-art solvers for BNSL.⁶

We consider the following solvers and their parameterizations. We refer to all of the solvers for each approach as a *solver family*.

ILP We use the GOBNILP solver [17] as a state-of-the-art representative of the ILP-based approaches to BNSL. GOBNILP uses the SCIP framework [1] and an external linear program solver; we chose the open source SoPlex solver [69] bundled with the SCIP Optimization Suite. We consider the most recent version, GOBNILP 1.6.2, which uses SCIP 3.2.0 with SoPlex 2.2.0, as well as GOBNILP 1.4.1 (SCIP 3.0.1, SoPlex 1.7.1). For both versions we consider two parameterizations: the default configuration, which searches for BNSL-specific cutting planes using graph-based cycle finding algorithms, and a second configuration, “-nc” (“no cycle-finding”), which only uses nested integer programs. We call these parameterizations `ilp-141`, `ilp-141-nc`, `ilp-162`, and `ilp-162-nc`, respectively, for short.

A* We use the URLearning solver [72] as a state-of-the-art representative approach to BNSL based on the A* search method. We consider three parameterizations: A*-ed3, which uses dynamic pattern databases, A*-ec, which uses a combination of dynamic and static pattern databases, and A*-comp which uses a strongly connected component-based decomposition [18].

CP We use the CPBayes solver [5] as the most recent state-of-the-art representative approach to BNSL based on branch-and-bound style constraint programming search with problem-specific filtering (search-space pruning)

⁶ In a preliminary version of this work [51], we also considered an earlier proposed branch-and-bound approach [11], which we found to be always dominated by ILP; therefore, we dropped it from consideration. Furthermore, the earlier proposed dynamic programming approach [37] is clearly dominated by A*. We have also discarded some parameterizations of both ILP- and A*-based solvers that were found to be uncompetitive.

techniques. This solver does not expose any parameters to control its behavior, so we apply the solver in our experiments in its default configuration, `cpbayes`.

The non-default parameterizations of the solvers were suggested to us by the solver developers. While we use both an “up-to-date” version (1.6.2) and an older version (1.4.1) of GOBNILP, it is important to note that, generally, the choice of parameters and the solver version can at times have a noticeable effect on the per-instance runtimes of the resulting solver—so much so that one could consider the solvers different.⁷

4.2 Training Data

To train our models we first obtained a collection of datasets from various sources. For each dataset we then evaluated one or more scoring functions to produce a collection of BNSL instances. We used datasets from the following three categories.⁸

REAL Real-world datasets obtained from machine learning repositories: the UCI repository [2], the MLData repository (<http://mldata.org/>), and the Weka distribution [27]. We searched primarily for datasets of fully or mostly categorical data and a reasonable number of variables (16–64) to produce instances that are feasible but non-trivial to solve. Every dataset found and matching these criteria was included. While some of the datasets have originally been designed for supervised learning, they have been regularly included also in studies of unsupervised learning. These datasets are summarized in more detail in Table 9 of Appendix A.

SAMPLED Datasets sampled from benchmark Bayesian networks, obtained from <http://www.cs.york.ac.uk/aig/sw/gobnilp/>. These datasets are widely used for evaluating the performance of individual solvers, for example, recently in the context of optimal BNSL [4–6, 17–19, 51, 50, 60]. These datasets are summarized in Table 10 of Appendix A.

SYNTHETIC Datasets sampled from synthetic Bayesian networks. We generated random networks of varying number of binary variables (20–60) and maximum in-degree (2–8). For each network one dataset was produced by sampling a random number (100–10,000) of records.

We preprocessed each dataset by removing unique identifiers (to avoid overfitting) and trivial variables that only take on one value. Continuous variables as well as other variables with very large domains were either removed or discretized using a normalized maximum likelihood approach [38] when possible. The maximum number of records per dataset was limited to 60,000 to make the evaluation of scoring functions feasible.

⁷ For corroborating evidence on this, see, e.g., empirical data provided [17] for different parameterizations and versions of GOBNILP.

⁸ The main motivations for including both more real and, on the other hand, synthetic datasets in the study are two-fold: (i) We aimed at a notably heterogeneous set of benchmarks for the study, yielding insights into the prediction task on a wide range of datasets with different properties; and (ii) the three-way categorization has analogies in the benchmark categorization used in the SAT domain [36].

Table 2 Number of source datasets, instances generated from the source datasets, and instances used in training and testing the models.

Category	Datasets	All Instances	Training & Testing
REAL	39	637	486
SAMPLED	19	317	283
SYNTHETIC	477	477	410

We considered five different scoring functions:⁹ the BDeu score with the Equivalent Sample Size parameter selected from $\{0.1, 1, 10, 100\}$ and the BIC score. For each dataset in the REAL and SAMPLED categories we produced multiple instances by considering all scoring functions and varying upper bounds on the size of each candidate parent set, ranging from 2 to 6, as well as the unbounded case. For each dataset in SYNTHETIC we produced one instance, choosing both the scoring function and the parent limit at random. For larger datasets, evaluating the scores was feasible only up to lower values of the maximum parent set size. The total number of datasets and BNSL instances produced is summarized in Table 2.

For running all solvers on these instances we used a cluster of Dell PowerEdge M610 computing nodes equipped with two 2.53-GHz Intel Xeon E5540 CPUs and 32-GB RAM. For each individual run we used one CPU core, with a timeout of two hours and a 30-GB memory limit. We treat the runtime of any instance as two hours if a solver exceeds either the time or memory limit.

For training the models, we used a subset of all instances obtained by removing very easy instances, solved within five seconds by all solvers, as well as instances on which all solvers failed.¹⁰ We call these the *training* instances (see Table 2) and focus on them in the following sections.

4.3 Feature Computation

In order to train the models we computed the features detailed in Section 3.1 for all training instances. Table 3 summarizes the time spent to compute these features separately for each feature category. We observe that the computation takes around 16 seconds per instance on average and about 26 seconds in the worst case. Further, most of the time is spent on probing, while features of all other categories are computed in less than one second. In other words, a time limit needs to be enforced only for computing the probing features. As witnessed by the maximum feature computation times, probing occasionally exhibits higher running times than the limit of 5 seconds to finish a preprocessing step. This can be caused by overhead resulting, for example, from memory deallocation operations. We gave an additional 5 seconds for probing to finish on those specific instances. If the probing solver was still not completed within this time, it was terminated.

⁹ In our experiments, the results were not very sensitive to the scoring function, except through its effect on the number of candidate parent sets and other features, so our results can generalize to other decomposable scores as well.

¹⁰ This is in line with related work on portfolio construction in other domains such as SAT [34] as well as the SAT Competitions where a similar criterion is used to filter out “too easy” instances from the competition benchmark sets [3]. Solver selection for very easy instances is trivial, as any choice of a solver is essentially a good one.

Table 3 The runtime of feature computation for each feature category in seconds, shown as the average, median, minimum, and maximum runtime over all training instances.

Feature set	Average	Median	Min	Max
Basic	0.00	0	0	0
Basic extended	0.00	0	0	0
Lower bounding	0.00	0	0	0
Greedy probing	2.53	2	0	6
A* probing	4.61	5	0	7
ILP probing	3.94	5	0	10
CP probing	4.49	6	0	10
All	15.57	18	0	26

All in all, the overhead from computing the features is negligible from a portfolio perspective, as our main interest is in choosing the fastest solver for harder instances that take several minutes or even hours to solve. The easiest instances by contrast are often solved already in the probing phase.¹¹

4.4 Availability of Experiment Data

To facilitate open access and further analysis of the data produced in the experiments of this work, we have made the full solver runtime data, as well as the models learned for runtime prediction, available at

<http://bnportfolio.cs.helsinki.fi/>.

Furthermore, the runtime and feature data are available as a scenario in the ASlib Algorithm Selection Library [8] for further benchmarking purposes at

http://github.com/coseal/aslib_data/tree/master/BNSL-2016.

5 Portfolios for BNSL

This section focuses on the construction of practical BNSL solver portfolios in order to address question Q1. Optimal portfolio behavior is to always select the best-performing solver for a given instance. As the main result, we will show that, perhaps somewhat surprisingly, it is possible to construct a practical BNSL solver portfolio that vastly outperforms any single solver using only the **Basic** features.

5.1 Solver Performance

As the basis of this work, we ran all the solvers on all the BNSL instances, as described in Section 4. A comparison of solver performance is shown in Figure 3, in terms of the number of instances for which a particular solver was empirically faster than all other solvers on the considered benchmarks. Tables 4 and 5 show an alternative comparison in terms of the total number of instances that were successfully solved within the given computational resources as well as the total

¹¹ The benchmark set used was not filtered based on probing results.

CPU time required to either solve an instance or run out of time or memory. The results are given in comparison to the Virtual Best Solver (VBS), which is the theoretically optimal portfolio that always selects the best solver, constructed by selecting *a posteriori* the fastest solver for each input instance. Essentially, a theoretical lower bound on the runtime of any portfolio approach using a fixed set of k solvers is the runtime of the VBS. Furthermore, by interleaving the executions of the solvers until the best solver for a specific instance terminates, a theoretical upper bound of k times the runtime of the VBS is obtained.

We observe that among the ILP parameterizations, the two default configurations, `ilp-141` and `ilp-162`, are empirically best-performing on the considered benchmarks, while in terms of total runtime all four show fairly similar performance empirically. Among the A* parameterizations, A*-comp does best on average, while A*-ec outperforms A*-ed3 on nearly all instances and is also often the fastest parameterization in the REAL category, even though its total performance is worse than that of A*-comp.

In terms of the relative performance of the solvers, Figure 4 shows the pairwise correlations between the solvers on all instances. Unsurprisingly, different parameterizations within the same solver family correlate strongly with each other. Within the A* family, the strongest correlation is between A*-ec and A*-ed3, while all ILP parameterizations are strongly correlated, though mildly less so between different versions of the solver. Between solver families, A* and ILP correlate with each other the least, while CP exhibits mild correlation with ILP and moderate correlation with A*. Interestingly, A*-comp correlates more with CP than with the other A* parameterizations.

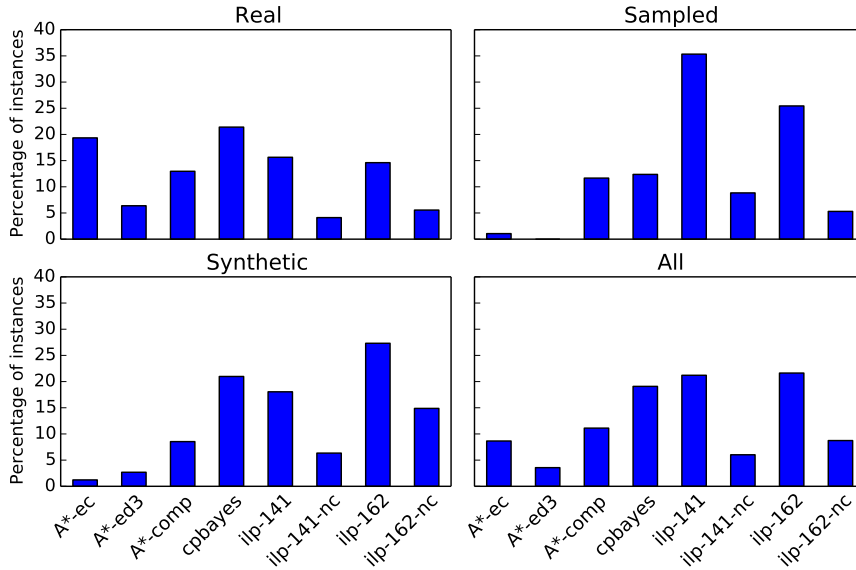


Fig. 3 The number of training instances for which a solver was fastest. Ties between solvers are broken at random.

Table 4 The performance of all solvers as well as the Virtual Best Solver (VBS) and four portfolios on all training instances, measured as the number of instances solved and the overall runtime. Instances that were not successfully solved within the given resources count as 7,200 seconds in the runtimes.

Solver	Instances solved (%)		Runtime (s)		
			Cumulative	Average	Median
VBS	1179	100	259,440	220	7.33
VBS without CP	1164	98	368,690	313	9.40
VBS without A*	1157	98	475,032	403	8.96
VBS without ILP	937	79	2,022,296	1,715	33.35
portfolio-basic	1141	96	540,384	458	12.30
autofolio-basic	1146	97	548,030	465	18.34
portfolio-all	1152	97	488,093	414	27.70
autofolio-all	1152	97	501,146	425	23.84
ilp-141	1036	87	1,364,855	1,158	36.39
ilp-141-nc	1034	87	1,384,022	1,174	41.83
ilp-162	1029	87	1,453,932	1,233	29.56
ilp-162-nc	1026	87	1,494,879	1,268	32.18
cpbayes	896	75	2,423,547	2,056	85.83
A*-comp	768	65	3,152,809	2,674	185.79
A*-ec	519	44	4,866,797	4,128	7,200.00
A*-ed3	478	40	5,163,876	4,380	7,200.00

While the ILP approach appears to be the best-performing measured in the total runtime and the number of instances solved on the set of benchmarks considered, the results suggest that the performance of ILP on a per-instance basis is quite orthogonal to that of both CP and A* (recall Fig. 1). We will now show that a BNSL solver portfolio can closely capture the best-case performance of *all eight* of the considered solver parameterizations in terms of empirical runtimes.

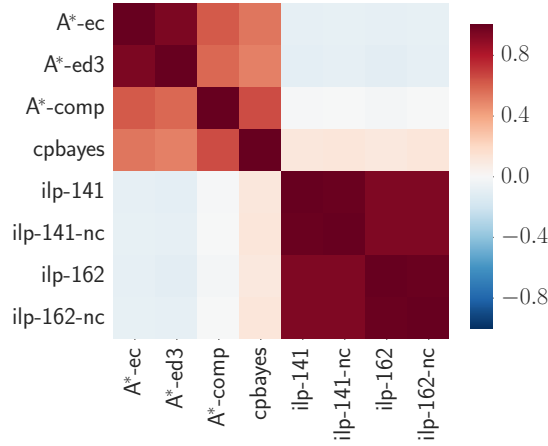


Fig. 4 Pairwise (Pearson) correlations between the runtimes of individual solvers.

Table 5 The performance of all solvers and portfolios within each instance category.

Solver	Solved (%)		Runtime (s)			Category
			Cumulative	Average	Median	
VBS	486	100	92,165	78	2.69	REAL
VBS without CP	480	98	141,833	120	5.13	
VBS without A*	469	96	244,625	207	5.60	
VBS without ILP	448	92	370,212	314	8.48	
portfolio-basic	470	96	209,490	178	4.60	
autofolio-basic	469	96	232,889	198	9.78	
portfolio-all	475	97	175,555	149	16.66	
autofolio-all	474	97	197,599	168	16.00	
ilp-141	396	81	800,432	679	55.68	
ilp-141-nc	396	81	799,734	678	56.78	
ilp-162	382	78	882,431	748	44.24	
ilp-162-nc	382	78	887,222	753	48.25	
cpbayes	427	87	549,230	466	14.85	
A*-comp	382	78	860,025	729	65.98	
A*-ec	311	63	1,300,350	1,103	156.30	
A*-ed3	281	57	1,523,034	1,292	523.43	
VBS	283	100	62,010	53	5.62	SAMPLED
VBS without CP	278	98	92,511	78	6.31	
VBS without A*	280	98	97,027	82	5.95	
VBS without ILP	227	80	453,422	385	33.07	
portfolio-basic	274	96	131,034	111	9.02	
autofolio-basic	277	97	123,468	105	17.15	
portfolio-all	278	98	115,254	98	23.97	
autofolio-all	280	98	97,502	83	19.08	
ilp-141	256	90	253,298	215	9.54	
ilp-141-nc	254	89	266,871	226	13.91	
ilp-162	257	90	280,990	238	13.57	
ilp-162-nc	252	89	309,674	263	15.07	
cpbayes	212	74	603,795	512	91.45	
A*-comp	182	64	749,656	636	145.95	
A*-ec	81	28	1,488,628	1,263	7,200.00	
A*-ed3	71	25	1,558,424	1,322	7,200.00	
VBS	410	100	105,264	89	14.98	SYNTHETIC
VBS without CP	406	99	134,346	114	16.15	
VBS without A*	408	99	133,380	113	15.90	
VBS without ILP	262	63	1,198,662	1,017	357.74	
portfolio-basic	397	96	199,860	170	25.44	
autofolio-basic	400	97	191,674	163	26.44	
portfolio-all	399	97	197,284	167	38.68	
autofolio-all	398	97	206,045	175	36.77	
ilp-141	384	93	311,125	264	45.16	
ilp-141-nc	384	93	317,417	269	50.39	
ilp-162	390	95	290,512	246	30.32	
ilp-162-nc	392	95	297,984	253	29.48	
cpbayes	257	62	1,270,522	1,078	758.21	
A*-comp	204	49	1,543,127	1,309	7,200.00	
A*-ec	127	30	2,077,819	1,762	7,200.00	
A*-ed3	126	30	2,082,419	1,766	7,200.00	

5.2 Portfolios for BNSL

As a main observation reported on in this section, we found that using only the **Basic** features (number of variables, n , and mean number of candidate parent sets, m/n) is enough to construct an efficient BNSL solver portfolio. We emphasize that, while on an intuitive level the importance of these two features may be to some extent unsurprising, such intuition does not directly translate into an actual predictor that would close-to-optimally predict the best-performing solver.

We create two portfolios that select a solver based on the runtime predictions from a random forest and preprocessor with hyperparameters optimized by AUTO-SKLEARN [20], as described in Section 3.2. We denote these portfolios (i) **portfolio-basic** and (ii) **portfolio-all**, using (i) the **Basic** features only and (ii) the full feature set, respectively, to make the algorithm selections.

Tables 4 and 5 show the performance of these two portfolios compared to each individual solver parameterization as well as the Virtual Best Solver. The reported portfolio runtimes include both the time required to run the selected solver and the time spent to compute the features used by the portfolio. Figures 5–8 present a more detailed view of portfolio performance, measured as the number of instances solved within a specific time, for the full benchmark set (All; Fig. 5), as well as the individual benchmark categories: REAL (Fig. 6), SAMPLED (Fig. 7), and SYNTHETIC (Fig. 8). Again, the time required to compute the necessary features is included in the solving time. We observe that **portfolio-basic** solves over 96% of the instances in the full benchmark set, with a cumulative runtime roughly twice that of the VBS. It also greatly outperforms every individual solver; the fastest solvers overall are the ones in the ILP solver family, which all solve 87% of the instances and are over five times slower than the VBS. The portfolio using only the **Basic** features is only slightly worse than **portfolio-all**, which solves a handful more instances and has a somewhat lower cumulative runtime. The difference between the two portfolios is more pronounced within the REAL and SAMPLED categories, while within SYNTHETIC their performance is almost equal. This is presumably due to both portfolios heavily leveraging the ILP family, which alone exhibits very good performance in SYNTHETIC, solving 95% of the instances.

For understanding the marginal contributions of the considered solvers, we consider the Shapley value [61] as a measure for the contribution of a specific solver to a portfolio, following Fréchette *et al.* [22]. In this framework, one considers constructing a portfolio by adding solvers incrementally and measuring the value of each solver as the increase in the portfolio’s performance when the solver is added. As these values greatly depend on the order in which solvers are added, the Shapley value of a solver is defined as its average value over all possible solver permutations. Table 6 shows the Shapley values for all solver parameterizations, using the total number of instances solved as the measure of portfolio performance. Within each of the solver families, we observe that **ilp-162**, **cpbayes**, and **A*-comp**, respectively, have the highest Shapley values on the considered benchmarks.

Given the good runtime performance of the portfolios obtained using runtime predictions from random forests as the underlying algorithm selection strategy, it is interesting to investigate to what extent the choice of algorithm selection strategy impacts portfolio performance using the same set of BNSL features. For comparison, we consider AUTOFOLIO [47], a state-of-the-art algorithm selection

Table 6 The contribution of each solver to the VBS and the two portfolios measured as the Shapley value in terms of the average number of additional instances solved after adding the indicated solver to the portfolio.

Solver	VBS	portfolio-all	portfolio-basic
ilp-162	184.53	181.82	178.75
ilp-141	184.12	179.78	181.86
ilp-141-nc	182.48	179.62	178.79
ilp-162-nc	181.50	178.96	177.44
cpbayes	160.42	152.18	149.37
A*-comp	136.24	131.60	127.62
A*-ec	78.28	77.72	77.10
A*-ed3	71.43	70.34	70.08

system¹², for constructing the portfolios **autofolio-basic** (using AUTOFOLIO on the **Basic** feature set) and **autofolio-all** (using AUTOFOLIO on the full feature set).¹³

AUTOFOLIO [47] trains a binary classifier for each pair of solvers which selects the better-performing for a given instance; the instances are weighted based on the difference in performance for the two solvers. Further, AUTOFOLIO selects among the feature sets to use during testing to minimize the overall solution time. A Bayesian optimization strategy is used to optimize the classifier hyperparameters, feature set and preprocessing choices.¹⁴ For an unseen instance, each of the trained classifiers votes for a solver; the solver with the most votes is used for that instance. The training and testing splits were the same for both AUTOFOLIO and AUTO-SKLEARN. For AUTOFOLIO, we also used an “outer” 10-fold cross-validation scheme to ensure it does not use testing instances during training.

The two portfolios produced by AUTOFOLIO perform very similarly on the benchmark set as those based on predicting runtimes with random forests. In more detail, **autofolio-all** solves more instances than **portfolio-all** within the first thirty seconds for all instance types; this is because AUTOFOLIO does not always use all of the feature sets, so it spends less time computing features during test time. After this initial phase, the number of instances solved under a given per-instance timeout was very similar for **portfolio-all** and **autofolio-all**. As Table 4 shows, though, in total, **portfolio-all** has a slightly lower cumulative runtime than **autofolio-all**; the detailed breakdown in Table 5 clarifies that this is largely due to better performance of **portfolio-all** on the REAL instances.

On the other hand, **portfolio-basic** solves more instances than **autofolio-basic** in the thirty second time limit. Indeed, **portfolio-basic** consistently outperforms *all* of the other portfolios and individual solvers within this time limit for all instance types. Eventually, **autofolio-basic** solves 5 more instances than **portfolio-basic**, albeit with a higher average and median runtime. In total, we do not see significant differences between the portfolios based on AUTOFOLIO and AUTO-SKLEARN. This may be at least partially due to the fact that, internally, they both use the SMAC Bayesian optimization engine [33] and similar model classes and preprocessors.

¹² In particular, we use an updated version recommended by the author, <https://github.com/mlindauer/AutoFolio>.

¹³ We thank an anonymous reviewer for proposing this comparison with AUTOFOLIO.

¹⁴ The AUTOFOLIO implementation includes a pre-solving component [29]. We disabled that feature for purposes of this comparison in order to strictly consider how well the models capture solver behavior; however, a similar strategy could be used to include a pre-solver for the AUTO-SKLEARN-based approach, as well.

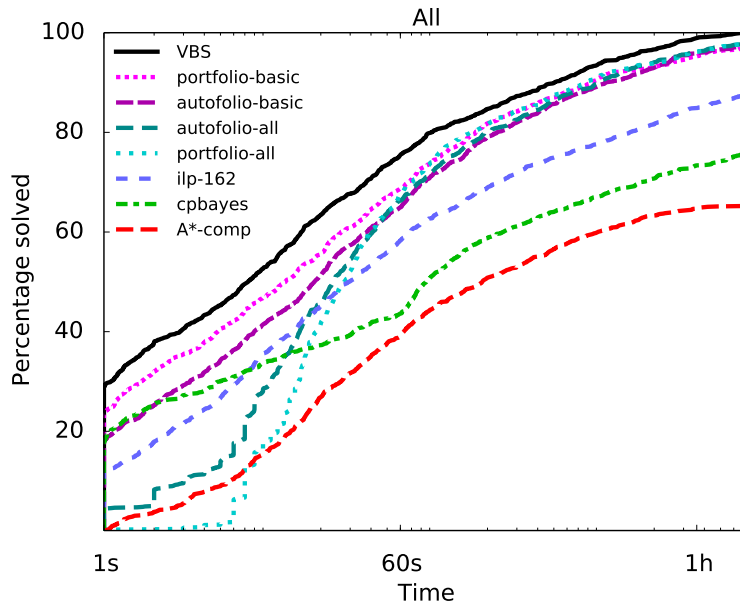


Fig. 5 Fraction of instances solved by the VBS, the portfolios, and individual solvers within a given amount of time.

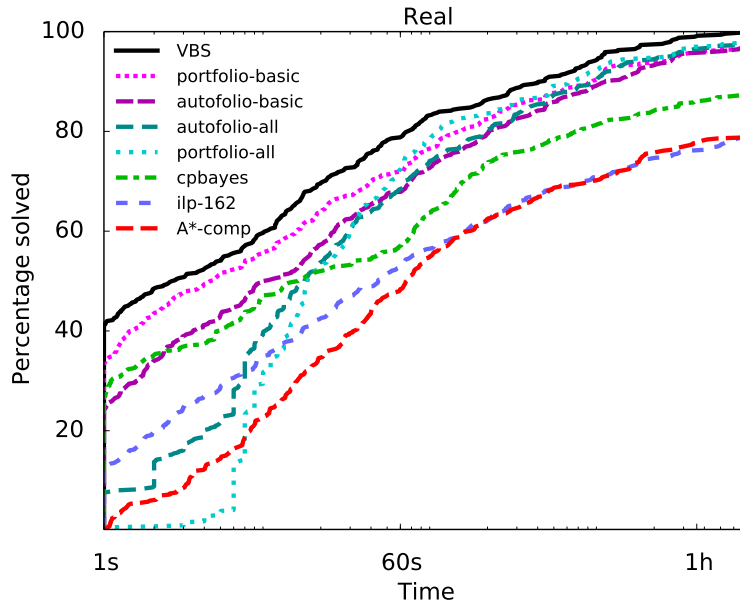


Fig. 6 Fraction of instances of the REAL category solved by the VBS, the portfolios, and individual solvers within a given amount of time.

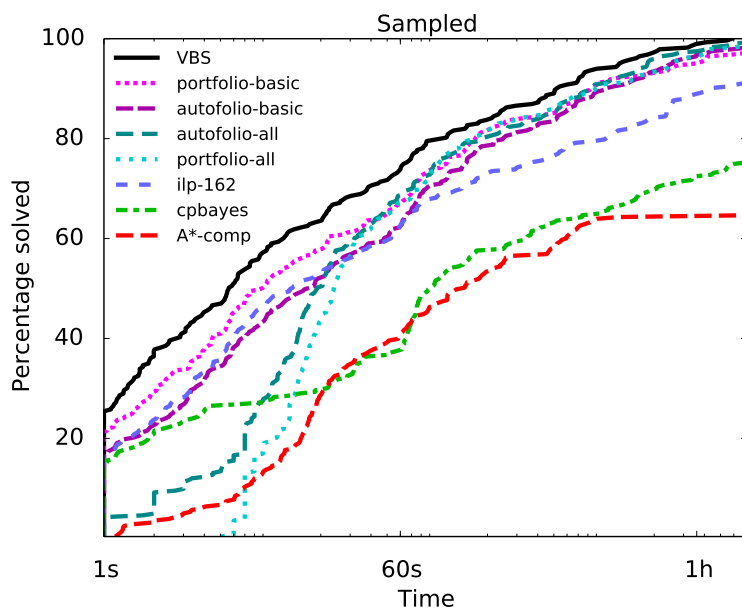


Fig. 7 Fraction of instances of the SAMPLED category solved by the VBS, the portfolios, and individual solvers within a given amount of time.

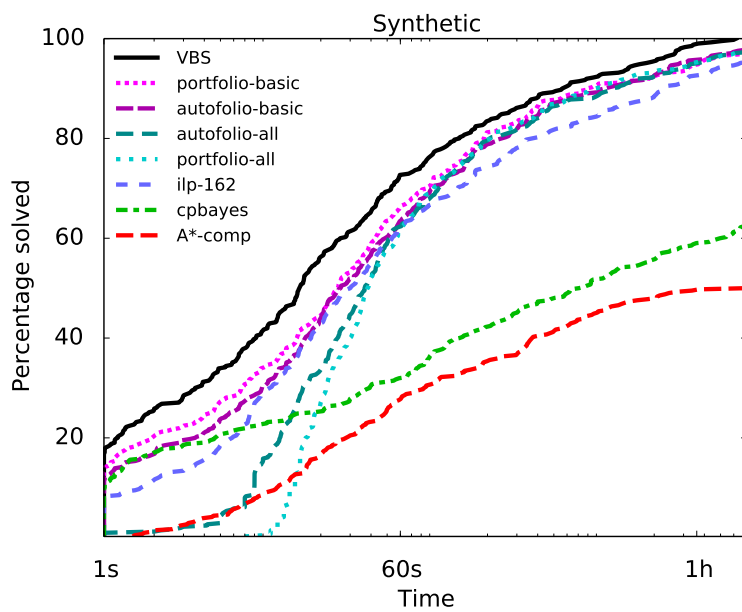


Fig. 8 Fraction of instances of the SYNTHETIC category solved by the VBS, the portfolios, and individual solvers within a given amount of time.

5.3 Basic Features and Solver Performance

As the **Basic** features yield efficient BNSL portfolios, we look more closely at the effect of the per-instance **Basic** feature values on solver performance. Figure 9 reinforces the orthogonal strengths of different solver families in the space spanned by these two features. Specifically, we observe that ILP parameterizations can fairly reliably solve instances up to around 1,000 candidate parent sets per variable, regardless of the number of variables. In comparison, the A* family consistently solves benchmark instances up to 30 variables, and many up to 40, even with tens of thousands of candidate parent sets per variable. Our results show that CP takes a middle ground between the two, solving many instances at the high end of either of the **Basic** features, albeit less consistently than either A* or ILP.

In particular, Figure 9 (top left) demonstrates why the **Basic** features result in strong portfolio behavior; namely, the instances which are optimally solved by the different solver families are nearly linearly separable in this space. The figure also supports the rough characterization (recall Section 1) of the computational limitations of state-of-the-art solvers: none of the state-of-the-art solvers are able to solve the benchmark instances where both of the **Basic** features are very large.

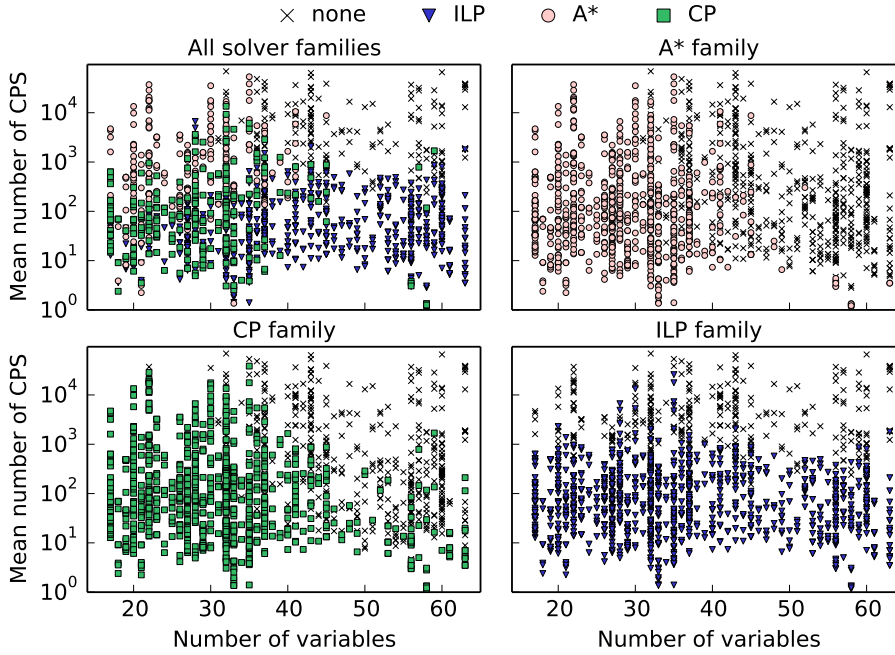


Fig. 9 All benchmark instances plotted in the space of the two **Basic** features, the number of variables and the mean number of candidate parent sets (CPS). Each instance is marked according to which solver was the fastest to solve it, specifically, whether the fastest solver was from the A*, ILP, or CP family, or whether none of the solvers could solve the instance. The comparison is presented for all solver families together (top left) and individually for each single family, highlighting their limitations as either or both features grow too large.

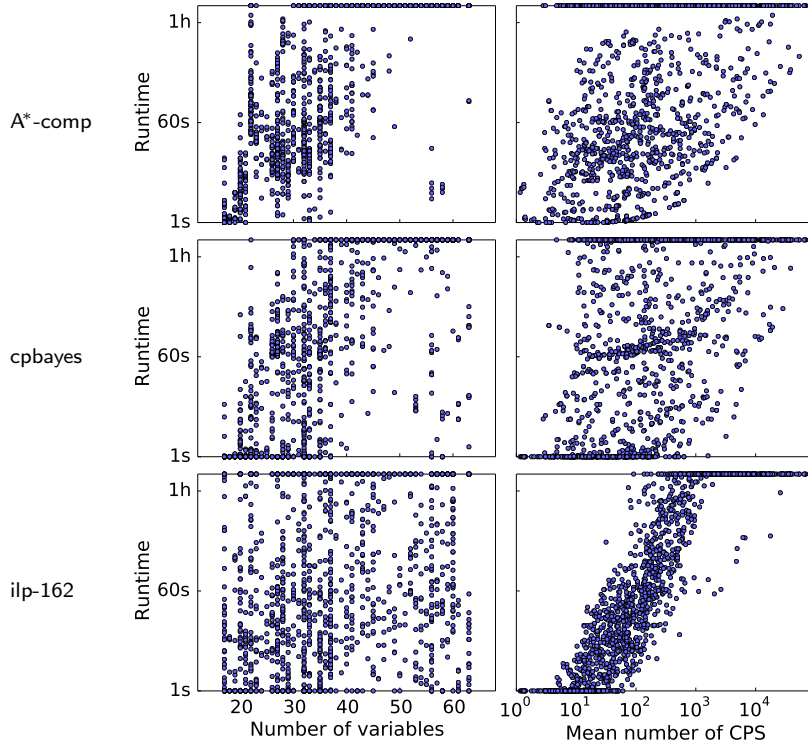


Fig. 10 Relationship between the **Basic** features, the number of variables and the mean number of candidate parent sets (CPS), and the runtimes of solvers.

Finally, we look deeper into the relationship between each feature independently and the specific solvers. Here we focus on **A*-comp**, **cpbayes**, and **ilp-162** since they have the highest Shapley value within the respective solver families for **portfolio-all**; we observed very similar trends for all solvers in each solver family. Figure 10 illustrates that the runtimes for **ilp-162** and the number of candidate parent sets are strongly related (coefficient of determination, that is, explained variance, $R^2 \approx 0.78$).¹⁵ On the other hand, the number of variables better explains the variance in the runtimes of **cpbayes** ($R^2 \approx 0.39$) and **A*-comp** ($R^2 \approx 0.47$). Conversely, **ilp-162** appears not to depend heavily on the number of variables ($R^2 \approx 0.0004$), while **A*-comp** and **cpbayes** seem able to solve instances irrespective of the number of candidate parent sets ($R^2 \approx 0.01$, $R^2 \approx 0.09$, respectively).

6 Predicting Runtimes

In this section, we turn to the arguably harder problem of predicting per-instance runtimes of individual solvers. Apart from pure scientific interest, accurate runtime

¹⁵ R^2 ranges from 0 to 1, where 0 indicates that the feature is completely uninformative about runtime, and 1 indicates that all of the variance in runtime is explained by the respective feature.

predictions on a per-instance basis are useful for job schedulers as computing clusters often require an estimated job time. In our case specifically, such predictions could also facilitate development of improved BNSL solvers. For example, a model could be exploited as a heuristic estimate for subproblem hardness during search within a parallel BNSL solver. As a further motivation, model-based algorithm configuration [33] crucially relies on runtime predictions in order to guide search for better configurations in the algorithm configuration space. In such contexts, note also that runtime is a primary resource to predict, as running out of other resources such as memory directly imply running out of time as well.

As shown in Section 5, the **Basic** features can effectively distinguish between solvers to use on a particular instance of BNSL. We will now address question Q2, that is, whether the use of additional features (cf. Section 3.1) improves the accuracy of the runtimes predicted by the random forests learned with AUTO-SKLEARN.

6.1 Predictions with Added Features

Figure 11 depicts the actual runtimes of solvers compared to the runtimes predicted by the random forests learned with AUTO-SKLEARN. We again use **A*-comp**, **cpbayes**, and **ilp-162** as representatives of their solver families (recall Section 5.2; similar conclusions hold for all solvers within the respective families). On the left we see this comparison for models trained using the **Basic** features only. Even though these predictions allow for good portfolio behavior, the considerable amount of prediction error makes them less useful for obtaining accurate estimates of the runtime. The right side, on the other hand, shows the same comparison when using **All**, where the predictions are more concentrated near the diagonal. In other words, the larger, more sophisticated feature set results in more accurate runtime predictions. Table 7 presents a numerical measure of the improvement in terms of change in the approximation factor, defined as $\rho = \max\{\frac{a}{p}, \frac{p}{a}\}$, where a and p are the actual and predicted runtimes, respectively. In particular, smaller approximation factors are better.

Additionally, we show the coefficient of determination (R^2) values of the predictions in Table 8. These values show that the observed variances in the actual runtimes are well-explained by the predictions. As expected, R^2 is always higher (better) when using **All** features compared to only the **Basic** ones. This offers another view which shows that the more sophisticated features improve prediction accuracy.

Table 7 The percentage of instances with an approximation factor within the given ranges of ρ , when predicting runtimes based on either **Basic** or **All** features. Higher percentages with lower approximation values indicate more accurate predictions.

Range of ρ	A*-comp		cpbayes		ilp-162	
	Basic	All	Basic	All	Basic	All
< 2	48%	60%	45%	67%	59%	71%
[2, 5)	22%	22%	27%	20%	29%	22%
[5, 10)	14%	7%	13%	7%	7%	4%
> 10	17%	11%	15%	6%	4%	3%

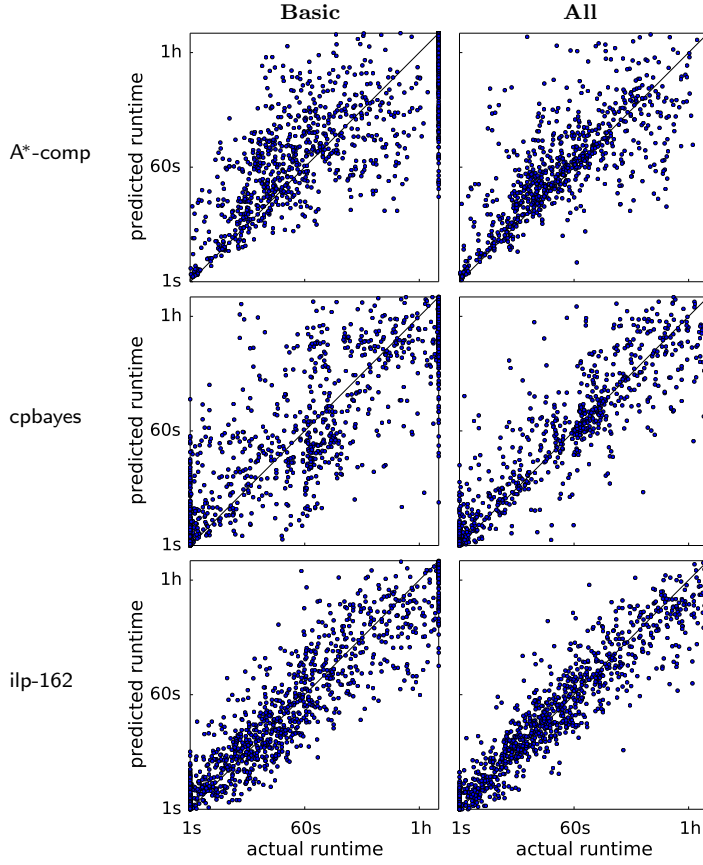


Fig. 11 The actual runtimes of solvers compared to the predicted runtimes when using **Basic** (left) or **All** (right) features.

Table 8 The coefficient of determination (R^2) for the actual runtime given the predicted runtime

Solver	A*-ec	A*-ed3	A*-comp	cpbayes	ilp-141	ilp-141-nc	ilp-162	ilp-162-nc
Basic	0.71	0.79	0.57	0.51	0.67	0.69	0.73	0.72
All	0.86	0.89	0.66	0.65	0.76	0.78	0.81	0.79

We also evaluated the impact of incrementally adding sets of features. Figures 12 and 13 show how the prediction error changes as we add **Basic** (features 1–2), **Basic extended** (1–23), **Upper bounding** (1–38), the relevant probing features for A* (1–38, 51–62), CP (1–38, 75–86), and ILP (1–38, 63–74), and finally **All** (1–86) for every solver. The results show that predictions using the **Basic** features are typically worse than those incorporating the other features, although this behavior is more pronounced for some solvers, feature sets and instance categories than others. The plots also suggest that some features help more than others for the different solvers. For instance, **Upper bounding** features greatly improve the predictions of A* compared to the **Basic** and **Basic extended** fea-

tures. In hindsight, this is relatively unsurprising since the efficacy of the upper bounding directly impacts the performance of A^* , showing that AUTO-SKLEARN effectively exploits features we intuitively expect to characterize the empirical hardness. Probing offers a glimpse at the true runtime behavior of the algorithms, and AUTO-SKLEARN leverages this information to further improve prediction accuracy. For both A^* and ILP, probing with the respective solvers alone is informative, while the other probing strategies (**All** features) yield little improvement and even weaken some of the predictions. In contrast, surprisingly, for CP the predictions modestly benefit from probing with other solvers as well. Out of the three solver families CP predictions improve most from added features in general.

Finally, we evaluate the root mean squared error (RMSE) of the predictions for each solver as we incrementally add feature sets. Figure 14 echoes the results from Figures 12 and 13. We again see that **Upper bounding** improves predictions on all A^* parameterizations. The respective probing features greatly improve the prediction accuracy for A^* -ec and A^* -ed3; relevant probing modestly improves the accuracy for the other solvers, as well.

6.2 Preprocessing Characteristics

We now turn to more qualitative analysis based on the preprocessor and single random forest with optimized hyperparameters learned by AUTO-SKLEARN.

First, we examine preprocessor choices. As shown in Figure 15, the choice of preprocessor often reflects the amount of information inherently available in the feature sets. Furthermore, Figure 15 includes a clustering of the solvers and feature sets based on the choice of preprocessor. In the clustering, we see that the families of solvers tend to cluster together.

The **Basic** feature set (dark tan) almost always result in a preprocessor which increases the dimensionality, either the polynomial expansion or random forest embedding technique; we interpret this to mean that the features alone do not provide sufficient information for accurate prediction, so AUTO-SKLEARN attempts to increase the information with preprocessing. Likewise, many of the “mildly informative” feature sets, such as **Simple UB** (dark teal), almost exclusively result in polynomial expansion for preprocessing the input features. Interestingly, the **Basic extended** feature set (light tan) results in polynomial expansion, a dimensionality expansion strategy, and feature agglomeration, a dimensionality reduction strategy, in roughly equal proportions for all solvers.

On the other hand, for the A^* algorithms with the larger feature sets like **All** (light brown), AUTO-SKLEARN has “too much” information, so it uses feature aggregation, as well as model-based and percentile-based feature selection, to combine or remove uninformative features; these choices typically are statistically significant. Preprocessing is usually not used for predicting most of the ILP runtimes using “informative” feature sets, such as **All** and **ILP Probing** (light teal); again, almost all of these choices are statistically significant.

This analysis demonstrates that the choice of preprocessing strategy by AUTO-SKLEARN largely agrees with intuition. For small, relatively uninformative feature sets, feature expansion strategies like polynomial expansion are often used; when more informative features are available, they are relatively unchanged. Finally,

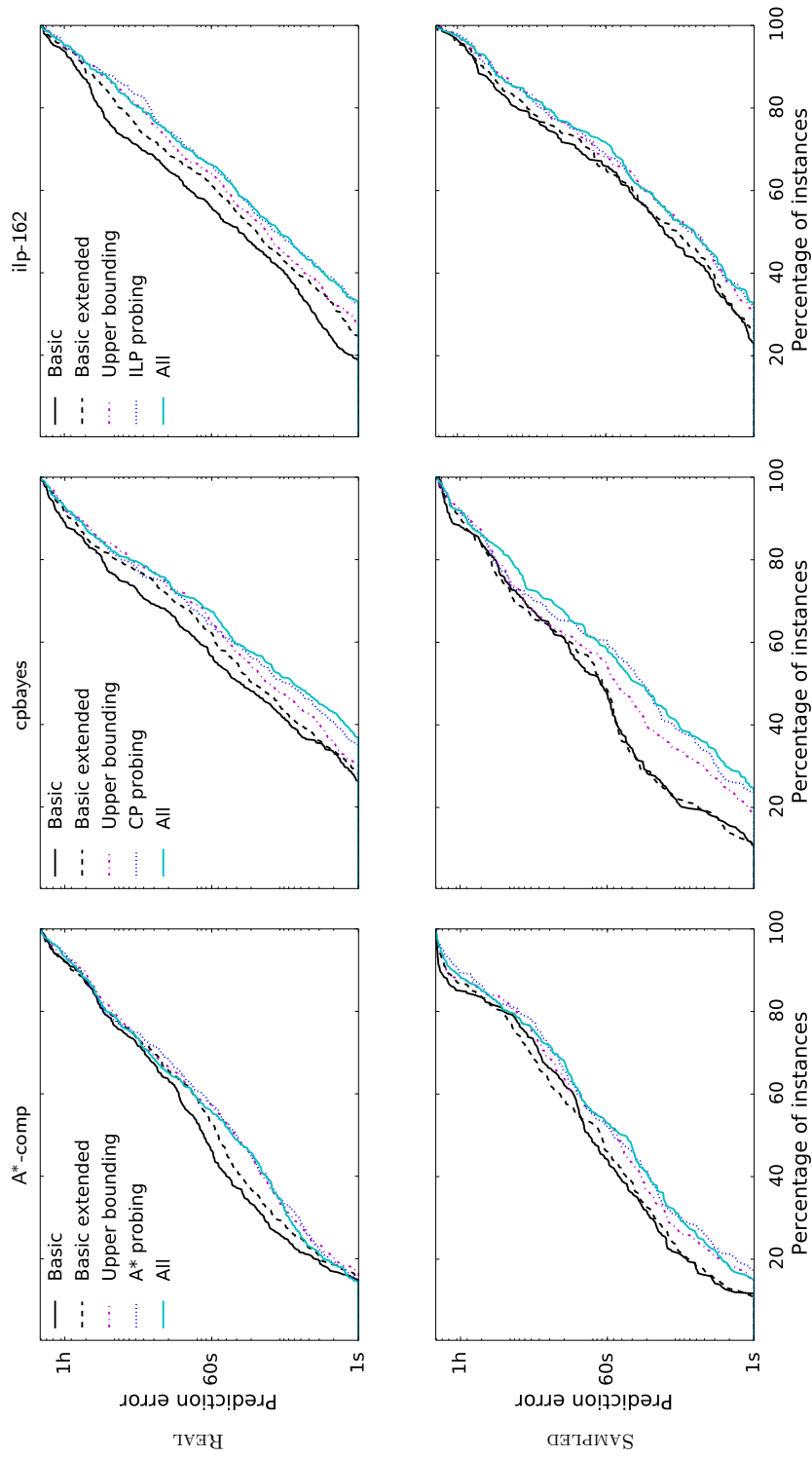


Fig. 12 The absolute prediction errors on REAL and SAMPLED instances using different sets of features, sorted in increasing order.

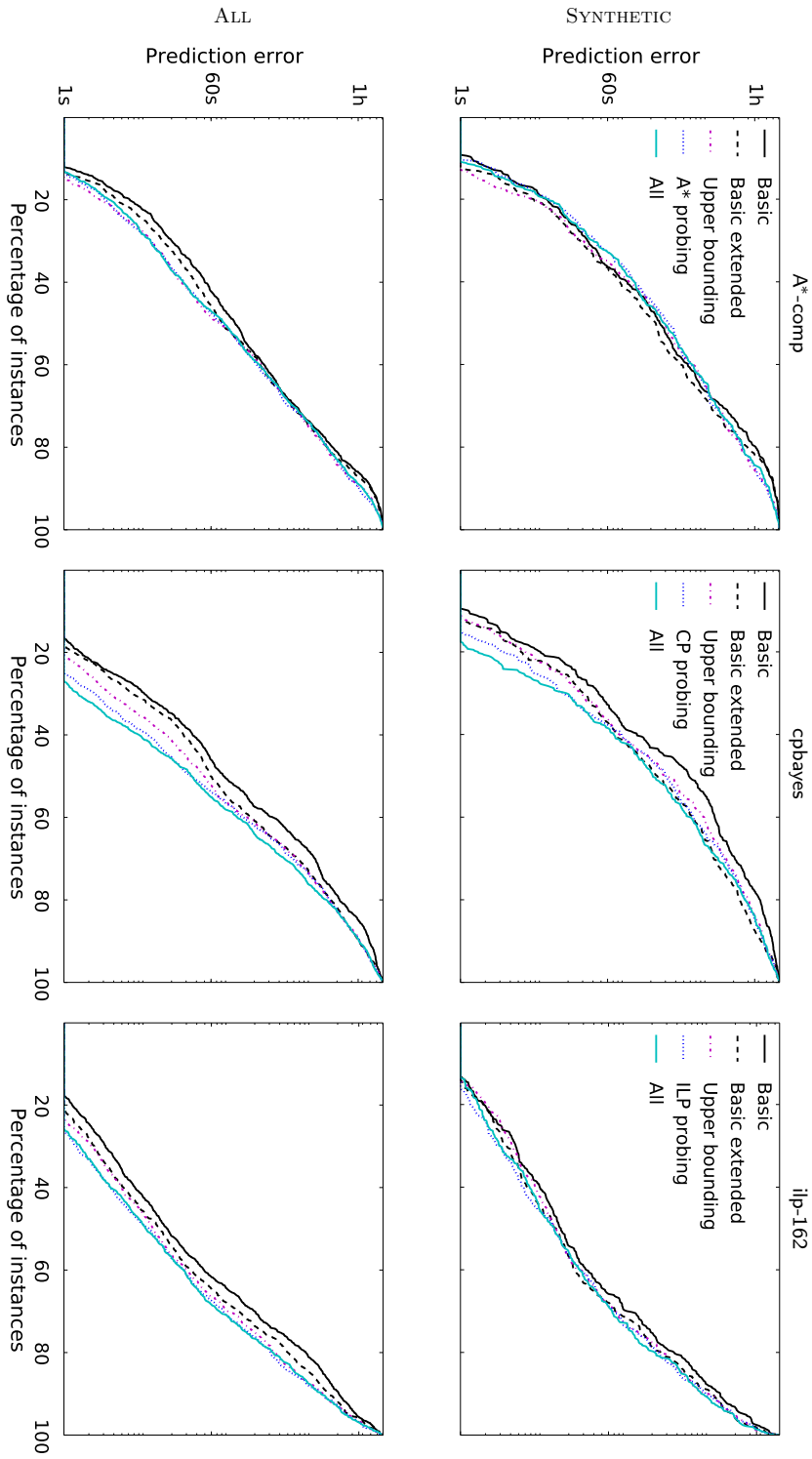


Fig. 13 The absolute prediction errors SYNTHETIC and ALL instances using different sets of features, sorted in increasing order.

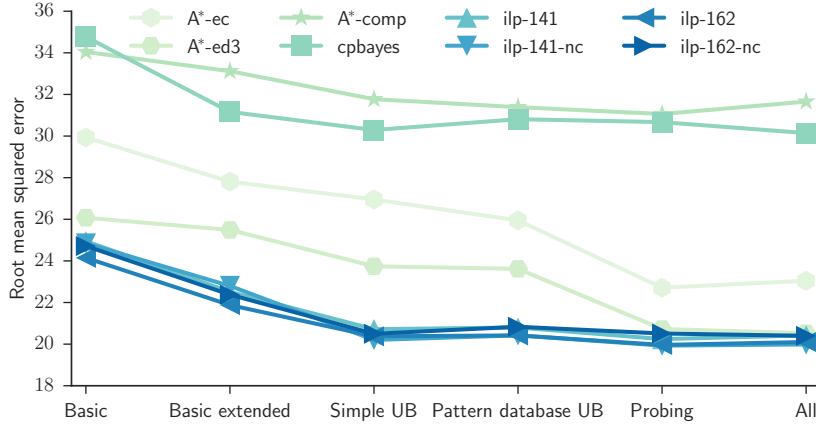


Fig. 14 The improvement of the root mean squared error of the runtime predictions as the more sophisticated features are used. “Probing” refers to the appropriate probing feature set for the respective solver, such as **A* probing** for the A*-ec solver.

when “too much” information is present, sophisticated feature selection strategies are used to retain useful features while removing noise.

6.3 Model Complexity

We additionally analyzed the complexity of the learned random forests, in terms of the mean size of the regression trees composing them. As expected, Figure 16(a) shows that the trees learned using the **Basic** features are the smallest. Other simpler feature sets, such as **Basic extended** and **Simple UB** also resulted in small trees for all solvers.

Somewhat surprisingly, though, the regression trees for the various ILP solvers are much larger than those for the **cpbayes** and **A*** family of solvers for the **A* probing**, **Pattern database UB**, **All** and **CP probing** feature sets. As shown in Figure 15, AUTO-SKLEARN often forewent preprocessing in these cases for ILP. On the other hand, it used sophisticated preprocessing, like the model-based approach, for **A*** and **cpbayes** a significant amount of the time. Thus, these results suggest an implicit tradeoff in AUTO-SKLEARN between resources used for preprocessing and the model itself.

Also unexpectedly, the trees for ILP without the graph-based cutting plane routines (the “-nc” parameterizations) are much larger than those using it with the **ILP probing** feature set. We hypothesize this is due to differences in the ILP implementation used for probing and the “-nc” solvers; namely, the ILP implementation used in probing *does* use the graph-based cutting plane routines. AUTO-SKLEARN uses preprocessing only sparingly in all of these cases, so it again appears that a more complex model is used to handle the noise in the features.

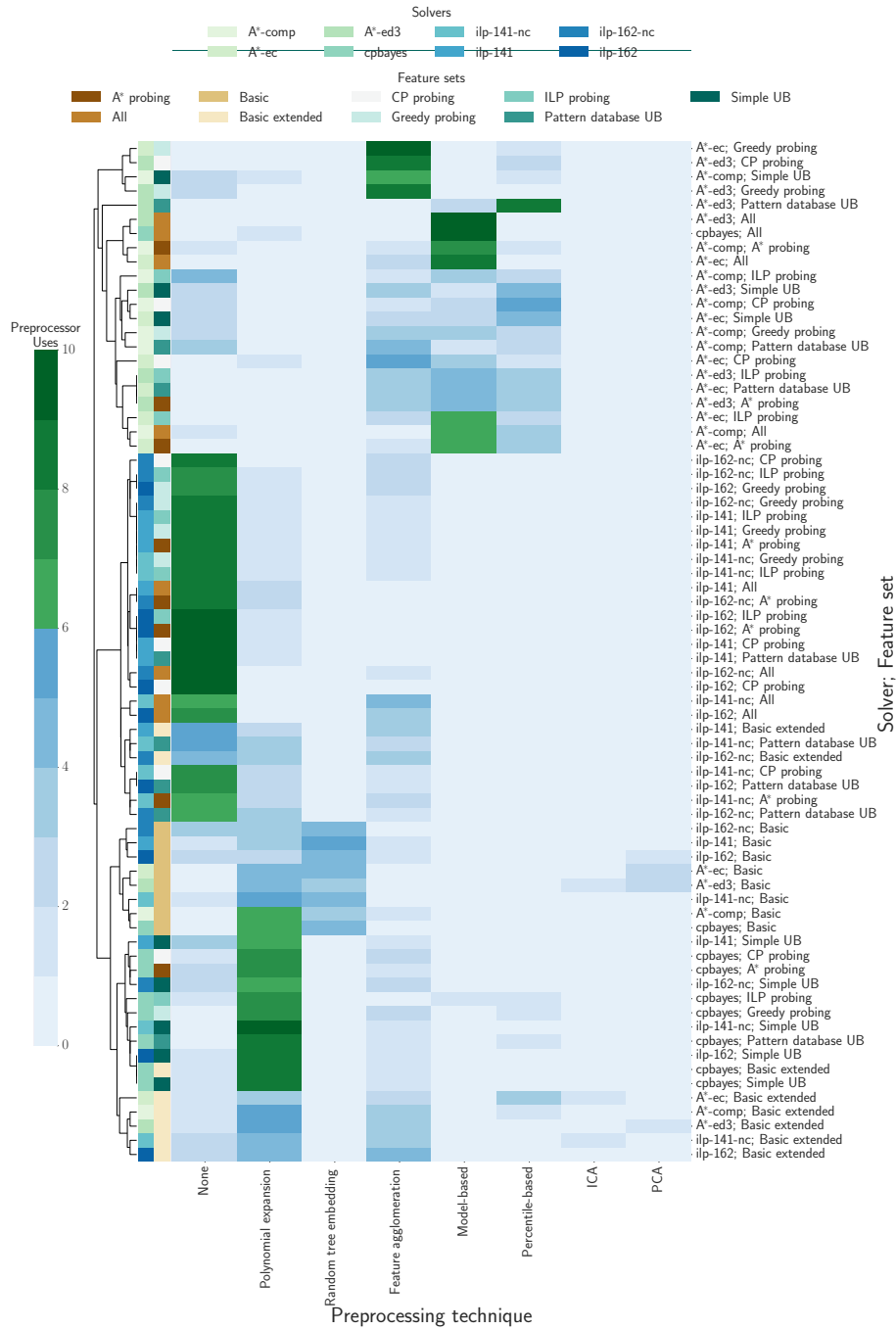


Fig. 15 The preprocessing techniques used by AUTO-SKLEARN for each combination of solver and feature set. The blue-green column of colors on the left indicate the solver in that row, and the green-brown column indicates the feature set; the text on the right also gives this information. Each cell shows the number of times the respective preprocessing technique was selected in one of the 10 cross-validation folds for the associated (solver, feature set) pair. The UPGMA algorithm [64] with a Euclidean distance metric was used for clustering. Cells shaded in green indicate statistically significantly high choices ($p < 0.01$, one-sided binomial test comparing to a uniform distribution, Benjamini-Hochberg multiple test correction).

6.4 Important Features

Finally, we computed the Gini importance [9] of each feature for predicting each solver while using the appropriate **Probing** features. The importance for a particular feature is calculated using a standard two-step technique [9]. First, the feature is corrupted with noise to create a new dataset. Then, the new dataset is used for training and testing as usual. The normalized increase in error when using the noisy feature is taken as its importance. For the random forests, this procedure is performed for all trees in the forest. The feature importance is then the average across all trees. Finally, we average the feature importances across each cross-validation fold.

Figure 16(b) shows important features for the different solvers. Several of the importances are unsurprising; the number of variables in the dataset determines the size of the search space for A^* , and that was the most important feature for all parameterizations. Similarly, the size of the linear program solved by ILP is directly determined by the number of candidate parent sets, and its most important features describe these sets. Likewise, the respective probing error bound features were typically somewhat important for ILP and CP. This is sensible because these features indicate when a solver can quickly converge to a nearly-optimal solution; however, as could be seen from Figure 14, the overall improvement to RMSE is modest with the addition of the probing features.

Figure 16(b) shows that the CP and A^* family models share many important features. For example, CP uses the pattern database relaxation which also guides the A^* search, and pattern database node degree features are indeed important for both CP and A^* models.

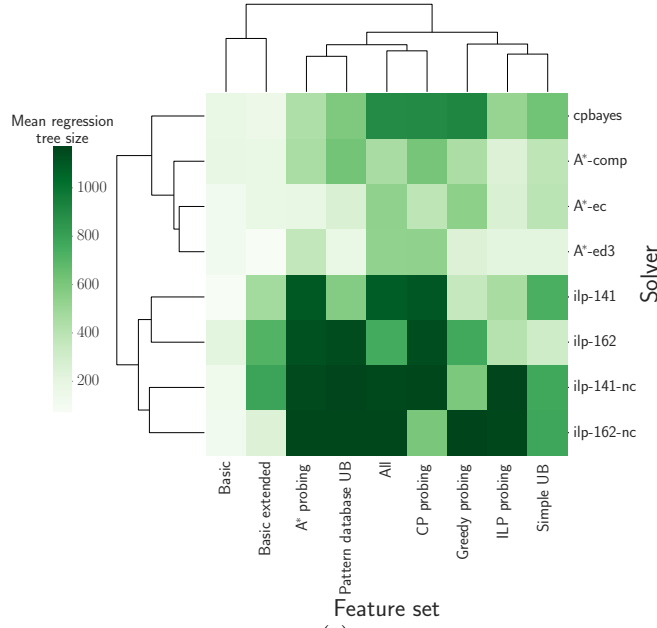
In contrast to ILP and CP, A^* -comp is the only A^* parameterization for which probing was an important feature. Coupled with the minimal improvement to RMSE shown in Figure 14 when using probing, this suggests that the runtime characteristics of the anytime variant of A^* are different enough from the A^* family of solvers included in the portfolio that it adds significant noise to learning.

Another somewhat unexpected result concerning A^* is that many **Simple UB** features are quite important. Previous experimental results [72] show that the pattern database bounding approach is much more informative *during* the A^* search. However, the solvers construct their pattern databases differently than those used for extracting features, so the structural properties, such as the number of non-trivial SCCs, of the constructed graphs may not reflect the difficulty of the problem for the solver.

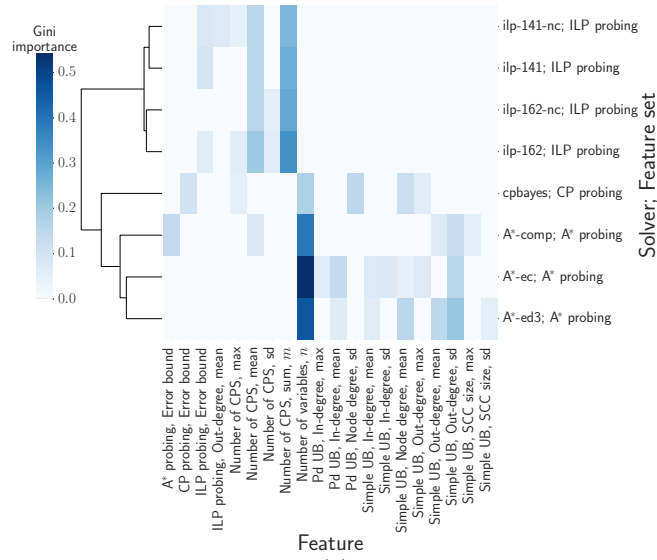
In general, the results presented in Figure 16(b) reveal that a small number of features were consistently important for any particular solver; this is in line with previous work [42, 44]. Qualitatively, this implies that most of the trees were based on the same small set of features.

7 Conclusions

We have investigated the *empirical hardness* of BNSL, the Bayesian network structure learning problem, in relation to several state-of-the-art complete solvers based on A^* search, integer linear programming, and constraint programming. While each of these solvers always finds an optimal Bayesian network structure (with respect



(a)



(b)

Fig. 16 (a) The average size of the regression trees in the random forests learned by AUTO-SKLEARN for each solver and feature set. (b) The Gini importance [9] of features in the learned random forest models for each solver using the respective **Probing** feature set. Only features with an importance of at least 0.05 for at least one solver are included. We use the abbreviations “CPS” for candidate parent sets, “Pd” for pattern database, and “sd” for standard deviation. The UPGMA algorithm [64] with a Euclidean distance metric was used for clustering in both cases; the features in (b) were not clustered.

to a given scoring function), the runtimes of the solvers can vary greatly even within instances of the same size. Moreover, on a given instance, some solvers may run very fast, whereas others require considerably longer time, sometimes by several orders of magnitude. We validated this general view, which has emerged from a series of recent studies, by conducting the most elaborate evaluation of state-of-the-art solvers to date. We have made the rich evaluation data publicly available¹⁶ in order to facilitate possible further analyses that go beyond the scope of the present work.

As the second contribution, we applied machine learning methods to construct *empirical hardness models* from the data obtained by the solver evaluations. Instantiating the general methodology of empirical hardness models [58, 46], we proposed several *features*, that is, real-valued functions of BNSL instances, which are potentially informative about solver runtimes and which go beyond the basic parameters of instance size.

We used two approaches, AUTO-SKLEARN and AUTOFOLIO, for building BNSL portfolio solvers, to directly address the algorithm selection problem. Additionally, we studied in more detail the runtime prediction accuracy of the models learned with AUTO-SKLEARN. Both of these state-of-the-art systems use Bayesian optimization to optimize model class, preprocessing and relevant hyperparameters, for the respective models.

The learned models allowed us to answer two basic questions concerning prediction of the solvers’ relative and absolute performance without actually running the solvers. The first question (Q1) asked whether the basic parameters of input size suffice for reliably predicting which of the solvers is the fastest on a given problem instance. We answered this question in the affirmative by showing that whenever a solver is significantly slower than the fastest solver on a given instance, the slower one is very rarely predicted as the fastest one. We compared the performance of portfolios based on models learned by both AUTOFOLIO and AUTO-SKLEARN, and observed that these two approaches yielded very similar portfolio runtime performance. For varying distributions of instances, our portfolio solver using a very basic set of BNSL features resulted in the fastest solver overall, exhibiting cumulative runtimes within two times that of the Virtual Best Solver (VBS). In contrast, the cumulative runtime of the best individual solver is over five times that of the VBS. As a result, the proposed solver portfolio is currently the fastest algorithm for solving BNSL when averaged over a large heterogeneous set of instances.

Our answer was affirmative also to the second question (Q2) of whether the runtimes of each of the solvers can be predicted more accurately by extending the set of features. We observed that, in general, the more high-quality the features, the more accurate the predictions. For algorithm selection, however, the more accurate runtime predictions translated only to a small improvement. This was somewhat expected since the selections based on the basic features already achieved very good performance.

Via the extensive empirical evaluation presented as part of this work, we managed to answer some of the key basic questions about the empirical hardness of BNSL. This first study opens several avenues for future research. First, we believe the proposed collection of features is not complete—presumably, there are even

¹⁶ <http://bnportfolio.cs.helsinki.fi/>,
http://github.com/coseal/aslib_data/tree/master/BNSL-2016

more informative, albeit possibly slower-to-compute, features yet to be discovered. For example, while not considered here, one straightforward possibility would be to use summary statistics for the BNSL features that are less susceptible to outliers, for example, medians. The question of how to efficiently trade informativeness for computational efficiency is relevant also more generally for the algorithm selection methodology; probing features [34], as applied in this work to the context of BNSL, provide just one, rather generic technique. Second, the empirical hardness model and its evaluated performance obviously depend on the distribution of the training and test instances. While this dependency is unavoidable, it is an intriguing question to what extent the dependency can be weakened by considering appropriate distributions and sufficiently large samples of instances.

Finally, we note that while in this work we focused on the runtime behavior of complete BNSL solvers, that is, exact algorithms that provide provably-optimal solutions to given BNSL instances, the techniques studied and developed in this paper could also be extended to cover in-exact local-search style, greedy, and approximate algorithmic approaches to BNSL. While such approaches typically exhibit better scalability than the exact approaches studied here, the fact that in-exact approaches cannot give guarantees of optimality on the produced solutions brings new challenges in terms of portfolio construction and prediction, specifically in understanding the interplay between solution quality and runtimes. Another potentially interesting direction for further study—although a somewhat secondary aspect compared to runtime behavior—would be to understand and predict the memory usage of exact approaches. Furthermore, it would be interesting to expand the study in the future by including additional datasets, for example, from OpenML [68].

Acknowledgements The authors thank James Cussens for discussions on GOBNILP and the anonymous reviewers for valuable suggestions that helped improve the manuscript.

References

1. Achterberg, T.: SCIP: Solving constraint integer programs. *Mathematical Programming Computation* **1**(1), 1–41 (2009)
2. Bache, K., Lichman, M.: UCI machine learning repository (2013). URL <http://archive.ics.uci.edu/ml>
3. Balint, A., Belov, A., Järvisalo, M., Sinz, C.: Overview and analysis of the SAT Challenge 2012 solver competition. *Artificial Intelligence* **223**, 120–155 (2015)
4. Bartlett, M., Cussens, J.: Integer linear programming for the Bayesian network structure learning problem. *Artificial Intelligence* **In press** (2015)
5. van Beek, P., Hoffmann, H.: Machine learning of Bayesian networks using constraint programming. In: Proceedings of the 21st International Conference on Principles and Practice of Constraint Programming (CP 2015), *Lecture Notes in Computer Science*, vol. 9255, pp. 429–445. Springer (2015)
6. Berg, J., Järvisalo, M., Malone, B.: Learning optimal bounded treewidth Bayesian networks via maximum satisfiability. In: Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS 2014), *JMLR Workshop and Conference Proceedings*, vol. 33, pp. 86–95. JMLR (2014)
7. Bielza, C., Larrañaga, P.: Discrete bayesian network classifiers: A survey. *ACM Comput. Surv.* **47**(1), 5:1–5:43 (2014)
8. Bischl, B., Kerschke, P., Kotthoff, L., Lindauer, M.T., Malitsky, Y., Fréchet, A., Hoos, H.H., Hutter, F., Leyton-Brown, K., Tierney, K., Vanschoren, J.: ASlib: A benchmark library for algorithm selection. *Artificial Intelligence* **237**, 41–58 (2016). DOI 10.1016/j.artint.2016.04.003. URL <http://dx.doi.org/10.1016/j.artint.2016.04.003>

9. Breiman, L.: Random forests. *Machine Learning* **45**, 5–32 (2001)
10. Buntine, W.: Theory refinement on Bayesian networks. In: *Proceedings of the 7th Conference on Uncertainty in Artificial Intelligence (UAI 1997)*, pp. 52–60. Morgan Kaufmann Publishers Inc. (1991)
11. de Campos, C., Ji, Q.: Efficient learning of Bayesian networks using constraints. *Journal of Machine Learning Research* **12**, 663–689 (2011)
12. Carbonell, J., Etzioni, O., Gil, Y., Joseph, R., Knoblock, C., Minton, S., Veloso, M.: Prodigy: an integrated architecture for planning and learning. *SIGART Bulletin* **2**, 51–55 (1991)
13. Cheng, J., Greiner, R., Kelly, J., Bell, D.A., Liu, W.: Learning Bayesian networks from data: An information-theory based approach. *Artificial Intelligence* **137**(1-2), 43–90 (2002)
14. Chickering, D.: Learning Bayesian networks is NP-complete. In: *Learning from Data: Artificial Intelligence and Statistics V*, pp. 121–130. Springer-Verlag (1996)
15. Cooper, G., Herskovits, E.: A Bayesian method for the induction of probabilistic networks from data. *Machine Learning* **9**, 309–347 (1992)
16. Cussens, J.: Bayesian network learning with cutting planes. In: *Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence (UAI 2011)*, pp. 153–160. AUA Press (2011)
17. Cussens, J.: Advances in Bayesian network learning using integer programming. In: *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI 2013)*, pp. 182–191. AUA Press (2013)
18. Fan, X., Malone, B., Yuan, C.: Finding optimal Bayesian network structures with constraints learned from data. In: *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence (UAI 2014)*, pp. 200–209. AUA Press (2014)
19. Fan, X., Yuan, C.: An improved lower bound for Bayesian network structure learning. In: *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI 2015)*, pp. 3526–3532. AAAI Press (2015)
20. Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., Hutter, F.: Efficient and robust automated machine learning. In: *Advances in Neural Information Processing Systems* **28** (2015)
21. Fink, E.: How to solve it automatically: Selection among problem-solving methods. In: *Proceedings of the 4th International Conference on Artificial Intelligence Planning Systems (AIPS 1998)*, pp. 126–136. AAAI Press (1998)
22. Fr chet te, A., Kotthoff, L., Michalak, T.P., Rahwan, T., Hoos, H.H., Leyton-Brown, K.: Using the shapley value to analyze algorithm portfolios. In: D. Schuurmans, M.P. Wellman (eds.) *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, pp. 3397–3403. AAAI Press (2016)
23. Friedman, N., Koller, D.: Being Bayesian about network structure. a Bayesian approach to structure discovery in Bayesian networks. *Machine Learning* **50**, 95–125 (2003)
24. Gebruers, C., Hnich, B., Bridge, D.G., Freuder, E.C.: Using CBR to select solution strategies in constraint programming. In: *6th International Conference on Case-Based Reasoning (ICCBR 2005)*, *Lecture Notes in Computer Science*, vol. 3620, pp. 222–236. Springer (2005)
25. Giraud-Carrier, C., Vilalta, R., Brazdil, P.: Introduction to the special issue on meta-learning. *Machine Learning* **54**(3), 187–193 (2004)
26. Gomes, C.P., Selman, B.: Algorithm portfolios. *Artificial Intelligence* **126**(1-2), 43–62 (2001)
27. Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I.H.: The WEKA data mining software: an update. *SIGKDD Explorations* **11**(1), 10–18 (2009)
28. Heckerman, D., Geiger, D., Chickering, D.: Learning Bayesian networks: The combination of knowledge and statistical data. *Machine Learning* **20**, 197–243 (1995)
29. Hoos, H., Kaminski, R., Lindauer, M., Schaub, T.: aspeed: Solver scheduling via answer set programming. *Theory and Practice of Logic Programming* **15**(1), 117–142 (2015)
30. Hoos, H., Lindauer, M.T., Schaub, T.: claspfolio 2: Advances in algorithm selection for answer set programming. *Theory and Practice of Logic Programming* **14**(4-5), 569–585 (2014)
31. Horvitz, E., Ruan, Y., Gomes, C.P., Kautz, H.A., Selman, B., Chickering, D.M.: A Bayesian approach to tackling hard computational problems. In: *Proceedings of the 17th Conference on Uncertainty in Artificial Intelligence (UAI 2001)*, pp. 235–244. Morgan Kaufmann (2001)

32. Hurley, B., Kotthoff, L., Malitsky, Y., O'Sullivan, B.: Proteus: A hierarchical portfolio of solvers and transformations. In: Proceedings of the 11th International Conference on Integration of AI and OR Techniques in Constraint Programming (CPAIOR 2014), *Lecture Notes in Computer Science*, vol. 8451, pp. 301–317. Springer (2014)
33. Hutter, F., Hoos, H.H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: Selected Papers of the 5th International Conference on Learning and Intelligent Optimization (LION 5), *Lecture Notes in Computer Science*, vol. 6683, pp. 507–523. Springer (2011)
34. Hutter, F., Xu, L., Hoos, H.H., Leyton-Brown, K.: Algorithm runtime prediction: Methods & evaluation. *Artificial Intelligence* **206**, 79–111 (2014)
35. Jaakkola, T.S., Sontag, D., Globerson, A., Meila, M.: Learning Bayesian network structure using LP relaxations. In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS 2010), *JMLR Proceedings*, vol. 9, pp. 358–365. JMLR.org (2010)
36. Järvisalo, M., Le Berre, D., Roussel, O., Simon, L.: The international SAT solver competitions. *AI Magazine* **33**(1), 89–92 (2012)
37. Koivisto, M., Sood, K.: Exact Bayesian structure discovery in Bayesian networks. *Journal of Machine Learning Research* pp. 549–573 (2004)
38. Kontkanen, P., Myllymäki, P.: MDL histogram density estimation. In: In Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics (AISTATS 2007), *JMLR Proceedings*, vol. 2, pp. 219–226. JMLR.org (2007)
39. Kotthoff, L.: Algorithm selection for combinatorial search problems: A survey. *AI Magazine* **35**(3), 48–60 (2014)
40. Kotthoff, L., Gent, I.P., Miguel, I.: An evaluation of machine learning in algorithm selection for search problems. *AI Communications* **25**(3), 257–270 (2012)
41. Kotthoff, L., Kerschke, P., Hoos, H., Trautmann, H.: Improving the state of the art in inexact TSP solving using per-instance algorithm selection. In: Revised Selected Papers of the 9th International Conference on Learning and Intelligent Optimization (LION 9), *Lecture Notes in Computer Science*, vol. 8994, pp. 202–217. Springer (2015)
42. Lee, J.W., Giraud-Carrier, C.G.: Predicting algorithm accuracy with a small set of effective meta-features. In: Proceedings of the 7th International Conference on Machine Learning and Applications (IEEE ICMLA 2008), pp. 808–812. IEEE Computer Society (2008)
43. Leite, R., Brazdil, P., Vanschoren, J.: Selecting classification algorithms with active testing. In: Proceedings of the 8th International Conference on Machine Learning and Data Mining in Pattern Recognition (MLDM 2012), *Lecture Notes in Computer Science*, vol. 7376, pp. 117–131. Springer (2012)
44. Leyton-Brown, K., Hoos, H.H., Hutter, F., Xu, L.: Understanding the empirical hardness of NP-complete problems. *Commun. ACM* **57**(5), 98–107 (2014)
45. Leyton-Brown, K., Nudelman, E., Shoham, Y.: Learning the empirical hardness of optimization problems: The case of combinatorial auctions. In: 8th International Conference on Principles and Practice of Constraint Programming (CP 2002), *Lecture Notes in Computer Science*, vol. 2470, pp. 556–572. Springer (2002)
46. Leyton-Brown, K., Nudelman, E., Shoham, Y.: Empirical hardness models: Methodology and a case study on combinatorial auctions. *Journal of the ACM* **56**(4) (2009)
47. Lindauer, M.T., Hoos, H.H., Hutter, F., Schaub, T.: AutoFolio: An automatically configured algorithm selector. *Journal of Artificial Intelligence Research* **53**, 745–778 (2015)
48. Lobjois, L., Lemaître, M.: Branch and bound algorithm selection by performance prediction. In: Proceedings of the 15th National Conference on Artificial Intelligence (AAAI 1998), pp. 353–358. AAAI Press (1998)
49. Madigan, D., York, J.: Bayesian graphical models for discrete data. *International Statistical Review* **63**, 215–232 (1995)
50. Malone, B., Järvisalo, M., Myllymäki, P.: Impact of learning strategies on the quality of Bayesian networks: An empirical evaluation. In: Proceedings of the 31st Conference on Uncertainty in Artificial Intelligence (UAI 2015), pp. 362–371. AUAI Press (2015)
51. Malone, B., Kangas, K., Järvisalo, M., Koivisto, M., Myllymäki, P.: Predicting the hardness of learning Bayesian networks. In: Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI 2014), pp. 2460–2466. AAAI Press (2014)
52. Malone, B.M., Yuan, C.: Evaluating anytime algorithms for learning optimal Bayesian networks. In: Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI 2013). AUAI Press (2013)

53. Ott, S., Imoto, S., Miyano, S.: Finding optimal models for small gene networks. In: Proceedings of the Pacific Symposium on Biocomputing 2004, pp. 557–567. World Scientific (2004)
54. Parviainen, P., Koivisto, M.: Finding optimal Bayesian networks using precedence constraints. *Journal of Machine Learning Research* **14**, 1387–1415 (2013)
55. Pearl, J.: Probabilistic reasoning in intelligent systems: networks of plausible inference. Morgan Kaufmann (1988)
56. Perrier, E., Imoto, S., Miyano, S.: Finding optimal Bayesian network given a superstructure. *Journal of Machine Learning Research* **9**, 2251–2286 (2008)
57. Pulina, L., Tacchella, A.: Treewidth: A useful marker of empirical hardness in quantified Boolean logic encodings. In: Proceedings of the 15th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning (LPAR 2008), *Lecture Notes in Computer Science*, vol. 5330, pp. 528–542. Springer (2008)
58. Rice, J.: The algorithm selection problem. *Advances in Computers* **15**, 65–118 (1976)
59. Rijn, J.N., Abdulrahman, S.M., Brazdil, P., Vanschoren, J.: Fast algorithm selection using learning curves. In: Proceedings of the 14th International Symposium on Advances in Intelligent Data Analysis (IDA 2015), *Lecture Notes in Computer Science*, vol. 9385, pp. 298–309. Springer (2015)
60. Saikko, P., Malone, B., Järvisalo, M.: MaxSAT-based cutting planes for learning graphical models. In: Proceedings of the 12th International Conference on Integration of Artificial Intelligence and Operations Research Techniques in Constraint Programming (CPAIOR 2015), *Lecture Notes in Computer Science*, vol. 9075, pp. 345–354. Springer (2015)
61. Shapley, L.S.: A value for n -person games. *Contributions to the theory of games* **2**, 307–317 (1953)
62. Silander, T., Myllymäki, P.: A simple approach for finding the globally optimal Bayesian network structure. In: Proceedings of the 22nd Conference in Uncertainty in Artificial Intelligence (UAI 2006), pp. 445–452. AUAI Press (2006)
63. Singh, A., Moore, A.: Finding optimal Bayesian networks by dynamic programming. Tech. rep., Carnegie Mellon University (2005)
64. Sokal, R.R., Michener, C.D.: A statistical method for evaluating systematic relationships. *The University of Kansas Science Bulletin* **38**(2), 1409–1438 (1958)
65. Spirtes, P., Glymour, C., Schemes, R.: Causation, Prediction, and Search. Springer, New York (1993)
66. Tamada, Y., Imoto, S., Miyano, S.: Parallel algorithm for learning optimal Bayesian network structure. *Journal of Machine Learning Research* **12**, 2437–2459 (2011)
67. Teyssier, M., Koller, D.: Ordering-based search: A simple and effective algorithm for learning Bayesian networks. In: Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence (UAI 2005), pp. 584–590. AUAI Press (2005)
68. Vanschoren, J., van Rijn, J.N., Bischl, B., Torgo, L.: OpenML: networked science in machine learning. *SIGKDD Explorations* **15**(2), 49–60 (2013)
69. Wunderling, R.: Paralleler und objektorientierter Simplex-Algorithmus. Ph.D. thesis, Technische Universität Berlin (1996)
70. Xu, L., Hutter, F., Hoos, H., Leyton-Brown, K.: SATzilla: Portfolio-based algorithm selection for SAT. *Journal of Artificial Intelligence Research* **32**, 565–606 (2008)
71. Yuan, C., Malone, B.: An improved admissible heuristic for finding optimal Bayesian networks. In: Proceedings of the 27th Conference in Uncertainty in Artificial Intelligence (UAI 2012), pp. 924–933. AUAI Press (2012)
72. Yuan, C., Malone, B.: Learning optimal Bayesian networks: A shortest path perspective. *Journal of Artificial Intelligence Research* **48**, 23–65 (2013)

Appendix A Details on the Data Sets

The numbers of variables and records in each of the data sets used in the experiments are shown in Tables 9 and 10 for REAL and SAMPLED, respectively.

Table 9 Sizes of the datasets in REAL.

Dataset	#Variables	#Records
letter	17	20,000
voting	17	435
zoo	17	101
lymph	19	148
eucalyptus	20	736
hepatitis	20	155
credit-g	21	1,000
hypothyroid	22	3,772
mushroom	22	8,124
spect	23	267
autos	26	205
colic	28	368
pyrim	28	74
flag	29	194
trains	30	10
anneal	32	898
backache	32	180
marketing	33	364
student-mat	33	395
student-por	33	649
turkiye	33	5,820
dermatology	35	366
soybean	36	307
kr-vs-kp	37	3,196
stemmatology	37	1,208
abscisic	41	5,456
diabetes	41	60,000
connect-4_6000	43	6,000
connect-4_60000	43	60,000
covtype_60000	43	60,000
sponge	45	76
wiki4he	53	913
lung-cancer	57	32
promoters	58	106
triazines	59	186
splice	61	3,190
audiology_63	63	226
optdigits	63	5,620
plants_63	63	34,781

Table 10 Sizes of the datasets in SAMPLED.

Dataset	#Variables	#Records
kredit	18	1,000
insurance	27	100; 1,000; 10,000
water	32	100; 1,000; 10,000
mildew	35	100; 1,000; 10,000
alarm	37	100; 1,000; 10,000
hailfinder	56	100; 1,000; 10,000
carpo	60	100; 1,000; 10,000