AUTO-ENSEMBLE STRUCTURE LEARNING OF LARGE GAUSSIAN BAYESIAN NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Learning the structure of Bayesian networks (BNs) from data is challenging, especially for datasets involving a large number of variables. The recently proposed divide-and-conquer (D&D) strategies present a promising approach for learning large BNs. However, they still face a main issue of unstable learning accuracy across subproblems. In this work, we introduce the idea of employing structure learning ensemble (SLE), which combines multiple BN structure learning algorithms together, to consistently achieve high learning accuracy across various problems. We further propose an automatic approach called Auto-SLE for constructing near-optimal SLEs, addressing the challenge of manually designing effective SLEs. The automatically constructed SLE is then integrated into a D&D framework. Extensive experiments firmly show the superiority of our method over existing methods in learning large BNs, achieving accuracy improvement usually by $30\% \sim 225\%$ on datasets involving 10,000 variables. Furthermore, our method generalizes very well to datasets with many more variables and different network characteristics than those present in the training data for constructing the SLE. These results indicate the significant potential of employing automatic construction of SLEs for BN structure learning.

026 027 028

029

024

025

004

010 011

012

013

014

015

016

017

018

019

020

021

1 INTRODUCTION

Learning the structure of Bayesian networks (BNs) (Pearl, 1985) from data has attracted much re-031 search interest, due to its wide applications in machine learning, statistical modeling, and causal 032 inference (Pearl, 1988; Kitson et al., 2023). Various methods have been proposed to tackle this 033 problem, including constraint-based methods (Colombo et al., 2014), score-based methods (Ramsey 034 et al., 2017), and hybrid methods (Tsamardinos et al., 2006). However, most previous studies primarily dealt with a relatively small number of variables. For example, the *bnlearn* repository (Scutari, 2010), which is widely used in the literature, contains mostly networks with only a few dozen nodes 037 (variables). In contrast, in real-world applications such as alarm events analysis (Cai et al., 2022), 038 MRI image interpretation (Ramsey et al., 2017), and human genome analysis (Schaffter et al., 2011), it is common to generate and collect data from thousands of variables and beyond. Unfortunately, as the number of variables increases, many of the existing methods would slow down dramatically 040 and become much less accurate (Zhu et al., 2021). 041

Recently, Gu and Zhou (2020) introduced a divide-and-conquer (D&D) framework named partitionestimation-fusion (PEF) for the structure learning of large BNs. Specifically, PEF consists of three steps: partitioning nodes into clusters (partition), learning a subgraph on each cluster (estimation), and fusing all subgraphs into a single BN (fusion). It is noteworthy that PEF is flexible as any existing structure learning algorithm can be used in the estimation step. Additionally, due to the smaller number of nodes in each cluster compared to the total number of nodes, the overall structure learning is accelerated. Moreover, the structure learning processes for different subproblems can be parallelized in a straightforward way, leading to further improvement in computational efficiency.

However, despite the evident advantages provided by PEF, it still faces a main issue of unstable
structure learning accuracy across subproblems. The root cause for this is that the partition step
of PEF tends to yield subproblems with significantly different characteristics, e.g., varying node
numbers. When a single structure learning algorithm is used to solve all subproblems, as in existing PEF-based methods (Gu & Zhou, 2020), achieving stable learning accuracy across different



Figure 1: An overview of P/SLE. The SLE is automatically constructed by Auto-SLE and integrated into the estimation step of PEF.

054

056

058

059

060

061

062

063

064

065 066

067

068

069

071

subproblems becomes challenging. In fact, even for the same algorithm, different hyperparameter values can lead to significant variations in its behavior, thereby affecting its suitability for solving specific subproblems. For instance, when employing the well-known fast greedy equivalence search (fGES) (Ramsey et al., 2017) algorithm with Bayesian information criterion (BIC) as the score function, its penalty coefficient should ideally increase with the sparsity of the underlying BN. However, given a BN structure learning problem, determining the optimal penalty coefficient in advance is difficult as the underlying BN is unknown.

082 Inspired by the success of ensemble methods like AutoAttack (Croce & Hein, 2020) in the field of 083 adversarial robustness, which utilizes an attack ensemble to achieve more reliable robustness evaluation compared to individual attacks, we introduce the idea of employing structure learning ensemble 084 (SLE) to achieve stable learning accuracy in BN structure learning. Specifically, a SLE comprises 085 several structure learning algorithms, dubbed member algorithms. When applied to a BN structure learning problem, a SLE runs its member algorithms individually and chooses the best of their out-087 puts. Similar to how AutoAttack integrates diverse attacks, a high-performing SLE should consist 088 of complementary member algorithms that excel at solving different types of problems. However, 089 in contrast to AutoAttack where the attacks are manually designed and selected, we propose to auto-090 matically construct SLEs. This can significantly reduce the reliance on human expertise and effort, 091 as manually constructing SLEs typically requires domain experts to explore the vast design space of 092 SLEs, which can be both laborious and intricate.

Specifically, we first formulate the problem of automatic SLE construction, and then present Auto-094 SLE, a simple yet effective approach to address this problem. The approach is implemented with a 095 design space containing several candidate algorithms; each algorithm has a set of hyperparameters. 096 To construct a SLE, Auto-SLE starts from an empty ensemble, and repeatedly adds the algorithm and its associated hyperparameter values to the ensemble such that the ensemble improvement is 098 maximized. In addition to its conceptual simplicity, Auto-SLE is appealing also because it yields SLEs provably within a constant factor of optimal for the training problem set. To achieve our goal of constructing a single SLE that can generalize well across network sizes and characteristics, we 100 utilize a training set comprising diverse networks of varying sizes. Finally, the constructed SLE is 101 integrated into the PEF framework, resulting in P/SLE (see Figure 1 for an overview). 102

From extensive comparisons with the state-of-the-art baselines, it is found that P/SLE consistently achieves significantly higher accuracy in learning large BNs, and its superiority becomes even more pronounced as the network size further increases. On datasets involving 10,000 variables, P/SLE typically achieves accuracy improvement by 30%~225%. Further experiments show that P/SLE, without any additional tuning or adaptation, generalizes well to datasets with much larger number (e.g., 30,000) of variables and possessing different network characteristics than the training data.

108 It is worth mentioning that Auto-SLE itself is a general SLE construction approach that can be applied to various types of BNs (Gaussian and non-Gaussian) and even other types of causal models. 110 In this work, we focus on Gaussian BNs, which are arguably the most widely studied type of BNs 111 and have a rich set of baselines, allowing us to thoroughly assess the potential of Auto-SLE. Besides, 112 it should be noted that PEF-based methods, such as P/SLE, are most suitable for learning BNs with a block structure to some extent, meaning the connections between subgraphs are relatively weak. 113 It is quite common for a large network to exhibit such a block structure, due to the underlying 114 heterogeneity among the nodes (Holland et al., 1983; Airoldi et al., 2008; Abbe et al., 2015). On 115 the other hand, our results also indicate that, even for large BNs without a block structure, P/SLE 116 can still achieve competitive learning accuracy. The promising performance of P/SLE not only 117 demonstrates the potential of applying SLEs to BN structure learning, a largely unexplored area, but 118 also highlights the effectiveness of Auto-SLE as a simple and easy-to-use approach for constructing 119 SLEs. This is expected to further facilitate the advancement of SLEs in this field. 120

120 121 122

123

124

132 133

137 138

144

2 PRELIMINARIES AND RELATED WORK

2.1 THE BN STRUCTURE LEARNING PROBLEM

The structure of a BN for d random variables $X_1, ..., X_d$ is represented by a directed acyclic graph (DAG), denoted as $\mathcal{G} = (V, E)$. Here, $V = \{1, ..., d\}$ is the set of nodes corresponding to the random variables and $E = \{(j, i) \in V \times V : j \rightarrow i\}$ is the directed edge set. Define $V_{\text{pa}(i)} = \{j \in V : (j, i) \in E\}$ as the parent node set of node *i*, and $X_{\text{pa}(i)}$ as the set of corresponding random variables. The joint probability density function *f* of $X_1, ..., X_d$ is factorized according to the structure of \mathcal{G} :

$$f(X_1, X_2, ..., X_d) = \prod_{i=1}^d f(X_i | X_{\text{pa}(i)}),$$
(1)

where $f(X_i|X_{pa(i)})$ is the conditional probability density of X_i given $X_{pa(i)}$. In this work, we focus on Gaussian BNs for continuous data. Specifically, the conditional distributions are specified by the following linear structural equation model:

$$X_i = \phi\left(X_{\text{pa}(i)}\right) + \varepsilon_i, \quad i = 1, ..., d,$$
(2)

where $\phi(\cdot)$ denotes a linear function and $\varepsilon_i \sim \mathcal{N}(0, \sigma_j^2)$. Suppose we have obtained m iid observations of X_1, \ldots, X_d , denoted as $D \in \mathbb{R}^{m \times d}$. Given D, the goal is to learn a DAG structure $\mathcal{G} = (V, E)$ that accurately reflects the conditional dependencies among X_1, \ldots, X_d . In practice, accuracy metrics such as F1 score and structural Hamming distance (SHD) are typically used to assess the quality of the learned BN.

145 2.2 RELATED WORK

146 The BN structure learning problem has been proven to be NP-hard (Chickering et al., 2004), leading 147 to main research efforts on developing approximation methods to solve it. These methods can be 148 broadly classified into constraint-based, score-based, and hybrid methods. Constraint-based meth-149 ods, such as PC (Spirtes et al., 2000), MMPC (Tsamardinos et al., 2003a), and PC-Stable (Colombo 150 et al., 2014), use conditional independence tests on observations to identify relationships among 151 variables. In comparison, score-based methods explore the space of DAGs or Markov equivalence 152 classes (MECs) using search heuristics such as tabu search (TS) (Bouckaert, 1995), genetic algorithm (GA) (Larranaga et al., 1996), and greedy search (Chickering, 2002). These methods also 153 employ score functions, such as BDeu (Akaike, 1974), BIC (Schwarz, 1978), and K2 (Cooper & 154 Herskovits, 1992), to guide the search. It is worth mentioning a recent research line of score-based 155 methods including NOTEARS (Zheng et al., 2018) and LEAST (Zhu et al., 2021) that reformulate 156 the structure learning as a continuous optimization problem. Finally, hybrid methods integrate both 157 constraint-based and score-based techniques. For example, MMHC (Tsamardinos et al., 2006) uses 158 MMPC to build the graph skeleton and utilizes TS to determine the final BN. 159

However, as the number of variables increases, many of the existing methods would slow down
 dramatically and become much less accurate (Zhu et al., 2021). Actually, based on our preliminary
 testing of 15 existing methods, fGES (Ramsey et al., 2017), which is a variant of the greedy search

algorithm (Chickering, 2002), and PC-Stable (Colombo et al., 2014), are the only methods capable
of maintaining relatively stable learning accuracy (e.g., achieving F1 score of around 0.5) within
reasonable time when the variable number reaches 1000. Moreover, when the number of variables
reaches 10,000, all methods, even after running for a quite long period of time (e.g., 24 hours), are
unable to output solutions.

PEF Framework To address the challenge of learning large BNs, Gu and Zhou (2020) introduced a partition-estimation-fusion (PEF) framework. It comprises the following three steps.

- Partition: The d nodes are partitioned into clusters with a hierarchical clustering algorithm.
- Estimation: An existing structure learning algorithm is applied to estimate a subgraph on each cluster of nodes.
- Fusion: Merge estimated subgraphs into one DAG containing all the *d* nodes.

The details of PEF can be found in Appendix A. While PEF has dramatically enhanced the capabil ities of handling large BNs, it still faces a main issue of unstable structure learning accuracy across
 subproblems. In the following we will describe the use of automatically constructed SLE in the
 estimation step to address this issue.

180 181

182

183

192 193

197 198

167

168

169 170

171

172

173 174

175

3 AUTOMATIC CONSTRUCTION OF SLES

3.1 PROBLEM FORMULATION

We first formulate the SLE construction problem. Formally, a SLE with k member algorithms is denoted as $\mathbb{A} = \{\theta_1, \dots, \theta_k\}$, where θ_i represents the *i*-th member algorithm of \mathbb{A} . Let $T = \{D_1, D_2, \dots\}$ denote a training problem set, where each D_i represents a BN structure learning problem with known ground truth. When using \mathbb{A} to solve a problem $D \in T$, one straightforward strategy is to run all member algorithms of \mathbb{A} individually in parallel, and the best solution among all the found solutions in terms of a quality measure Q (e.g., F1 score) is returned. Let $Q(\mathbb{A}, D)$ and $Q(\theta_i, D)$ denote the performance of \mathbb{A} and θ_i on D in terms of Q, respectively. Without loss of generality, we assume a larger value is better for Q. Then we have:

$$Q(\mathbb{A}, D) = \max_{\theta \in \mathbb{A}} Q(\theta_i, D).$$
(3)

Then the performance of \mathbb{A} on T in terms of Q, denoted as $Q(\mathbb{A}, T)$, is the average value of \mathbb{A} 's performance on the training problems in $T(Q(\mathbb{A}, T) = 0 \text{ when } \mathbb{A} \text{ is empty})$:

$$Q(\mathbb{A},T) = \frac{1}{|T|} \sum_{D \in T} Q(\mathbb{A},D).$$
(4)

199 For the automatic construction of \mathbb{A} , the member algorithms of \mathbb{A} are not manually determined but 200 automatically selected from an algorithm configuration space Θ . Specifically, suppose we have a 201 candidate algorithm pool $\{A_1, ..., A_n\}$, which can be constructed by collecting existing algorithms. 202 Each candidate algorithm has some hyperparameters. Let Θ_i denote the hyperparameter configu-203 ration space of A_i , where a configuration $\theta \in \Theta_i$ refers to a setting of A_i 's hyperparameters, such 204 that its behaviors is completely specified. Then, the algorithm configuration space Θ is defined as $\Theta = \Theta_1 \cup \Theta_2 \cdots \cup \Theta_n$, and each $\theta \in \Theta$ represents a specific candidate algorithm along with the 205 specific values of its hyperparameters. As presented in Definition 1, the SLE construction problem 206 is to select k member algorithms from Θ to form a SLE \mathbb{A}^* , such that its performance on the training 207 set T in terms of Q is maximized. 208

Definition 1. Given T, Q, Θ , and k, the SLE construction problem is to find $\mathbb{A}^* = \{\theta_1^* \dots \theta_k^*\}$ that maximizes $Q(\mathbb{A}^*, T)$, s.t. $\theta_i^* \in \Theta$ for $i = 1 \dots k$.

- 211
- 212 3.2 AUTO-SLE: A GREEDY APPROACH 213
- 214 We now introduce Auto-SLE, a simple yet effective approach for automatically constructing SLEs. 215 As shown in Algorithm 1, Auto-SLE starts with an empty ensemble \mathbb{A} (line 1) and finds the candidate algorithm and its hyperparameter values denoted as θ° that, if included in \mathbb{A} , maximizes ensemble

2	1	7
2	1	ε
2	1	ç

Algorithm 1: Auto-SLE

Input: quality measure Q, training set T, algorithm configuration space Θ , ensemble size k 8 **Output:** A $1 \mathbb{A} \leftarrow \emptyset, i \leftarrow 1;$ ² while $i \leq k$ do 221 $\theta^{\circ} \leftarrow \arg \max_{\theta \in \Theta} Q (\mathbb{A} \cup \{\theta\}, T) - Q (\mathbb{A}, T);$ 3 222 if $Q(\mathbb{A} \cup \{\theta^{\circ}\}, T) = Q(\mathbb{A}, T)$ then return \mathbb{A} ; $\mathbb{A} \leftarrow \mathbb{A} \cup \{\theta^{\circ}\}, i \leftarrow i+1;$ 5 224 6 end 225 7 return A

226 227 228

229

230

231

239

240

254 255

256

257 258

259

improvement in terms of Q (line 3). Ties are broken arbitrarily here. After this, θ° is subject to the following procedure: if adding it to \mathbb{A} does not improve performance, which means the construction process has converged, Auto-SLE will terminate and return A (line 4); otherwise θ° is added to A (line 5). The above process will be repeated until k member algorithms have been found (line 2).

232 Let $\Delta(\theta|\mathbb{A}) = Q(\mathbb{A} \cup \{\theta\}, T) - Q(\mathbb{A}, T)$ denote the performance improvement brought by adding 233 θ to \mathbb{A} . Noticing that each iteration of Auto-SLE needs to find θ° that maximizes $\Delta(\theta|\mathbb{A})$, when the 234 algorithm configuration space Θ is large or even infinite (e.g., candidate algorithms have continuous 235 hyperparameters), using enumeration to find θ° is impractical. In practice, we employ hyperparame-236 ter optimization procedures, such as Bayesian optimization (Lindauer et al., 2022), to approximately maximize $\Delta(\theta|\mathbb{A})$, and the runs of these procedures account for the vast majority of the total com-237 putational costs of Auto-SLE. 238

3.3 THEORETICAL JUSTIFICATIONS

241 We now theoretically analyze the performance of Auto-SLE on a give training problem set. For 242 notational simplicity, henceforth we omit the T in $Q(\cdot, T)$ and directly use $Q(\cdot)$. Our analysis is 243 based on the following key fact that $Q(\cdot)$ is monotone and submodular. 244

Fact 1. $Q(\cdot)$ is a monotone submodular function, i.e., for any two SLEs $\mathbb{A}, \mathbb{A}' \subset \Theta$ and any $\theta \in \Theta$, 245 it holds that $Q(\mathbb{A}) \leq Q(\mathbb{A} \cup \mathbb{A}')$ and $Q(\mathbb{A} \cup \mathbb{A}' \cup \{\theta\}) - Q(\mathbb{A} \cup \mathbb{A}') \leq Q(\mathbb{A} \cup \{\theta\}) - Q(\mathbb{A})$. 246

247 The proof is straightforward and can be found in Appendix B. Intuitively, Q exhibits a diminishing 248 returns property that the marginal gain of adding θ diminishes as the ensemble size increases. Based 249 on Fact 1, Theorem 1 holds.

250 **Theorem 1.** Using a hyperparameter optimization procedure that, in each iteration of Auto-SLE, 251 returns θ within ϵ -absolute error of the maximum of $\Delta(\theta|\mathbb{A})$, i.e., $\Delta(\theta|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) - \epsilon$, then the 252 quality $Q(\mathbb{A})$ of the SLE constructed by Auto-SLE is bounded by 253

$$Q(\mathbb{A}) \ge (1 - 1/e) \cdot Q(\mathbb{A}^*) - k\epsilon, \tag{5}$$

where \mathbb{A}^* is the optimal SLE to the SLE construction problem in Definition 1. Alternatively, if θ is within ϵ -relative error of $\Delta(\theta^{\circ}|\mathbb{A})$, i.e., $\Delta(\theta|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) \cdot (1-\epsilon)$, then the quality $Q(\mathbb{A})$ of the SLE constructed by Auto-SLE is bounded by

$$Q(\mathbb{A}) \ge (1 - 1/e^{1 - \epsilon}) \cdot Q(\mathbb{A}^*).$$
(6)

260 The proof (see Appendix B) is a slight extension of the classical derivations in maximizing $\Delta(\theta|\mathbb{A})$ (Nemhauser et al., 1978). Based on Theorem 1, Auto-SLE achieves (1 - 1/e)-approximation for 261 the optimal quality when given a perfect hyperparameter optimization procedure with $\epsilon = 0$. Sub-262 optimal hyperparameter optimization procedures result in worse outcomes but small errors ϵ do not 263 escalate. This is important because with large algorithm configuration space, it cannot be expected 264 for blackbox optimization procedures to find θ° in realistic time. However, at least for some scenar-265 ios with a few parameters, widely-used Bayesian optimization techniques such as SMAC (Lindauer 266 et al., 2022) have empirically been shown to yield performance close to optimal within reasonable 267 time budget. 268

Note that Theorem 1 only provides a performance guarantee for the training problem set used for 269 constructing the SLE. By employing the same techniques introduced by Liu et al. (2020), we can

recover guarantees for an independent test set that is sampled or generated from the same distribution
 as the training set, given a sufficiently large training set. Nevertheless, our experimental results (see
 Section 4.2) demonstrate that the constructed SLE performs well beyond the training set.

273 274 275

3.4 APPLYING THE CONSTRUCTED SLE TO A NEW PROBLEM

Importantly, when presented with a testing problem to solve, it cannot be assumed that the ground truth is known. Thus, quality measures that do not require ground truth, such as the BIC score adopted in our experiments, are used to select the best output from the outputs of member algorithms. It is worth mentioning that although this work focuses on utilizing the constructed SLE to enhance PEF, the SLE itself can also serve as a complete and independent BN structure learning method.

281 282 283

284 285

286

287

289

290

4 EXPERIMENTS

Extensive experiments are conducted to answer two research questions (RQs).

- RQ1: Does the SLE constructed by Auto-SLE enhance PEF in learning large BNs?
- **RQ2**: Can the SLE generalize to larger problem sizes and different network characteristics than that present in the training set?
- 291 292 4.1 EXPERIMENTAL SETUP

293 Benchmarks and Evaluation Metrics We generate diverse and large-scale benchmarks follow-294 ing the approach in (Gu & Zhou, 2020). Specifically, the approach consists of four steps: (i) select 295 an existing network structure and replicates it until a predefined variable number is reached; (ii) 296 connect the replicas by adding 10% of edges between them randomly, while ensuring the final net-297 work remains a DAG; (iii) utilize the complete DAG and Eq. (2) to generate observations, where the 298 weights in the linear function and the standard deviations of Gaussian noises are sampled uniformly from $[-1, -0.5] \cup [0.5, 1]$ and [0, 1], respectively; (iv) re-scale the observations such that all data 299 columns have the same mean and standard deviation. We select 10 networks from the *bnlearn* repos-300 itory (Scutari, 2010) with node numbers ranging from 5 to 441. Based on each of them, we use the 301 above approach to generate testing problems with around 1000 and 10,000 variables. Following (Gu 302 & Zhou, 2020), the sample size m in each testing problem is set to 1000. 303

The widely-used F1 score and SHD are adopted as the metrics for assessing learning accuracy. In line with previous comparative study (Ramsey et al., 2017), we use two specific variants of F1 score, i.e., F1 Arrowhead (F1^{\rightarrow}) that considers direction and F1 Adjacent (F1⁻) that ignores direction, since some methods (PC-Stable and fGES) output MECs of DAGs which may contain edges without directions. Moreover, the wall-clock runtime of the methods is reported. For F1 metrics, a higher value is better; for SHD and runtime, a lower value is better.

310

Constructing the SLE with Auto-SLE A diverse training problem set is beneficial for construct-311 ing a SLE with good performance across problem sizes and network characteristics. Specifically, 312 the training set T comprises 100 problems, with variable numbers ranging from 5 to 1000, generated 313 using the above approach based on a network randomly selected from the 32 networks in *bnlearn*. 314 These problems are exclusively used for constructing the SLE and are independent of the testing 315 problems. For the algorithm configuration space Θ , we consider two algorithms PC-stable (with 316 two hyperparameters) and fGES (with three hyperparameters) as candidate algorithms, which out-317 perform others in our preliminary testing. We use their implementations from the causal discovery 318 tool box *TETRAD* (Ramsey et al., 2018). The sum of $F1^-$ and $F1^-$ is adopted as the quality mea-319 sure Q, and ensemble size k is set to 4 as running more iterations of Auto-SLE brings a negligible 320 improvement to the SLE's performance on the training set (see Appendix C.3). SMAC (version 321 3) (Lindauer et al., 2022), a Bayesian optimization tool, is used to maximize $\Delta(\theta|\mathbb{A})$ in each iteration of Auto-SLE, with a time budget 12 hours per run. Consequently, Auto-SLE consumes 322 approximately 48 hours in total to construct the SLE. Then, the SLE is integrated into PEF, resulting 323 in P/SLE, which is evaluated in subsequent experiments without further tuning or adaptation.

324 Table 1: Results on testing problems with 1000 variables, in terms of F1 Adjacent (F1⁻), F1 Arrow-325 head (F1^{\rightarrow}), SHD, and runtime (T). On each network, the mean ± std performance obtained by each 326 method on 10 problems is reported. The best performance in terms of accuracy metrics is marked 327 with an underline, and the performance that is not significantly different from the best performance (according to a Wilcoxon signed-rank test with significance level p = 0.05) is indicated in **bold**. Let 328 B be the best performance achieved among the baselines and A be the performance of P/SLE. The 329 improvement (Impro.) ratio is calculated as (A-B)/B for F1 metrics (a higher value is better), and is 330 calculated as (B-A)/B for SHD (a lower value is better). 331

Problem (V , E)		Alarm (1036,1417)	Asia (1000,1100)	Cancer (1000,880)	Child (1000,1375)	Earthquake (1000,880)	Hailfinder (1008,1307)	Healthcare (1001,1416)	Mildew (1015,1468)	Pigs (1323,1954)	Survey (1002,1103)
P/SLE	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T (s) \end{array}$	0.81±0.02 0.68±0.04 625.7±73.6 7.4±0.2	0.96±0.00 0.74±0.01 125.2±15.2 5.7±0.1	0.98±0.00 0.94±0.01 43.1±6.9 5.1±0.1	0.86±0.01 0.60±0.02 581.1±23.3 8.4±0.3	0.98±0.00 0.94±0.01 46.1±8.1 5.1±0.1	0.80±0.02 0.61±0.03 571.3±62.7 243.3±249.0	0.89±0.01 0.59±0.01 <u>485.0±38.4</u> 6.4±0.2	0.68±0.02 0.54±0.03 1135.4±71.8 8.7±0.6	0.70±0.02 0.61±0.02 1262.8±84.4 726.9±1176.6	0.92±0.01 0.79±0.02 309.9±29.8 6.0±0.1
P/SLE(D)	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.68±0.01 0.57±0.03 1234.3±77.9 26.0±1.8	0.78±0.01 0.64±0.01 710.5±24.7 20.4±3.0	0.74±0.01 0.73±0.01 611.1±31.2 17.1±2.7	0.75±0.01 0.56±0.01 1055.5±39.4 21.6±5.0	0.74±0.01 0.72±0.01 617.5±27.1 15.4±2.3	0.66±0.01 0.57±0.02 1158.4±65.2 366.0±287.4	0.77±0.01 0.56±0.01 1051.6±31.0 9.7±1.8	0.56±0.02 0.45±0.03 1826.2±100.4 14.2±2.2	0.53±0.01 0.47±0.02 2442.4±115.6 5697.5±1870.6	0.75±0.01 0.62±0.02 871.6±32.0 7.8±0.1
P/SLE(R)	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.72±0.02 0.42±0.02 1070.1±56.9 19.6±1.1	0.87±0.00 0.49±0.01 511.7±17.2 12.2±1.9	0.77±0.01 0.32±0.03 697.1±44.0 11.4±1.9	0.83±0.01 0.51±0.02 680.7±34.7 13.7±1.7	0.77±0.01 0.30±0.02 703.7±34.5 10.4±1.4	0.66±0.01 0.44±0.02 1152.3±38.7 33.3±8.3	0.87±0.01 0.44±0.01 750.8±31.6 12.9±3.1	0.57±0.02 0.31±0.01 1647.1±60.2 15.5±1.8	0.58±0.04 0.35±0.03 1934.3±60.6 5886.7±1671.0	0.84±0.01 0.46±0.01 675.9±28.1 8.3±2.0
P/PC-Stable	$\begin{array}{c} F1^- \\ F1^- \\ SHD \\ T (s) \end{array}$	0.76±0.02 0.48±0.02 883.5±42.2 5.5±0.8	0.91±0.00 0.56±0.01 342.7±9.5 6.1±0.1	0.85±0.01 0.47±0.03 414.8±23.8 6.0±0.1	0.85±0.01 0.52±0.02 585.3±29.1 6.2±0.9	0.85±0.01 0.43±0.02 416.8±16.4 5.7±0.4	0.71±0.01 0.48±0.02 922.3±41.1 19.1±6.3	0.89±0.01 0.46±0.01 648.0±20.3 2.8±0.2	0.62±0.02 0.35±0.02 1382.5±42.3 6.6±0.6	0.62±0.05 0.40±0.04 1627.4±84.8 4806.5±2311.3	0.89±0.01 0.58±0.03 458.0±26.2 5.3±0.8
P/fGES	$\begin{array}{c} F1^- \\ F1^\to \\ SHD \\ T \ (s) \end{array}$	0.68±0.01 0.57±0.03 1234.3±77.9 9.3±1.1	0.78±0.01 0.64±0.01 710.5±24.7 9.0±0.4	0.74±0.01 0.73±0.01 611.1±31.2 8.1±0.4	0.75±0.01 0.56±0.01 1055.5±39.4 9.0±1.3	0.74±0.01 0.72±0.01 617.5±27.1 7.9±0.4	0.66±0.01 0.57±0.02 1158.4±65.2 342.6±283.8	0.77±0.01 0.56±0.01 1051.6±31.0 4.4±0.2	0.56±0.02 0.45±0.03 1826.2±100.4 11.7±1.7	0.53±0.01 0.47±0.02 2442.4±115.6 570.3±1084.5	0.75±0.01 0.62±0.02 871.6±32.0 6.9±0.6
PC-Stable	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T (s) \end{array}$	0.72±0.01 0.50±0.01 1256.9±26.9 383.4±65.6	0.73±0.01 0.37±0.01 1189.2±37.0 167.7±14.5	0.54±0.01 0.18±0.00 1771.5±42.8 176.4±19.8	0.82±0.01 0.53±0.02 792.2±36.3 309.2±57.3	0.54±0.01 0.25±0.00 1700.1±49.5 165.9±12.4	0.67±0.01 0.38±0.01 1520.4±42.6 453.8±103.8	0.80±0.01 0.61±0.01 724.8±24.7 172.6±14.6	0.61±0.01 0.23±0.01 2144.2±47.3 569.4±150.5	0.45±0.36 0.29±0.23 1718.2±41.7 43121.0±35902.9	0.68±0.01 0.44±0.01 1165.8±40.2 169.7±15.3
fGES	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.57±0.01 0.49±0.02 2226.1±48.7 769.1±104.8	0.51±0.00 0.42±0.01 2301.3±17.7 823.2±26.6	0.43±0.00 0.42±0.00 2320.3±16.1 787.0±19.9	0.55±0.00 0.42±0.01 2337.9±29.4 907.6±35.2	0.43±0.00 0.42±0.01 2326.9±20.0 813.5±28.0	0.47±0.01 0.42±0.01 2512.1±37.0 689.4±44.5	0.58±0.00 0.43±0.01 2331.5±25.1 832.4±27.6	0.57±0.01 0.48±0.02 2259.6±49.2 646.5±28.6	0.53±0.00 0.48±0.01 3143.3±51.1 1257.0±74.2	0.50±0.00 0.40±0.01 2396.8±22.3 770.2±22.8
Impro. ratio	$F1^-$ $F1^{\rightarrow}$ SHD	7.4% 19.0% 29.2%	4.7% 16.3% 63.5%	14.6% 29.5% 89.6%	0.2% 6.8% 0.7%	14.6% 30.0% 88.9%	13.2% 6.7% 38.1%	-0.1% -3.3% 25.2%	10.4% 12.3% 17.9%	12.5% 28.2% 22.4%	3.2% 26.0% 32.3%

351 352

353 **Baselines and Settings** Six baselines are considered, including existing state-of-the-art methods 354 and PEF-based methods. Specifically, given the rich literature on BN structure learning (Kitson 355 et al., 2023), we collect 15 existing methods, including score-based, constraint-based, and hybrid methods, and conduct a preliminary testing of them (see Appendix C.2). The results indicate that 356 fGES and PC-Stable are the only methods capable of maintaining F1 scores of around 0.5 within rea-357 sonable time when the variable number reaches 1000. Therefore, we choose these two as baselines. 358 Besides, fGES and PC-Stable are also integrated into the estimation step of PEF, resulting in two 359 new baselines: P/fGES and P/PC-Stable. Furthermore, to validate the effectiveness of Auto-SLE, 360 we consider two alternative ensemble construction approaches: (i) default SLE, which contains the 361 default fGES and PC-Stable, as well as variants with randomly chosen hyperparameter values for 362 each of them; (ii) random SLE, which contains two variants with randomly chosen hyperparameter 363 values for each of fGES and PC-Stable. Both of these SLEs consist of four member algorithms 364 (same as our constructed SLE) and are integrated into PEF, yielding two baselines P/SLE(D) and 365 P/SLE(R).

366 To prevent the compared methods from running prohibitively long, a runtime limit of 24 hours is set 367 on each testing problem. All the experiments are conducted on a Linux server with an Intel Gold 368 6336Y CPU @ 2.40GHz, 96 cores, and 768GB of memory. Precise details of the experimental setup, 369 including the preliminary testing results, SLE construction, benchmarks, metrics, and baselines, are 370 in Appendix C. The codes for repeating our experiments are available in the supplementary.

371

372 4.2 **RESULTS AND ANALYSIS** 373

374 The partition step of PEF typically results in subproblems with $5\% \sim 10\%$ variables of the original 375 problem. For testing problems with 1000 and 10,000 variables, the subproblems have $50 \sim 100$ and $500 \sim 1000$ variables, respectively, which are problem sizes covered by training data. However, as 376 the number of variables increases further (e.g., to 30,000), the subproblems will be much larger than 377 the training problems. We first examine the performance on testing problems with 1000 and 10,000

402

403

404

405

380												
381	Problem (V , E)		Alarm (10027, 13713)	Asia (10000, 11000)	Cancer (10000, 8800)	Child (10000, 13750)	Earthquake (10000, 8800)	Hailfinder (10024, 12996)	Healthcare (10003, 14148)	Mildew (10010, 14472)	Pigs (10143, 14978)	Survey (10002, 11003)
382	P/SLE	$\begin{array}{c} F1^- \\ F1^- \end{array}$	0.80 0.66	0.94 0.72	0.96 0.92	0.85 0.61	0.96 0.92	<u>0.81</u> <u>0.66</u>	0.88 0.59	0.69 0.55	0.76 0.69	<u>0.90</u> <u>0.77</u>
383		SHD T (s)	<u>6366</u> 95.5	<u>1717</u> 35.6	<u>705</u> 24.4	<u>6017</u> 80.8	<u>683</u> 24.2	5400 828.6	<u>5361</u> 50.0	10911 360.3	7468 2120.5	3334 39.2
384	D/SI E/D)	$F1^{-}$ $F1^{\rightarrow}$	0.34 0.28	0.37 0.30	0.29 0.28	0.42 0.32	0.31 0.30	0.34 0.30	0.43 0.31	0.30 0.23	0.33 0.29	0.35 0.28
385	FISEE(D)	SHD T (s)	42805 4084.4	36831 2754.9	40464 3308.6	35660 3383.0	37413 2810.6	40151 3865.9	37297 3739.6	50139 4532.6	41958 7040.1	39476 2931.4
386	D/OLD (D)	$F1^-$ $F1^-$	0.49 0.25	0.59 0.25	0.42 0.15	0.69 0.41	0.43 0.14	0.45 0.28	0.71 0.28	0.38	0.02	0.55 0.19
387	P/SLE(R)	SHD T (s)	23673 676.8	18141 386.7	26652 667.9	12844 343.2	26065 538.6	25024 628.6	15483 319.6	32409 5563.9	15551 7048.2	21367 379.0
388		$F1^-$ $F1^{\rightarrow}$	0.59	0.71	0.54	0.77	0.55	0.55	0.80	0.44	0.12	0.66
389	P/PC-Stable	SHD T (s)	16842 210.1	11262 145.0	17280 147.3	8932 145.2	16869 133.3	17125 219.0	10815 156.8	24847 508.7	16005 7005.0	14046 125.4
390	DISCES	$F1^-$ $F1^{\rightarrow}$	0.34 0.28	0.37 0.30	0.29 0.28	0.42 0.32	0.31 0.30	0.34 0.30	0.43 0.31	0.30 0.23	0.33 0.29	0.35 0.28
391	P/IGES	SHD T (s)	42805 4037.6	36831 2796.0	40464 3331.4	35660 3369.7	37413 2833.4	40151 3669.7	37297 3755.3	50139 4439.0	41958 4179.0	39476 2905.7
392		F1 [−]		-	-		-	-	-	-		-
393	PC-Stable	SHD T (s)	- 86400.0	- 86400.0	- 86400.0	- 86400.0	- 86400.0	86400.0	- 86400.0	86400.0	86400.0	86400.0
394		F1 [−]	-	-	-	-	-	-	-	-	-	-
395	fGES	F1 SHD T (s)	- - 86400.0	- - 86400.0	- 86400.0	- - 86400.0	- - 86400.0	- - 86400.0	- - 86400.0	- - 86400.0	- - 86400.0	- 86400.0
396		F1-	36.0%	31.4%	78.8%	10.8%	75.7%	47.6%	10.4%	54.5%	133.0%	36.0%
397	Impro. ratio	F1 ⁻³ SHD	115.6% 62.2%	127.2% 84.8%	225.4% 95.9%	33.9% 32.6%	208.0% 96.0%	94.7% 68.5%	70.9% 50.4%	135.0% 56.1%	135.8% 52.0%	180.0% 76.3%

378 Table 2: Results on problems with 10,000 variables. "-" means a solution is not found within the 379 budget of 24 hours. The best performance in terms of accuracy is marked in **bold** and an underline.

variables to answer RQ1. Specifically, based on each network selected from *bnlearn*, we generate 10 testing problems (with different random seeds), test each method on these problems, and report the mean \pm std in terms of the evaluation metrics as well as stastical test results in Table 1. For the performance evaluation involving 10,000 variables, due to the very long runtime of the baselines, we generate one test problem based on each network and report the testing results in Table 2.

406 The first observation from Table 1 is that P/SLE consistently achieves significantly higher accuracy 407 across all three metrics compared to the baselines, except for the F1 metrics on Healthcare testing 408 problems, where P/SLE performs slightly worse than PC-Stable. Table 2 shows that the superior-409 ity of P/SLE becomes more pronounced on testing problems with 10,000 variables. In call cases 410 it achieves substantially higher accuracy than all baselines across all accuracy metrics. Notably, 411 compared to the best performance achieved by the baselines, P/SLE often achieves improvements 412 in F1 Adjacent of over 30% and up to 133%, and improvements in F1 Arrowhead and SHD of over 50% and even up to 225%. Since the difference between P/SLE and P/PC-Stable (P/fGES) lies in 413 the use of a SLE in the estimation step instead of a single algorithm, the consistent advantages of 414 P/SLE over them confirm that using SLEs can stably achieve high learning accuracy across sub-415 problems. On the other hand, we also observe that P/SLE(D) and P/fGES obtain identical learning 416 accuracy. This is because, in the default ensemble, the output of fGES always has the best BIC score 417 among all member algorithms, thus making it consistently being chosen as the final output. While 418 P/SLE(R), with a randomly constructed SLE, can avoid this issue, it fails to achieve satisfactory 419 learning accuracy, which in some cases is even worse than P/PC-Stable and P/fGES that do not use 420 SLEs. Therefore, the advantages of P/SLE over P/SLE(D) and P/SLE(R) highlight the effectiveness 421 of Auto-SLE in producing high-quality SLEs with complementary member algorithms.

422 The second observation is that on testing problems with 1000 variables, P/PC-Stable and P/fGES 423 often achieve higher learning accuracy than PC-Stable and fGES, respectively. Moreover, when 424 the number of variables reaches 10,000, PC-Stable and fGES are unable to find solutions within 24 425 hours, while P/PC-Stable and P/fGES are able to. These findings show that PEF framework can 426 indeed enhance the capabilities of handling large BNs, which is consistent with the observations 427 in (Gu & Zhou, 2020). Finally, all PEF-based methods generally consume much less runtime than 428 non-PEF-based methods, attributed to the underlying divide-and-conquer strategy. Among PEF-429 based methods, P/SLE often has the shortest or close to the shortest runtime, and for most testing problems, it consistently outputs the final solution within a reasonable time (less than 1000 seconds). 430 In summary, all the above findings affirmatively answer RQ1, i.e., the SLE constructed by Auto-SLE 431 substantially improves PEF in learning the structure of large BNs.



Figure 2: Performance curves on Alarm and Asia problems with up to 30,000 variables. SHD and runtime are plotted on log scale.

Table 3: Testing results on Yeast and WS problems. "P/SLE-b" represents the best performance achieved among P/SLE(R), P/SLE(D), P/fGES, and P/PC-Stable. The best performance in terms of accuracy metrics is marked with an underline, and the performance that is not significantly different from the best performance is indicated in **bold**.

$\begin{array}{c} \text{Problem} \\ (V , E) \end{array}$		WS (1000, 2000)	WS (10000, 20000)	Yeast (4441, 12873)
P/SLE	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.80±0.01 0.51±0.01 <u>1240.2±47.7</u> 9.8±0.3	0.79 0.49 12611 83.2	0.106 0.084 20912 7105.4
P/SLE-b	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.79±0.01 0.51±0.01 1252.4±32.2 3.7±0.3	0.72 0.40 15760 363.1	0.003 0.001 33162 5771.32
PC-Stable	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.75±0.01 0.42±0.01 1437.4±21.2 176.2±6.5	- - 86400.0	86400.0
fGES	$\begin{array}{c} F1^{-} \\ F1^{\rightarrow} \\ SHD \\ T \ (s) \end{array}$	0.67±0.00 0.43±0.01 2401.0±22.3 1433.3±61.4	- - 86400.0	0.066 0.056 29847 178257.4
Impro. ratio	$\begin{array}{c} F1^- \\ F1^- \\ SHD \end{array}$	2.2% -0.6% 1.0%	10.6% 21.1% 20.0%	60.6% 50.0% 29.9%

4.3 **GENERALIZATION TO LARGER PROBLEMS**

We now investigate RQ2. Specifically, we generate Alarm and Asia testing problems with 20,000 and 30,000 variables, and plot in Figure 2 the performance of the compared methods as the variable number ranges from 1000 to 30,000 (detailed results can be found in Appendix C.5). Note that when the variable number exceeds 20,000, the subproblems resulted from the partition step of PEF would be much larger in size than the training problems. It can be seen from Figure 2 that P/SLE generalizes well to larger problems, maintaining relatively stable learning accuracy. In contrast, the performance of all baselines deteriorates rapidly as the number of variables increases.

GENERALIZATION TO PROBLEMS WITH NO BLOCK STRUCTURE 4.4

The above results have demonstrated that P/SLE can stably achieve high learning accuracy for net-works with a block structure. Of course, there are many real-world networks without any block

486 structure (Olesen & Madsen, 2002). Although PEF-based methods are not specifically designed 487 for such networks, it is worth testing P/SLE on them to provide a complete spectrum of its perfor-488 mance. Specifically, we apply P/SLE for gene expression data analysis, a traditional application of 489 structure leanring for BN. We use the largest publicly available gene dataset Yeast (Schaffter et al., 490 2011) involving 4,441 nodes and 12,873 edges, where the underlying networks are commonly referred to as gene regulatory networks. Besides, we generate small-world networks with 1000 and 491 10,000 nodes, using igraph (Csardi et al., 2006) based on the Watts-Strogatz (WS) model (Watts & 492 Strogatz, 1998). We choose the WS model because: (i) it has no block structure, and (ii) it is not 493 included in the *bnlearn* repository, thereby enabling evaluation of the generalization of P/SLE to 494 network characteristics beyond the training set. As before, we generate 10 testing problems based 495 on the network with 1000 nodes and one testing problem based on the network with 10,000 nodes. 496 The testing results are presented in Table 3. 497

One can observe that P/SLE still achieves competitive learning accuracy on these networks. On 498 WS problems with 1000 variables, it always achieves the best performance or the performance not 499 significantly different from the best, across all accuracy metrics. On WS problems with 10,000 vari-500 ables and Yeast, the advantages of P/SLE become pronounced, similar to the previous observations 501 on networks with block structures. In summary, these findings show the generalization ability of the 502 constructed SLE across network characteristics beyond the training problem set. 503

5 CONCLUSION

504

505 506

507

521 522

In this work, we introduced the idea of using SLEs for BN structure learning and proposed Auto-SLE, an automatic approach that can largely reduce human efforts in building high-quality SLEs. 508 Extensive experiments showed that our method P/SLE could consistently achieve high accuracy in 509 learning large BNs and generalize well across problem sizes and network characteristics. 510

511 **Limitations** There are two main limitations of this work. First, while using an SLE by running its 512 member algorithms in parallel would not significantly increase the wall-clock runtime compared to 513 running a single algorithm, executing them sequentially in the absence of multi-core compute leads 514 to significantly longer runtime. A potential solution is to train a selection model that predicts the 515 best-performing algorithm in the SLE for a given problem, and runs that algorithm only. Second, the 516 generalization ability of P/SLE relies on a training set with diverse networks of varying sizes. If such a set cannot be collected in practice, then it may not generalize well. Moreover, as aforementioned, 517 PEF-based methods are most suitable for learning BNs with a block structure to some extent. For 518 BNs with no block structure at all, P/SLE may not be the best choice, and in these cases the SLE is 519 more suitable as a standalone learning method rather than being integrated into the PEF framework. 520

- REFERENCES
- 523 Emmanuel Abbe, Afonso S Bandeira, and Georgina Hall. Exact recovery in the stochastic block 524 model. IEEE Transactions on Information Theory, 62(1):471-487, 2015. 525
- Edo M Airoldi, David Blei, Stephen Fienberg, and Eric Xing. Mixed membership stochastic block-526 models. Advances in Neural Information Processing Systems, 21, 2008. 527
- 528 Hirotugu Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic* 529 Control, 19(6):716–723, 1974. 530
- Bryon Aragam and Qing Zhou. Concave penalized estimation of sparse gaussian bayesian networks. 531 *The Journal of Machine Learning Research*, 16(1):2273–2328, 2015. 532
- 533 Albert-László Barabási and Eric Bonabeau. Scale-free networks. Scientific american, 288(5):60-69, 534 2003.
- Remco Ronaldus Bouckaert. Bayesian Belief Networks: From Construction to Inference. PhD 536 thesis, 1995. 537
- Ruichu Cai, Siyu Wu, Jie Qiao, Zhifeng Hao, Keli Zhang, and Xi Zhang. Thps: Topological hawkes 538 processes for learning causal structure on event sequences. IEEE Transactions on Neural Networks and Learning Systems, 2022.

540	
540	David Maxwell Chickering. Learning equivalence classes of bayesian-network structures. <i>Journal</i> of Machine Learning Research, 2:445–498, 2002.
542	
543 544	Max Chickering, David Heckerman, and Chris Meek. Large-sample learning of bayesian networks is np-hard. <i>Journal of Machine Learning Research</i> , 5:1287–1330, 2004.
545	
546	Diego Colombo, Marloes H Maathuis, et al. Order-independent constraint-based causal structure learning. <i>Journal of Machine Learning Research</i> , 15(1):3741–3782, 2014.
547	
548 549	Gregory F Cooper and Edward Herskovits. A bayesian method for the induction of probabilistic networks from data. <i>Machine learning</i> , 9:309–347, 1992.
550	
551	Francesco Croce and Matunias Hein. Renable evaluation of adversarial robusiness with an ensem-
552	Machine Learning, pp. 2206–2216, 2020.
222	Cabor Ceardi Tamas Nenuez, et al. The igraph software package for complex network research
555	InterJournal, complex systems, 1695(5):1–9, 2006.
556	Paul Frdős Alfréd Rényi et al. On the evolution of random graphs Publ math inst hung acad
557 558	sci, 5(1):17–60, 1960.
559	Nir Friedman, Iftach Nachman, and Dana Pe'er. Learning bayesian network structure from massive
560	datasets: The" sparse candidate" algorithm. <i>arXiv preprint arXiv:1301.6696</i> , 2013.
561	
562	Maxime Gasse, Alex Aussem, and Haytham Elghazel. A hybrid algorithm for bayesian network
563	structure learning with application to multi-label learning. Expert Systems with Applications, 41 (15):6755, 6772, 2014
564	(15).0755-0772, 2014.
565	Jiaying Gu and Qing Zhou. Learning big gaussian bayesian networks: Partition, estimation and
566	fusion. Journal of Machine Learning Research, 21(1):6340–6370, 2020.
567	Paul W Holland Kathryn Blackmond Laskey and Samual Lainhardt, Stochastic blackmodals; First
568 569	steps. Social Networks, 5(2):109–137, 1983.
570	Neville Kenneth Kitson, Anthony C Constantinou, Zhigao Guo, Yang Liu, and Kiattikun Chobtham.
571	A survey of bayesian network structure learning. Artificial Intelligence Review, pp. 1-94, 2023.
572	Pedro Larranaga Cindy MH Kujipers Roberto H Murga and Vosu Vurramendi Learning bayesian
573	network structures by searching for the best ordering with genetic algorithms <i>IEEE Transactions</i>
574	on Systems. Man. and Cybernetics. 26(4):487–493, 1996.
575	
576	Marius Lindauer, Katharina Eggensperger, Matthias Feurer, André Biedenkapp, Difan Deng, Car-
577	oin Benjamins, 11m Kunkopi, Kene Sass, and Frank Hutter. Smac3: A versatile bayesian opti-
578	(1):2475 2483 2022
579	(1).27.5-2705, 2022.
580	Shengcai Liu, Ke Tang, Yunwei Lei, and Xin Yao. On performance estimation in automatic al-
581	gorithm configuration. In Proceedings of the 34th AAAI Conference on Artificial Intelligence,
582	AAAI'2020, pp. 2384–2391, New York, NY, Feb 2020.
583	Dimitris Margaritis et al Learning Revesion network model structure from date. DhD thesis School
584	of Computer Science Carnegie Mellon University Pittsburgh PA USA 2003
585	or comparer belonce, curregie menon emversity i fusburgh, 17t, 007t, 2005.
586	George L. Nemhauser, Laurence A. Wolsey, and Marshall L. Fisher. An analysis of approximations
587	for maximizing submodular set functions - I. Mathmatical Programming, 14(1):265–294, 1978.
588	Kristian G. Olesen and Anders I. Madsen Maximal prime subgraph decomposition of bayesian
589	networks. <i>IEEE Transactions on Systems, Man. and Cybernetics. Part B</i> , 32(1):21–31, 2002.
590	······································
591	Judea Pearl. Bayesian networks: a model of self-activated memory for evidential reasoning. In
592	Proceedings of the /th Conference of the Cognitive Science Society, pp. 329–334, 1985.
232	

Judea Pearl. Probabilistic reasoning in intelligent systems: networks of plausible inference. 1988.

- Jose M Pena. Learning gaussian graphical models of gene networks with false discovery rate control. In *European conference on evolutionary computation, machine learning and data mining in bioinformatics*, pp. 165–176. Springer, 2008.
- Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables and more: The fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *International Journal of Data Science and Analytics*, 3:121–129, 2017.
- Joseph D Ramsey, Kun Zhang, Madelyn Glymour, Ruben Sanchez Romero, Biwei Huang, Imme
 Ebert-Uphoff, Savini Samarasinghe, Elizabeth A Barnes, and Clark Glymour. Tetrad a toolbox
 for causal discovery. In *Proceedings of the 8th International Workshop on Climate Informatics*,
 2018.
- Thomas Schaffter, Daniel Marbach, and Dario Floreano. Genenetweaver: in silico benchmark generation and performance profiling of network inference methods. *Bioinformatics*, 27(16):2263–2270, 2011.
- 609 610 Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, pp. 461–464, 1978.
- Marco Scutari. Learning bayesian networks with the bnlearn r package. *Journal of Statistical Software*, 35:1–22, 2010.
- Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, Prediction, and Search.* MIT Press, 2000.
- Ioannis Tsamardinos, Constantin F Aliferis, and Alexander Statnikov. Time and sample efficient
 discovery of markov blankets and direct causal relations. In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 673–678, 2003a.
- Ioannis Tsamardinos, Constantin F Aliferis, Alexander R Statnikov, and Er Statnikov. Algorithms
 for large scale markov blanket discovery. In *FLAIRS conference*, volume 2, pp. 376–380. St. Augustine, FL, 2003b.
- Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing
 bayesian network structure learning algorithm. *Machine learning*, 65:31–78, 2006.
- Duncan J Watts and Steven H Strogatz. Collective dynamics of 'small-world' networks. *nature*, 393 (6684):440–442, 1998.
- Sandeep Yaramakala and Dimitris Margaritis. Speculative markov blanket discovery for optimal feature selection. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*, pp. 4–pp. IEEE, 2005.
 - Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Rong Zhu, Andreas Pfadler, Ziniu Wu, Yuxing Han, Xiaoke Yang, Feng Ye, Zhenping Qian, Jingren Zhou, and Bin Cui. Efficient and scalable structure learning for bayesian networks: Algorithms and applications. In *Proceedings of the 37th IEEE International Conference on Data Engineering*, pp. 2613–2624, 2021.
- 638 639

632

633

634

635

636

637

613

A THE PARTITION-ESTIMATION-FUSION (PEF) FRAMEWORK

640 641 642

643

644

645

646

647

This section presents the implementation details of the partition-estimation-fusion (PEF) Gu & Zhou (2020) framework. PEF consists of the following three steps.

- **Partition**: The nodes are divided into clusters using a modified hierarchical clustering (MHC) algorithm.
- Estimation: An existing structure learning method is applied to estimate a subgraph on each cluster of nodes.

661 662

663

664

665 666

667

672 673

677 678

684 685

687

690

Algorithm 2: Modified Hierarchical Clustering 649 Hierarchical clustering given the dissimilarity matrix $D = (d(i, j))_{d \times d}$; 650 Generate the dendrogram T_D of the hierarchical clustering; 651 Choose p by Eq. (9) and l by Eq. (10); 652 Relabel clusters in $C \leftarrow C_l$ so that $S_1 \leq \cdots \leq S_{d-l}$; 653 while |C| > p do 654 $(i^*, j^*) \leftarrow \arg\min_{(i,j)} \{ d(C_i, C_j) : i < j \land j > p \};$ 655 $C_{i^*} \leftarrow C_{i^*} \cup C_{i^*}, C \leftarrow C \setminus \{C_{i^*}\};$ 656 end 657 return $C = \{C_1, C_2, \dots, C_p\};$ 658

• Fusion: Merge estimated subgraphs into one DAG containing all the nodes.

Suppose we have observed m iid observations of random variables X_1, \ldots, X_d , denoted as $D \in$ $\mathbb{R}^{m \times d}$. Denote the *i*-th column of D as \mathbf{x}_i . Given D, the goal is to learn a DAG structure $\mathcal{G} = (V, E)$ that accurately reflects the conditional dependencies among X_1, \ldots, X_d .

A.1 PARTITION

The partition step (P-step) of the PEF involves partitioning nodes into clusters. From this procedure, 668 p clusters, denoted as C_i for i = 1, 2, ..., p, are generated. This utilizes a modified hierarchical clus-669 tering (MHC) approach, equipped with average linkage, that autonomously determines the number 670 of clusters p. The distance between two nodes i and j in PEF is defined by a specific equation, 671

$$d(i,j) = 1 - |r_{ij}| \in [0,1] \tag{7}$$

where $r_{ij} = cor(X_i, X_j)$ represents the correlation between X_i and X_j for i, j = 1, 2, ..., d. 674 The correlation is calculated using covariance $cov(\mathbf{x}_i, \mathbf{x}_i)$ and standard deviations $\sigma_{\mathbf{x}_i}, \sigma_{\mathbf{x}_i}$, and the 675 following equation. 676

$$cor(X_i, X_j) = \frac{cov(\mathbf{x}_i, \mathbf{x}_j)}{\sigma_{\mathbf{x}_i} \sigma_{\mathbf{x}_j}}$$
(8)

679 PEF mandates that the minimum cluster size should be 0.05d. For each $h = 0, 1, ..., d-1, C_h$ desig-680 nates the clusters formed during the h-th iteration of bottom-up hierarchical clustering. Specifically, 681 $C_0 = \{\{1\}, \{2\}, \dots, \{d\}\}$ consists of p singleton clusters, while $C_{p-1} = \{\{1, 2, \dots, d\}\}$ denotes a single cluster encompassing all d nodes. Let p_i signify the count of big clusters in C_i . PEF selects a 682 particular p according to the following equation: 683

$$p = \min\{p_{max}, \max_{0 \le i \le d-1} p_i\}$$

$$\tag{9}$$

686 where $p_{max} \leq 20$ represents a user-defined maximum count of big clusters.

Let l represent the topmost level on the dendrogram housing p big clusters, expressed as the follow-688 ing equation: 689

$$l = \underset{0 \le i \le d-1}{\arg\max\{i : p_i = p\}}$$
(10)

691 Clusters in C_l are then relabeled in descending order based on their size, so $S_1 \ge S_2 \ge \cdots \ge S_{d-l}$, 692 with $S_i = |C_i|$. The first p clusters are subsequently identified as the primary big clusters of interest. 693 PEF proceeds to allocate the leftover small clusters to these p big clusters, which is accomplished 694 by repetitively merging the two nearest clusters, provided one is a small cluster. The pseudocode of the MHC algorithm is detailed in Algorithm 2, where the dissimilarity matrix $D = (d(i, j))_{d \times d}$ is 696 calculated by Eq. (7). In the experiments, we limit the size of the cluster to be less than 10% of the 697 original problem to prevent the occurrence of excessively large subproblems.

698 699

700

- A.2 ESTIMATION
- During the estimation step (E-step), the PEF determines the structure of each subgraph individually. 701 Within the PEF framework, this step acts like a blackbox, allowing users to employ any structure

702 learning algorithm to estimate subgraphs without needing in-depth knowledge of its technical de-703 tails. Typically, this step yields p partial DAGs (PDAGs). It is noteworthy that both DAGs and 704 complete PDAGs (CPDAGs) are subsets of PDAGs. If the time complexity of a structure learning 705 technique surpasses $O(d^2)$, the time required to learn small subgraphs during the E-step becomes 706 considerably less than that needed to estimate an entire DAG. Assuming that during the partition step nodes were divided into p clusters C_1, C_2, \ldots, C_p and the duration to learn a PDAG on C_i 707 is t_i , parallelizing the learning of p subgraphs across p cores can reduce the E-step duration to 708 $\max\{t_i : i = 1, 2, \dots, p\}.$ 709

710

719 720

721 722

723 724

725

726 727

728

729

741

747 748

753

711 712

In the fusion step (F-step), a hybrid methodology is employed to learn the full DAG structure from the estimated subgraphs (obtained in the E-step). This step unfolds in two stages. First, PEF generates a candidate edge set A to restrict the search space. Through a series of statistical tests, PEF discerns a subset, A*, comprising candidate edges between subgraphs. Consequently, the candidate edge set A consists of A* and all edges learned in each subgraph from the E-step. Second, PEF optimize the DAG structure by iteratively updating edges within set A based on a modified BIC score. The final output of the F-step is a DAG.

B PROOFS

A.3 FUSION

B.1 PROOF OF FACT 1

Fact 2. $Q(\cdot)$ is a monotone submodular function, i.e., for any two SLEs \mathbb{A} , $\mathbb{A}' \subset \Theta$ and any $\theta \in \Theta$, it holds that $Q(\mathbb{A}) \leq Q(\mathbb{A} \cup \mathbb{A}')$ and $Q(\mathbb{A} \cup \mathbb{A}' \cup \{\theta\}) - Q(\mathbb{A} \cup \mathbb{A}') \leq Q(\mathbb{A} \cup \{\theta\}) - Q(\mathbb{A})$.

Proof. By definition, $Q(\mathbb{A}, D) = \max_{\theta \in \mathbb{A}} Q(\theta, D)$, then it holds that $Q(\mathbb{A} \cup \mathbb{A}', D) = \max_{\theta \in \mathbb{A} \cup \mathbb{A}'} Q(\theta, D) \ge \max_{\theta \in \mathbb{A}} Q(\theta, D)$. The monotonicity holds.

To prove submodularity, we have

$$Q\left(\mathbb{A} \cup \{\theta\}\right) - Q\left(\mathbb{A}\right) = \frac{1}{|T|} \sum_{D \in T} [Q(\mathbb{A} \cup \{\theta\}, D) - Q(\mathbb{A}, D)] \text{ by definition of } Q(\cdot)$$

$$= \frac{1}{|T|} \sum_{D \in T} [Q(\theta, D) - Q(\mathbb{A}, D)]^{+}$$

$$\geq \frac{1}{|T|} \sum_{D \in T} [Q(\theta, D) - Q(\mathbb{A} \cup \mathbb{A}', D)]^{+} \text{ by monotonicity}$$

$$= Q\left(\mathbb{A} \cup \mathbb{A}' \cup \{\theta\}\right) - Q\left(\mathbb{A} \cup \mathbb{A}'\right).$$
(11)

The proof is complete.

742 B.2 PROOF OF THEOREM 1 743

Theorem 2. Using a hyperparameter optimization procedure that, in each iteration of Auto-SLE, returns $\hat{\theta}$ within ϵ -absolute error of the maximum of $\Delta(\theta|\mathbb{A})$, i.e., $\Delta(\hat{\theta}|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) - \epsilon$, then the quality $Q(\mathbb{A})$ of the SLE constructed by Auto-SLE is bounded by

$$Q(\mathbb{A}) \ge (1 - 1/e) \cdot Q(\mathbb{A}^*) - k\epsilon, \tag{12}$$

where \mathbb{A}^* is the optimal SLE to the SLE construction problem in Definition 1. Alternatively, if $\hat{\theta}$ is within ϵ -relative error of $\Delta(\theta^{\circ}|\mathbb{A})$, i.e., $\Delta(\hat{\theta}|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) \cdot (1-\epsilon)$, then the quality $Q(\mathbb{A})$ of the SLE constructed by Auto-SLE is bounded by

$$Q(\mathbb{A}) \ge (1 - 1/e^{1 - \epsilon}) \cdot Q(\mathbb{A}^*).$$
(13)

754 755 *Proof.* Order the candidate algorithms in \mathbb{A}^* as $\{\theta_1^* \dots \theta_k^*\}$. We denote $\mathbb{A} = \{\theta_1, \theta_2, \dots, \theta_k\}$ where θ_i is the algorithm added to \mathbb{A} in the *i*-th iteration of Auto-SLE. Let $\mathbb{A}_i = \{\theta_1, \dots, \theta_i\}$ and

756 let $\Delta(\theta|\mathbb{A}) = Q(\mathbb{A} \cup \{\theta\}) - Q(\mathbb{A})$ denote the performance improvement brought by adding θ to 757 758 In the first case where $\Delta(\hat{\theta}|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) - \epsilon$, for all positive integers $i < l \leq k$, we have: 759 760 $Q(\mathbb{A}^*) < Q(\mathbb{A}^* \cup \mathbb{A}_i)$ by monotonicity 761 $= Q(\mathbb{A}_i) + \sum_{i=1}^{\kappa} \Delta(\theta_j^* | \mathbb{A}_i \cup \{\theta_1^*, \dots, \theta_{j-1}^*\}) \quad \text{by telescoping sum}$ 762 763 764 $\leq Q(\mathbb{A}_i) + \sum_{\theta \in \mathbb{A}^*} \Delta(\theta | \mathbb{A}_i)$ by submodularity 765 766 (14)by definition of θ° $\leq Q(\mathbb{A}_i) + \sum_{\substack{0 \leq h \\ n \neq i}} \Delta(\theta^{\circ} | \mathbb{A}_i)$ 767 768 769 $\leq Q(\mathbb{A}_i) + \sum_{\theta \in \mathbb{A}^*} \left(Q(\mathbb{A}_{i+1}) - Q(\mathbb{A}_i) + \epsilon \right) \qquad \text{by } \Delta(\theta^\circ | \mathbb{A}) \leq \Delta(\hat{\theta} | \mathbb{A}) + \epsilon$ 770 771 $\leq Q(\mathbb{A}_i) + k \left(Q(\mathbb{A}_{i+1}) - Q(\mathbb{A}_i) + \epsilon \right).$ 772 773 Let $\delta_i = Q(\mathbb{A}^*) - Q(\mathbb{A}_i)$, which allows us to rewrite the above equation as $\delta_i \leq k(\delta_i - \delta_{i+1} + \epsilon)$, 774 then $\delta_{i+1} \leq (1 - \frac{1}{k})\delta_i + \epsilon$. Hence, we have 775 776 $\delta_l \le (1 - \frac{1}{k})^l \delta_0 + k\epsilon \cdot \left[1 - (1 - \frac{1}{k})^l\right]$ 777 778 $\leq e^{-l/k}\delta_0 + k\epsilon$ by $1 - x \le e^{-x}$ for all $x \in \mathbb{R}$ (15)779 $= e^{-l/k} (Q(\mathbb{A}^*) - Q(\mathbb{A}_0)) + k\epsilon$ 780 by that $\mathbb{A}_0 = \emptyset$ and $Q(\mathbb{A}_0) = 0$ $=e^{-l/k}Q(\mathbb{A}^*)+k\epsilon$ 781 782 Rearranging $\delta_l = Q(\mathbb{A}^*) - Q(\mathbb{A}_l) \leq e^{-l/k}Q(\mathbb{A}^*) + k\epsilon$, we have 783 784 $Q(\mathbb{A}_l) \ge (1 - e^{-l/k}) \cdot Q(\mathbb{A}^*) - k\epsilon.$ 785 (16)786 Since the SLE found by Auto-SLE is \mathbb{A}_k , then we have 787 788 $Q(\mathbb{A}_k) > (1 - 1/e) \cdot Q(\mathbb{A}^*) - k\epsilon.$ (17)789 790 In the second case where $\Delta(\hat{\theta}|\mathbb{A}) \geq \Delta(\theta^{\circ}|\mathbb{A}) \cdot (1-\epsilon)$. Similarly, for all positive integers $i < l \leq k$, 791 we have: 792 $Q(\mathbb{A}^*) \leq Q(\mathbb{A}_i) + \sum_{\theta \in \mathbb{A}^*} \left(Q(\mathbb{A}_{i+1}) - Q(\mathbb{A}_i) \right) / (1-\epsilon) \quad \text{by } \Delta(\theta^\circ | \mathbb{A}) \leq \Delta(\hat{\theta} | \mathbb{A}) / (1-\epsilon)$ 793 794 (18) $= Q(\mathbb{A}_i) + \frac{k}{1-\epsilon} \left(Q(\mathbb{A}_{i+1}) - Q(\mathbb{A}_i) \right).$ 796 797 Similarly, let $\delta_i = Q(\mathbb{A}^*) - Q(\mathbb{A}_i)$ and use the above procedure, we have 798 799 $Q(\mathbb{A}_l) \ge (1 - e^{-l(1-\epsilon)/k}) \cdot Q(\mathbb{A}^*).$ (19)800

Let l = k, we have

$$Q(\mathbb{A}_k) \ge (1 - 1/e^{1-\epsilon}) \cdot Q(\mathbb{A}^*).$$
⁽²⁰⁾

The proof is complete.

C DETAILS OF THE EXPERIMENTS

807 808

801

802 803 804

805 806

809 Throughout the experiments, we set a random seed as 1024, ensuring the reproducibility of our experiments. The codes for repeating our experiments can be found in the supplementary.

810 C.1 EVALUATION METRICS

812 Let M1 be the true MEC of the DAG and M2 be the estimated MEC, the F1 score for adjacencies (F1 Adjacent) is calculated as the harmonic mean of precision and recall: 2TP/(2TP+FP+FN), where 813 TP is the number of adjacencies shared by M1 and M2, FP is the number of adjacencies in M2 but 814 not in M1, and FN is the number of adjacencies in M1 but not in M2. The F1 score for arrowhead 815 (F1 Arrowhead) is calculated in a similar way. An arrowhead is taken to be in M1 and M2 for each 816 variable A and B such that $A \rightarrow B$ in both M1 and M2, and an arrowhead is taken to be in one MEC 817 but not the other when for each variable A and B such that $A \rightarrow B$ in one but $A \leftarrow B$ in the other, or 818 A–B (no directions) in the other, or A and B are not adjacent in the other. 819

The structural Hamming distance (SHD) is defined as the number of edge insertions, deletions or flips in order to transform the learned DAG to the ground truth.

823 C.2 PRELIMINARY TESTING

822

827

828

829

830

831

832

833

834

835

853

854

855

856

858

861

863

- We collect 15 existing methods, listed below.
 - Score-based combinatorial search methods: HC Chickering et al. (2004), TABU Bouckaert (1995), CCDr Aragam & Zhou (2015), fGES Ramsey et al. (2017)
 - Score-based continuous optimization methods: NOTEARS Zheng et al. (2018), GOLEM Zhu et al. (2021)
 - Constraint-based methods: PC-Stable Colombo et al. (2014), GS Margaritis et al. (2003), IAMB Tsamardinos et al. (2003b), Fast-IAMB Tsamardinos et al. (2003b), IAMB-FDR Pena (2008), Inter-IAMB Yaramakala & Margaritis (2005)
 - Hybrid methods: MMHC Tsamardinos et al. (2006), RSMAX2 Friedman et al. (2013), H2PC Gasse et al. (2014)

We collect the open-source implementations of these methods. Most of the implementations are collected from the *bnlearn* Scutari (2010) repository ¹; CCDr ², NOTEARS ³, GOLEM ⁴ are collected from Github; fGES and PC-Stable are collected from *TETRAD* Ramsey et al. (2018) ⁵.

839 We generate testing problems based on two random graph models, Erdös-Rényi (ER) Erdős et al. 840 (1960) and scale-free (SF) Barabási & Bonabeau (2003), where the edge number is set to be two 841 times the node number. Specifically, each testing problem involves 1000 variables, has Gaussian 842 noise, and the observation number m = 1000. Based on each graph model, 10 different testing 843 problems are generated, resulting in a total of 20 testing problems. Each method is then applied to these testing problems, with a runtime limit of 3600 seconds for each problem. Table 4 presents the 844 average performance of these methods on all 20 testing problems, measured by F1 score (Adjacent) 845 and runtime. The results indicate that fGES and PC-Stable are the top-performing methods, consis-846 tently maintaining F1 scores above 0.5 when the variable count reaches 1000. While HC and TABU 847 also achieve F1 scores above 0.5, their runtime is significantly longer, making them impractical for 848 testing problems involving 10,000 variables. As a result, they are excluded from the comparison 849 experiments. 850

851 C.3 AUTOMATIC CONSTRUCTION OF THE SLE

Algorithm Configuration Space During the SLE construction process, the algorithm configuration space is defined by two candidate algorithms and their hyperparameters.

• PC-Stable Colombo et al. (2014) with two hyperparameters: significance threshold of CI tests within the interval $\alpha \in [0.01, 0.2]$ and the search's maximum depth interval $m \in [1, 1000]$

862 ³https://github.com/xunzheng/notears (Apache-2.0 license)

⁸⁵⁹ ¹https://www.bnlearn.com/bnrepository (Creative Commons Attribution-Share Alike License).

²https://github.com/itsrainingdata/ccdrAlgorithm(license not specified)

⁴https://github.com/ignavierng/golem(Apache-2.0 license)

⁵https://www.ccd.pitt.edu/tools (GNU General Public License (GPL) v2 license)

864	Table 4: Preliminary testing results.	"-" means not	returning solutions	s within a time b	udget of 3600s.
-----	---------------------------------------	---------------	---------------------	-------------------	-----------------

Method	F1	runtime (s)
PC-Stable	0.71	576.21
fGES	0.66	625.86
HC	0.528	3170.22
TABU	0.528	3181.777
CCDr	0.392	834.781
MMHC	0.331	1210.805
RSMAX2	0.326	1416.733
GOLEM	0.32	3622.433
IAMB-FDR	0.271	2021.871
Inter-IAMB	0.189	2865.472
IAMB	0.157	2943.536
Fast-IAMB	0.107	3351.702
NOTEARS	0.099	3604.744
GS	-	3600.32
H2PC	-	3610.257

• fGES Ramsey et al. (2017) with three hyperparameters: structural penalty of the BIC score within interval $\lambda \in [1.0, 1000.0]$, the maximum number of parents for a single node during the search process within interval $m \in [1, 1000]$, and the option to use the faithfulness assumption or not.

We use the implementations of them from the causal discovery tool box *TETRAD* Ramsey et al. (2018) ⁶. SMAC (version 3) Lindauer et al. (2022) ⁷ is used as the hyperparameter optimization procedure.

Training Problem Set For Auto-SLE, we used the data generation method and randomly pro-duced 100 training problem instances. In order to ensure that the training data is diverse, compre-hensive, and representative, each instance was derived by the following steps: (i) selecting a base network at random from the *bnlearn* Scutari (2010) repository⁸; (ii) replicating it a random num-ber of times, while making sure the total node count is less than or equal to 1000, and connecting the replicas by adding 10% of edges between them randomly, ensuring the final network remains a DAG; (iii) utilizing the complete DAG and $X_i = \phi(X_{\text{pa}(i)}) + \varepsilon_i$, to produce 1,000 observational data records. Here the weights in the linear function ϕ and the standard deviations of Gaussian noises are sampled uniformly from $[-1, -0.5] \cup [0.5, 1]$ and [0, 1], respectively; (iv) re-scaling the observation data so that all data columns possess the same mean and standard deviation. These 100 training problem instances are independent of testing data and are exclusively for the construction of the SLE.

The Constructed SLE The constructed SLE contains four member algorithms, as running more iterations of Auto-SLE brings negligible improvement (smaller than 0.1) in SLE's performance on the training set. Specifically, the training performance ()in terms of the sum of F1 Adjacent and F1 Arrowhead) progress is: ite 1 (143.4) \rightarrow ite 2 (156.1) \rightarrow ite 3 (158.1) \rightarrow ite 4 (158.5) \rightarrow ite 5 (158.5). The SLE is detailed in Table 5. It is interesting to find that the SLE only contains fGES, meaning PC-Stable has not defeated fGES in the construction process.

⁶https://www.ccd.pitt.edu/tools (GNU General Public License (GPL) v2 license)

⁷https://automl.github.io/SMAC3/main(3-clause BSD license)

^{917 &}lt;sup>8</sup>https://www.bnlearn.com/bnrepository (Creative Commons Attribution-Share Alike License)

Table 5: The SLE contructed by Auto-SLE.

Member Algorithm	Candidate Algorithm	Hyperparameter Values
1	fGES	$\lambda = 5.872727477498224, m = 185$, and without faithfulness assumption
2	fGES	$\lambda = 20.60452402709974, m = 17$, and without faithfulness assumption
3	fGES	$\lambda = 2.537674798471765, m = 91$, and without faithfulness assumption
4	fGES	$\lambda = 5.624595350676138, m = 11$, and without faithfulness assumption

C.4 DEFAULT SLE AND RANDOM SLE

Default SLE and random SLEs are constructed in alternative ways other than Auto-SLE. The former contains the default fGES and PC-Stable, as well as variants with randomly chosen hyperparameter values for each of them; The latter contains two variants with randomly chosen hyperparameter values for each of fGES and PC-Stable. According to the causal discovery toolbox *TETRAD* Ramsey et al. (2018), the default hyperparameter values for PC-Stable is: $\alpha = 0.05$ and m = 1000; the default hyperparameter values for fGES is $\lambda = 1.0, m = 1000$, and without faithfulness assumption. Finally, SLE (Default) and SLE (Random) are shown in Table 6 and Table 7.

Table 6: Default SLE

Member Algorithm	Candidate Algorithm	Hyperparameter Values
1	PC-Stable	$\alpha = 0.05$, and $m = 1000$
2	fGES	$\lambda = 1.0, m = 1000$, and without faithfulness assumption
3	PC-Stable	$\alpha = 0.08399128452994686$, and $m = 850$
4	fGES	$\lambda=797.254519880674,$ $m=871,$ and without faithfulness assumption

Table 7: Random SLE.

Member Algorithm	Candidate Algorithm	Hyperparameter Values
1	PC-Stable	$\alpha = 0.08399128452994686$, and $m = 850$
2	fGES	$\lambda = 797.254519880674, m = 871$, and without faithfulness assumption
3	PC-Stable	$\alpha = 0.10744707122087944$, and $m = 980$
4	fGES	$\lambda=792.8350933385117,$ $m=456,$ and with faithfulness assumption

C.5 TESTING RESULTS ON LARGER ALARM PROBLEMS AND ASIA PROBLEMS

Table 8 and Table 9 present the testing results of the compared methods on the Alarm and Asia problems with up to 30,000 variables, respectively.

C.6 COMPUTE RESOURCE

Unless otherwise indicated, all experiments in this work are conducted on a Linux server equipped with an Intel(R) Xeon(R) Gold 6336Y CPU @ 2.40GHz, 96 cores, and 768GB of main memory. The system version is Ubuntu 22.04.2 LTS.

070			1		1	1	I	1	
972	the the d in	00 (37)			_	_	_		
913	e m und atec	-300	33.6	5.9 319 63	1.8 .0 84.0 9.8	2.5 5.8 99.0	/ / 12.9	0.00	0.00
974	the dic a	arm 007,	78 56 56	25 1 200 72	11 6 722 722	32 16 545 685	410, 1 675	864	864
975	in link	A (3)							
970	wo der) is	0							
977	$\frac{\text{un}}{05}$	2000	9. 1 8 8	4 L 2 2 4	5.5.0	6 1.0 9.3	8 5.0 9.9	0.0	0.0
970	ch an = 0	arm- 17,	345 44 5	29. 15. 345 760	32. 16. 7160	50. 3649 3649	5.5 4.8 3006 6549	8640	8640
979	ea ith p =	Ala (200				1			
981	d w vel	<u>_</u> @							
982	e le	371		1 5 8.	0,000	0.0	- 0.6	0.0	0.0
983	nar nce	rm-1 27, 1	79.7 66.2 98.6	34.] 28.3 4280 5262	49.25.367.684.	58.6 30.3 6845 210.2	34.] 28.3 3896	6400	6400
984	run is 1 fica	Ala (100	Ū		1	-	4.0	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	~
985	ics snit								
986), al letr letr	9000 2347		e 80 8.		200	. v. o. v.	0.0	0.00
987	HD m / HD	1. 8, 1.	78.9 63.6 91.(35.9 29 3610 2874	50. 26.(371.	59.3 31.3 483(178.	35.9 29.0 6108 2585	6400	53.6 45.3 2128 4889
988	, S acy st v	41s (902	Ū		1	-	, w , ,	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	6 8
989	ead cun c te	<u>ہ</u>							
990	who f ac ank	800	0 8 0 0	5 95 55	0.6	6 9 0 0 0	5.0 0.0	0.0	7 5.9 5.9
991	s of d-r	arm- 29, 1	8 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9	39. 32. 175	52. 27. 219	61. 32. 88.	39. 32. 160(8640	53. 46. 1942 5030
992	l A gne	(80 A]							
993	H in te	8 (7			_	_			
994	ent e ii xor	, 961	9.0	2.1 35 892 8.5	3.8 8.1 21.0	2.9 4.2 63.0	2.1 5.0 1.4	0.00	3.8 5.0 45.0 65.7
995	jac anc lco	Alarn 7030	∞ <u>≈</u> <u>4</u> 4	21°74	2 7 7 7 N	9 % 101 2	218 80 80 80	864	5 4 1 4 S
996	Wi	₹.0							
997	F1 rrfo 5 a	000 248)	• -		. o			0.	9.0
998	of] c pe	rm-6 31, 8	80.7 65.9 817. 34.9	43.6 35.5 804 570.4	55.2 28.7 1931 74.8	64.8 35.2 33.7	43.6 35.5 3043 504.(2400	53.8 45.5 4678 5760
999	ns vest din	Ala (603	e e		=	~~~~		×	5 T
1000	teri ne h cor	9 (1							
1001	in T	-500	7.1.2	222	6.1 0.1 0.1	6.00	5.7	0.0	53.9 23.9 23.9
1002	es, ed.	larm 032,	311 80 35 35	46 38 133 45:	55 950 950	64 35 667 29	46 133, 46	8640	53 46 121 166
1003	abl	A ()							
1004	ref orn	000 516)	_					04	_ ∞.
1005	00 v i is verf	m-4 3, 5;	30.2 55.1 538.(50.4 41.3 9357 62.3	60.2 51.9 526.0 24.6	68.1 38.0 869.	50.4 41.3 357.1 31.7	36.5 36.5 1115.1	54.2 45.9 707.
1006	st p	Alaı (403	~~~~~			4			
1007	be be	0 @							
1008	p to pro the	-300	4 9 0 9	6 I 6 5	5.0.2.6	1 7 - +	6 0. e.	9.9 1.2	9.9 9.8 3.8
1009	h u 10 m	larm 334,	79. 18 64. 18 18	51 65/ 76	60 14 14 12 14 14 14 14 14 14 14 14 14 14 14 14 14	8, 351,98	51 66 60	58 37 619 1909	53 44 736 545
1010	on	A 9							
1011	ns od ent	83)				_	_		
1012	eth ffer	m-2(0.3 5.8 26.0	9.5 9.2 361 0.6	8.0 8.9 9.5	3.2 5.9 13.0 5.5	9.5 9.2 0.4	4.8 2.5 60.0 97.3	5.7 7.2 97.0 69.1
1014	din m	Alar 2035	× • £ -	v 4 w 4	9 6 75	191	N 4 6 0	33 4 6 53 3 4 6	5 4 4 5 23 4 4 5
1015	n p ach								
1016	ları y e icaı	1000	~ n •	0 e 3 c	0.00	0.0	N 0 0	° ° ° ° °	2 × 0. 6.
1017	d b nifi	arm- 36, 1	83.8 74.5 478. 7.2	69. 62. 108 12.	72.7 42.8 6.1	75. 46. 899. 3.3	69 62 4.8	72. 50. 1261 273.	58.5 52.5 600.
1018	ine sig	Alå (10	-						
1019	ults bta 10t		ad	ad	ad	ad	ad	ad	ad
1020	resi e o is I		acent whe e (s)	aceni whe e (s)	acent whe e (s)	acent whe e (s)	acent whe e (s)	aceni whe e (s)	acent whe e (s)
1021	ng anc hat		Adjé Arrc D	Adjé Arrc D 1time	Arrc D 1time	Adj [£] Arrc D	Adj ⁱ D 11imu	Adjé D 1time	Adji Arrc D
1022	stii rm: ;e tl	E	E E E	EI E	E E E	E E E	Run SH	E E E	FI Ru Ru
1023	Té. anc	V		_	_	ole		0	
1024	be Tm		ш	E(D)	E(R)	-Stał	ES	tablé	
1025	abl¢ std rfo bld .		TS/c	TS/c	JS/	2/PC		C-S	GES
	<u>н д</u> а								- <u> </u>

0.9.1 Setuing results on Asia problems with up to 30000 variables. In terms of F1 Adjacent. F1 Arrowhead, SH2, and runtime. On each network, the mean dependent on 10 problems is reported. The mesh performance in terms of a cacuacy metrics is marked with an indefine. and 10 effortance (base) performance in terms of a cacuacy metrics is marked with an indefine. and 10 efformance (base) performance (base) performance in terms of a cacuacy metrics is marked with an indefine. and 10 efformation on 10 problems is reported. The mesh performance in terms of a cacuacy metrics is marked with an indefine. and 10 effort the mean metric in terms of a cacuacy metrics is marked with an indefine. and 10 effort the mean indefine and 10 effort the mean indefine. and 10 effort the mean indefine and 10 effort the mean indefine. The mean indefine and 10 effort the mean indefine and 10 effort the mean indefine. The mean indefine and 10 effort the mean indefine and 10 effort the mean indefine and 10 effort the mean indefine. The mean indefine and 10 effort the mean inde	1026	= e +					I			
else : Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and nutrine. On each network, the me of accuracy metrics is marked with an undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of a counter with a undefine, and account with a metric significant accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, and of accuracy metrics is marked with a undefine, andefine, and account accuracy metrics is marked with a u	1027	an ed j	000	-		0_	0~	0 0		0
dip Testing results on Asia problems with pro 30000 variables, in terms of F1 Adjacem, F1 Arrowhead, SHD, and nuntime. On each network, the continue to the less performance in terms of accumacy metrics is marked with a marked w	1028	anc	-30(2.5 0.9 87.0 87.8	0.2 / 2963).2 / 962.	5.1 4.6 131.	/ 29.0	~ ~ ~ 00	~ ~ ~ 00
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	1029	be di	Asia 0000	2294	- <i>.</i> ,	32	0° 4 ×	33	86	86
Ability of the set performance in terms of FI Adjacent, FI Arrowhead, SHD, and runtime. On each network of performance obtained with up to 30000 variables, in terms of FI Adjacent, FI Arrowhead, SHD, and runtime. On each network of performance interaction is provided with an under term in the set performance in terms of FI Adjacent, FI Arrowhead, SHD, and runtime. On each network of performance interaction is provided with an under term in the set performance interaction is provided with a multiple of manacci in terms of FI Adjacent, FI Arrowhead, SHD, and runtime. On each network of the multiple of manacci in terms of FI Adjacent, FI Arrowhead, SHD, and runtime. On each network of the multiple of manacci in terms of accurate performance interaction. The multiple of multiple of manacci in terms of the multiple of multin and multin and multiple of multin and multiple of multiple of m	1030	s ii t	0							
4000 End (1) End (2) End (2) <thend (2)<="" th=""> <thend (2)<="" th=""> <thend< td=""><td>1031</td><td>'orl 1de 5) i</td><td>00</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></thend<></thend></thend>	1031	'orl 1de 5) i	00							
400 Filt Adjacent, Filt Arrowhead, SHD, and runtime. On each n 400	1032	0.0	2000	23 .0 .1 8	7.1 180 180	60.0 60.0	0 5.1 11.7	.5 19.0 7.6	0.00	0.00
400 51 Adjacent, F1 Arrowhead, SHD, and mutime. On each exponent of F1 Adjacent, F1 Arrowhead, SHD, and mutime. On each exponent of significance level processing results on Naia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and mutime. On each exponent of significance level processing results on Naia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and mutime. On each exponent of significance level processing results on Naia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and mutime. On each exponent of the loss processing results on Naia Processing Process	1033	a u u	sia- 000,	6 12 0 14 ⊟	24 57 301 723 723	45 19 452 715	252 26 61	8 6. 271 661	864	864
400 Constraint Constraint <td>1034</td> <td>$\frac{1}{p}$</td> <td>(20 A</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	1034	$\frac{1}{p}$	(20 A							
400 9.1 Fasting results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and nutnine. C 0 40 9.1 Fasting results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and nutnine. C 0 40 40 10	1035	hn e d w svel	1_6							
400 91 Flasting results on Asia problems with up to 30000 variables, in terms of F1 Adjacem, F1 Arrowhead, SHD, and runtime of performance that is not significantly different from the best performance in terms of F1 Adjacem, F1 Arrowhead, SHD, and runtime of performance that is not significantly different from the best performance in terms of F1 Adjacem, F1 Arrowhead, SHD, and runtime of a current so that is not significantly different from the best performance (according to a Wilcoxon signed-rank test with significance of according to a Wilcoxon signed-rank test with significance significance significance significance significance significance significance significance signit a wilcoxon sison of a wilcoxon sison of a wilcoxon sison of a w	1036	e le e le	100	~~ ° ~	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	× 1.0.8	7 0 0 0	8 90 11.0	0.0	0.0
Althor Adjacent, F1 Arrowhead, SHD, and runti able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, and runti and 10 able 9: Testing results on Asia problems with up to 30000 variables, in terms of according to a Wilcoxon signed-rank test with significs is to significantly different from the best performance (according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to the molecular signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with significs is to according to a Wilcoxon signed-rank test with signific significant signitant sist ana significant significant significant sist and sist	1037	me anc	ia-1 00,	93.8 71.8 35.	36. 368. 368.	58. 25. 814 374	71. 31. 126 143	36. 29. 683 2838	640 / / /	640 / / /
able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Arrowhead, SHD, and ru berformance that is not significantly different from the best performance (according to a Wilcoxon signed-rank test with significantly different from the best performance (according to a Wilcoxon signed-rank test with significant by a 2000 variables, in terms of F1 Arrowhead, SHD, and ru 0000 variables, in terms of F1 Arrowhead, SHD, and ru 00000 variables, in terms of F1 Arrowhead, SHD, and ru 00000 variables, in terms of F1 Arrowhead, SHD, and ru 00000000000000000000000000000000000	1038	nti is 1 fica	100			-	-	ю`,	~ ~	~
6400 6400 741 741 743 743 744 </td <td>1039</td> <td>l ru cs gni</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	1039	l ru cs gni								
Bible 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhead, SHD, 1 deformance that is not significantly different from the best performance in terms of accuracy m erformance that is not significantly different from the best performance in terms of accuracy m endomance that is not significantly different from the best performance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m endomance in terms of accuracy m erformance in terms of accuracy m erform accuracy m erformance in terms of accuracy m erform accuracy m erform accuracy m erform accuracy m erformance in terms of accuracy m erform accuracy m erformance in terms of accuracy m erform accuracy m erformance in terms of accuracy m erform erform accuracy m erform accuracy m erform accurac	1040	anc etri 1 si	000	~ ~ ° ~	- 2 8 5	400.0	60.7	2.08	0.0	0.0
Bible 9: Testing results on Asia problems with up to 30000 variables, in terms of FI Adjacent, FI Arrowhead, SHI dependenting to a Wilcoxon signed-rank test volume that is not significantly different from the best performance (according to a Wilcoxon signed-rank test volume that is not significantly different from the best performance in terms of accuracy efformance that is not significantly different from the best performance in terms of accuracy efformance in terms of a accuracy efformance in terms of a accuracy efforts in terms of accuracy efforts in terms of a accuracy efforts in terms of a accuracy efforts in terms of a accuracy effect in the effort of a state in terms of a accuracy efforts in terms of a accuracy efformance in terms of a accuracy efforts in terms of a ac	1041	vith D.	sia-9 00, 9	94.3 28.49	38. 31. 3171 2210	60. 26. 214	72. 33. 92.	38. 31. 171 1973	640	640
Problem Adar-1000 Adar-3000	1042	SHI acy st v	(Å 8			-		<i>с</i> о	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	~~~
Problem Asia-1000 Asia-3000 Variables, in terms of F1 Adjacent, F1 Arrowhea Provenance able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 Arrowhea 06000 20000 </td <td>1043</td> <td>d, 5 curs</td> <td>6</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	1043	d, 5 curs	6							
6400 9400 9400 9400 9400 9400 9400 9400 9400 9400 9400 9400 9400 10 9400 9400 9400 9400 9400 10 9400 10 9400 9400 9400 9400 9400 9400 10 10 940 10 940 940	1044	acc ank	880	9 9 9 9 9	1 4.1 8.6 8.6	.6 .2 66.0	6.7.6	.0 78.0 0.8	0.0	00.0
Bottomance that is not significantly up to 30000 variables, in terms of F1 Articles Articles Anti-Toto Anti-T	1045	of of d-r	Asia- 000,	24 C 130 24	33 4 251 153	61 133 133 133	73 33 830 68	41 33 251 [°] 137	864	864
Polylem Asia problems with up to 30000 variables, in terms of F1 Adjacent, F1 A djacent, F2 A ga, A di A djacent, F2 A ga, A d	1046	Arro ms	8							
Protein Provide and the set performance in terms of F1 Adjacent, F able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F 00000 variables, in terms of F1 Adjacent, F able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F 00000 variables, in terms of F1 Adjacent, F able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F 00000 variables, in terms of F1 Adjacent, F able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F 9441 able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Adjacent, F 9441 able 9: Testing results on Asia problems with up to 30000 state able of a without able of a state able o	1047	1 ⊿ teri siĝ	88			0				_
6001 4000 variables, in terms of FI Adjacent able 9: Testing results on Asia problems with up to 30000 variables, in terms of FI Adjacent 6001 able 9: Testing results on Asia problems with up to 30000 variables, in terms of FI Adjacent 6001 able 9: Testing results on Asia problems with up to 30000 variables, in terms of FI Adjacent 6001 able 9: Testing results on Asia problems with up to 30000 variables, in terms of FI Adjacent 6001 able 9: Testing results on Asia problems is reported. The best performance 6001 able 9: Testing results on Asia problems is reported. The best performance 6000 able 9: Testing results on Asia problems is reported. The best performance 6000 able 9: Testing results on Asia problems with up to 3000 variables, in terms of FI Adjacent 6001 able 9: Testing results on Asia problems is reported. The best performance 6000 able 9: Testing results on Asia problems is reported. The best performance 6000 able 9: Testing results on Asia problems is reported. The best performance 6000 able 9: Testing results on a sign able 9: Testing results on a sign able 9: Testing results 6000 able 9: Testing results on a sign able 9: Testing results 713 113 PARueline 2: Table 2: Tarrowhead 2: Table 2: Tarrowhead 2: Table 2: Ta	1048	in F	1-70(4.5 2.9 61.0 9.8	1.7 4.3 308 84.1	333 75 012.(5.0 4.9 89.0	-1.7 4.3 308.(/ / / 100.0	/ / / 400.0
6201 6201 <td< td=""><td>1049</td><td>ent ice cox</td><td>Asia 7000</td><td>96.01-</td><td>4 <i>ω</i> [2 %</td><td>9 11 0 0</td><td>Γ ω 8 4</td><td>4 ° 13 80</td><td>86</td><td>86</td></td<>	1049	ent ice cox	Asia 7000	96.01-	4 <i>ω</i> [2 %	9 11 0 0	Γ ω 8 4	4 ° 13 80	86	86
6.8001 with up to 30000 variables, in terms of F1 Ad able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 Ad 99001 efformance that is not significantly different from the best performance (according to a Vold) 99001 efformance obtained by each method on 10 problems is reported. The best performance (according to a Vold) 99001 old. Problem Asia-1000 Asia-2000 Asia-9000 Asia-6000 efformance obtained by each method on 10 problems is reported. The best performance (according to a Vold) 99001 99001 99001 old. Problem Asia-1000 Asia-2000 Asia-9000 Asia-9000 4900 5000.600 5000<	1050	jac nan Vil								
Problem Asia - 1000 Variation Log 5500 able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 69001 50000 variables, in terms of F1 able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 69001 50000 variables, in terms of F1 able 9: Testing results on Asia problems with up to 30000 variables, in terms of F1 69001 5000 50	1051	Ad orr a V	0000		. = 0	. 0 .			0.	0.
Formulation Asia and base of the set performance obtained by each method on 10 problems is reported. The best performance obtained by each method on 10 problems is reported. The best performance (according old. Problem Asia-1000 Asia-3000 Asia-3000 Asia-3000 Asia-5000 Asia-5	1052	F1 erf to	ia-6(00, 6	94.4 72.8 31.(47 38.4 502 524.2	65.4 29.1 688. 51.4	77.5 36.6 308. 24.4	47.0 38.4 5021 466.1	6	~ ~ ~ 5400
Formulation Asia - 2000	1053	of st p ing	(60C	S		×	<u>ې</u>	11	~ ×	×
6.001 event 2001 10000 variables, in terr able 9: Testing results on Asia problems with up to 30000 variables, in terr 69001 99001 99001 ad performance obtained by each method on 10 problems is reported. The erformance (accond) 99001 99001 99001 99001 addit Asia-1000 Asia-3000 Asia-3000 Asia-3000 4543 9901 addit Asia-1000 Asia-3000 Asia-3000 Asia-3000 4543 9001 addit FI Adjacent 554 341 734 724 724 PSLE FI Adjacent 737 724 734 724 734 724 PSLE(R) FI Adjacent 737 724 734 724 734 724 PSLE(R) FI Adjacent 737 724 734 724 734 724 PSLE(R) FI Adjacent 737 734 734 734 734 734 PSLE(R) FI Adjacent 737 734 734 73	1054	ns be: ord								
62001 42000 variables, in current contract obtained by each method on 10 problems is reported. T 99001 able 9: Testing results on Asia problems with up to 30000 variables, in d performance obtained by each method on 10 problems is reported. T 99001 addressing results on Asia problems with up to 30000 variables, in d performance that is not significantly different from the best performance (1000, 1100) 99001 addressing results on Asia-1000 Asia-2000 Asia-4000 Asia-4000 Problem Asia-1000 Asia-2000 Asia-4000 Asia-4000 PSLE FI Adjacent 54 73.4 94.7 94.7 94.7 PSLE FI Adjacent 54.7 73.4 94.7 94.7 94.7 PSLE FI Adjacent 54.7 73.4 94.7 94.7 94.7 PSLE(D) FI Adjacent 54.7 73.4 94.7 94.7 94.7 PSLE(D) FI Adjacent 53.7 73.4 73.4 73.4 73.4 PSLE(D) FI Adjacent 73.4 49.9 53.9 94.7 94.7 94.7 PSLE(R) </td <td>1055</td> <td>lhe acc</td> <td>5000</td> <td></td> <td>8; 4; 9; 6;</td> <td>5.00</td> <td>0.4.6.9</td> <td>8 4 9 </td> <td>0.0</td> <td>6.6 8.0 9.9</td>	1055	lhe acc	5000		8; 4 ; 9; 6;	5.00	0.4.6.9	8 4 9 	0.0	6.6 8.0 9.9
6000 variables, be successed by each method on 10 problems is reported erformance obtained by each method on 10 problems is reported tod. 0000 variables, properted erformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significantly different from the best performanc elformance that is not significant $\frac{13}{35.4}$ $\frac{1000}{31.4}$ $\frac{1000}{31.400}$ $\frac{1000}{31.400$	1056	e (; T	sia-:	94. 12.29	51 104 25(24 11 12 12 12 12 12 12 12 12 12 12 12 12	385.39	51. 42. 1041 214	8640	47. 39. 1287 6822
4000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 43000 44000 440000 44000 <td< td=""><td>1057</td><td>es, ted. anc</td><td>(50 A</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>	1057	es, ted. anc	(50 A							
689.01 69.01 69.01 able 9: Testing results on Asia problems with up to 30000 variation best performance obtained by each method on 10 problems is repertoned. 99.01 ad performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significantly different from the best performance that is not significant size size size size size size size size	1058	able	00						_	
68001 910 910 910 82001 910 910 910 910 92011 910 910 910 910 910 92011 910 910 910 910 910 910 92011 910	1059	rep	94	7.1 1.0 1.0	3.4 3.8 390 11.6	0.4 0.9 32.0 6.9	$\frac{1.5}{63.0}$	3.4 3.8 90.0	/ / / 00.0	7.9 9.8 08.0
682.01 2.001 5.001 able 9: Testing results on Asia problems with up to 30000 2.001 5.001 cf performance obtained by each method on 10 problems 5.001 5.001 cf performance obtained by each method on 10 problems 5.001 5.001 cf performance obtained by each method on 10 problems 5.001 5.001 5.001 old. Problem 5.4 9.41 9.47 5.001 PisLE FI Adjacent 5.4 9.41 9.47 9.47 PisLE FI Adjacent 7.2.6 3.2.40 4.00 5.901 10.6 PisLE FI Adjacent 7.7.7 7.2.6 3.2.4 9.4.1 9.4.7 PisLE FI Adjacent 7.7.7 7.2.6 3.7.1 10.6 PisLE FI Adjacent 7.7.7 3.2.4 9.4.1 0.66 PisLE FI Adjacent 7.7.3 5.7.3 5.2.8 5.0.3 6.1 PisLE FI Adjacent 7.2.9 5.7.7.3 5.7.7.3 5.7.7.3 5.7.7.3	1060	0 v i is	Asia 1000	9523	2.4 % 1	6 8 8 1	× 4 % ×	x 4 % []	864	4 8 2 4
68 2.001 68 2.001 68 2.001 68 2.001 6901 9.01 6901 9.01 6901 9.01 6901 9.01 6901 9.01 6901 9.001 601 9.001 601 9.001 601 9.001 601 9.001 601 9.001 601 7.37 73.0 9.001 73.1 9.10 73.2 9.11 73.3 9.001 73.1 9.100 73.2 9.100 73.3 9.10 73.1 9.100 73.2 9.11 73.3 9.11 73.4 9.11 73.5 9.12 73.6 9.1 73.7 9.1 73.8 9.1 73.9 9.1 73.0 9.1 73.1 9.1 73.1	1061)00 sms	1.3							
6400 2.200 6.200 82.01 2.200 6.200 82.01 2.200 6.200 82.01 2.200 6.200 82.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200 92.01 2.200 2.200	1062	o 30 bble le b	300	_	_	0	0	0	07	o 0.
6401 92.01 62.01 82.01 92.01 92.01 82.01 92.01 92.01 82.01 92.01 92.01 82.01 92.01 92.01 92.01 92.01 92.01 92.01 92.01 93.01 92.01 92.01 93.01 92.01 92.01 93.01 92.01 92.01 93.01 92.01 92.01 93.01 92.01 92.01 93.01 92.01 92.01 93.01 92.02 92.01 93.01 92.03 92.01 93.01 92.04 93.01 10.00 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01 92.05 92.01 93.01<	1063	o tc prc n th	a-30 0, 3:	73.0 73.0 10.6	61 50.8 52.8	77.2 37.5 697. 10.2	85.2 46.4 4.9	61.0 50.8 439. 34.7	56.8 27.2 321.	48.0 39.8 611. 452
6.001 82.01 92.01 abble 9: Testing results on Asia problems with ad performance obtained by each method on erformance that is not significantly different fi old. 92.01 0.01 P/SLE FI Adjacent 73.7 P/SLE FI Adjacent 73.7 72.6 P/SLE FI Adjacent 77.7 67.3 P/SLE(D) FI Adjacent 77.7 72.6 P/SLE FI Adjacent 73.7 72.6 P/SLE(R) FI Adjacent 77.7 72.6 P/SLE(R) FI Adjacent 73.7 72.6 P/SLE(R) FI Adjacent 73.7 72.6 P/SLE(R) FI Adjacent 77.7 72.6 P/SLE(R) FI Adjacent 77.7 72.6	1064	ton ron	Asi (300	2,1,4		Б. (····	4	38.0	1.12 21
6400 8200 8201 9201 8201 9201 8201 9201 8201 9201 8201 9201 9201 8901 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9201 9210 9201 9210 9201 9210 9201 9210 9201 9210 9201 9210 921 9200 9201 9210 921 92340 921 92340 921 92340 921	1065	vitł on 1 fi								
6.001 8.001 9.001 abble 9: Testing results on Asia problem 6.001 abble 9: Testing results on Asia 0.001 plot 73.7 3.4 P/SLE F1 Adjacent 5.4 9.1 P/SLE(D) F1 Adjacent 5.4 9.1 P/SLE(D) Runtime 5.4 9.1 P/SLE(R) Runtime 5.4 9.1 P/SLE(R) Runtime 5.4 9.1 P/SLE(R) Runtime 5.4 9.1 P/SLE(R) Runtime 5.4 9.2 P/SLE(R) Runtime 5.4 9.2 <t< td=""><td>1066</td><td>ns v od o rer</td><td>2200</td><td></td><td>82133</td><td>2 4 ^c 0.</td><td>0.0</td><td>ω 0 ^r. 4</td><td>2 6 2 . 0 2 .</td><td>0.0</td></t<>	1066	ns v od o rer	2200		82133	2 4 ^c 0.	0.0	ω 0 ^r . 4	2 6 2 . 0 2 .	0.0
6201able 9: Testing results on Asia probad performance that is not significantly dold.ci performance that is not significantly dold.pold. <tr< td=""><td>1067</td><td>llen Xthc Tiffe</td><td>sia-2</td><td>94. 72. 8.1</td><td>67. 55. 230</td><td>79. 40. 7.1</td><td>86. 86. 3.4 3.4</td><td>67. 55. 11.</td><td>62. 29. 3457 9984</td><td>48. 40. 7762</td></tr<>	1067	llen Xthc Tiffe	sia-2	94. 72. 8.1	67. 55. 230	79. 40. 7.1	86. 86. 3.4 3.4	67. 55. 11.	62. 29. 3457 9984	48. 40. 7762
	1068	ob] me y d	(20							
	1069	t pr tch	00							
able 9: Testing results on A able 9: Testing results on A able 9: Testing results on A erformance that is not signification old. P/SLE P/State P/SLE	1070	vsia ' ea fica	100	4 5 0 4	5.9 29	8.0 8.0	2.0	7.7 9.9 4.	3.1 5.6 74.0 2.5).5 .7 2.0
able 9: Testing results o able 9: Testing results o cd performance obtained erformance that is not sig old. P/SLE P/SLE </td <td>1071</td> <td>n A by gnij</td> <td>Asia- 000,</td> <td>86520</td> <td>2323</td> <td>50 2 84 86 50 2</td> <td>8 8 8 6</td> <td>595 w</td> <td>52 E</td> <td>50 231 231 220</td>	1071	n A by gnij	Asia- 000,	86520	2323	50 2 84 86 50 2	8 8 8 6	595 w	52 E	50 231 231 220
6201 able 9: Testing result able 9: Testing result cerformance that is nol old. P/SLE SHD Runtime P/SLE Runtime P/GES SHD Runtime FI Adjace	1072	s o ned sig	5							
able 9: Testing res able 9: Testing res ad performance that is old. PYSLE PATOWH PYSLE PYSLE PYSLE PYSLE PATOWH PYSLE PYSLE PATOWH PATOWH	1073	sult tair not		ent lead	ent lead	ent lead	ent lead	ent lead	ent lead	ent ead
able 9: Testing 2001 able 9: Testing 2001 old. P/SLE P/SLE FI Arr P/GES FI Arr P/GES FI Arr P/GES FI Arr P/GES FI Arr P/SLE FI Arr	1074	obi		iowh HD Hi	iowh HD h	djac owh HD	djac Owh HD	iowh THD time	djac owh time	HD him
able 9: Testi 9201 able 9: Testi 9201 erformance t 104 P/SLE 1	1075	ng Ice	티미	Arr Sl Rur	Arr SI Rur	Arr SI Run	Arr Sl Rur	Arr Arr SJ Rur	Art SJ Rur	Ari SI Rur
P/SLE(D) P/SLE(1076	esti nan ce t	oblei -,-	1 ¹¹ E	1 ¹¹ E	- ^н Е	^L E	Ξ	E	1 ¹ E
1028 able 9 berf P/SLE0 P/SLE0 erform	1077	orn ianc	-V-	(*)	<u> </u>	R	ible	s	sle	
	1078	e 9 verf vrm		/SLE	LE(LE(3-Sté	fGE	-Stat	GES
	1079	abl erfc old		d'	P/S	P/S	P/PC	P/	PC	- ÷