
A Conditional Independence Test in the Presence of Discretization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Testing conditional independence (CI) has many important applications, such as
2 Bayesian network learning and causal discovery. Although several approaches have
3 been developed for learning CI structures for observed variables, those existing
4 methods generally fail to work when the variables of interest can not be directly
5 observed and only discretized values of those variables are available. For example,
6 if X_1 , \tilde{X}_2 and X_3 are the observed variables, where \tilde{X}_2 is a discretization of the
7 latent variable X_2 , applying the existing methods to the observations of X_1 , \tilde{X}_2
8 and X_3 would lead to a false conclusion about the underlying CI of variables
9 X_1 , X_2 and X_3 . Motivated by this, we propose a CI test specifically designed to
10 accommodate the presence of discretization. To achieve this, a bridge equation
11 and nodewise regression are used to recover the precision coefficients reflecting
12 the conditional dependence of the latent continuous variables under the nonpara-
13 normal model. An appropriate test statistic has been proposed, and its asymptotic
14 distribution under the null hypothesis of CI has been derived. Theoretical analysis,
15 along with empirical validation on various datasets, rigorously demonstrates the
16 effectiveness of our testing methods.

17 1 Introduction

18 Independence and conditional independence (CI) are fundamental concepts in statistics. They are
19 leveraged for exploring queries in statistical inference, such as sufficiency, parameter identification,
20 adequacy, and ancillarity [9]. They also play a central role in emerging areas such as causal discovery
21 [18], graphical model learning, and feature selection [36]. Tests for CI have attracted increasing
22 attention from both theoretical and application sides.

23 Formally, the problem is to test the CI of two variables X_{j_1} and X_{j_2} given a random vector (a set
24 of other variables) \mathbf{Z} . In statistical notation, the null hypothesis is written as $H_0 : X_{j_1} \perp X_{j_2} \mid \mathbf{Z}$,
25 where \perp denotes “independent from.” The alternative hypothesis is written as $H_1 : X_{j_1} \not\perp X_{j_2} \mid \mathbf{Z}$,
26 where $\not\perp$ denotes “dependent with.” The null hypothesis implies that once \mathbf{Z} is known, the values of
27 X_{j_1} provide no additional information about X_{j_2} , and vice versa. Different tests have been designed
28 to handle different scenarios, including Gaussian variables with linear dependence [37, 25, 22, 26]
29 and non-linear dependence [16, 38, 31, 27, 1] (*For detailed related work, please refer to App. D*).

30 Given observations of X_{j_1} , X_{j_2} , and \mathbf{Z} , the CI can be effectively tested with existing methods.
31 However, in many scenarios, accurately measuring continuous variables of interest is challenging
32 due to limitations in data collection. Sometimes the data obtained are approximations represented as
33 discretized values. For example, in finance, variables such as asset values cannot be measured and are
34 binned into ranges for assessing investment risks (e.g., sell, hold, and strong buy) [7, 8]. Similarly,
35 in mental health, anxiety levels are often assessed using scales like the GAD-7, which categorizes

36 responses into levels such as mild, moderate, or severe [23, 17]. In the entertainment industry, the
 37 quality of movies is typically summarized through viewer ratings [29, 10].

38 When discretization is present, existing CI tests
 39 can fail to determine the CI of underlying con-
 40 tinuous variables. This issue arises because ex-
 41 isting CI tests treat discretized observations as
 42 observations of continuous variables, leading
 43 to incorrect conclusions about their CI relation-
 44 ships. More precisely, the problem lies in the
 45 discretization process, which introduces new dis-
 46 crete variables. Consequently, *although the in-
 47 tent is to test the CI of the underlying continuous
 48 variables, what is actually being tested is the CI
 49 involving a mix of both continuous and newly introduced discrete variables.* In general, this CI
 50 relationship is inconsistent with the one among the underlying continuous variables.

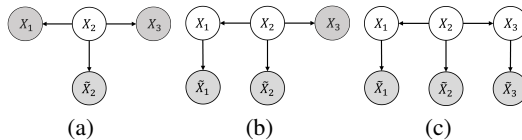


Figure 1: We illustrate different data generative processes with causal graphical models. The discretization process introduces new discrete variables which are denoted with a tilde (\sim).

51 As illustrated in Fig. 1, we show different data-generative processes using causal graphical models
 52 [24] in the presence of discretization. A gray node indicates an observable variable, while a white
 53 node indicates a latent variable. Variables denoted by X_j (without a tilde \sim) represent continuous
 54 variables, which may not be observed; while variables denoted by \tilde{X}_j represent observed discretized
 55 variables derived from X_j due to discretization. In Fig. 1(a), X_2 is latent, and only its discrete
 56 counterpart \tilde{X}_2 is observed. In this case, rather than observing X_1, X_2 , and X_3 , we only observe
 57 X_1, \tilde{X}_2 , and X_3 . Existing CI methods use these observations to test *whether* $X_1 \perp X_3 \mid \{X_2\}$, but
 58 what is actually being tested is *whether* $X_1 \perp X_3 \mid \{\tilde{X}_2\}$. In fact, according to the *causal Markov
 59 condition* [30], it can be inferred from Fig. 1(a) that $X_1 \perp X_3 \mid \{X_2\}$ and $X_1 \not\perp X_3 \mid \{\tilde{X}_2\}$.
 60 This mismatch leads to existing CI methods, that employ observations to check the CI relationships
 61 between X_1 and X_3 given X_2 , to reach incorrect conclusions. Due to the same reason, checking the
 62 CI also fails in Fig 1(b) and Fig 1(c).

63 In this paper, we design a CI test specifically for handling the presence of discretization. An appropri-
 64 ate test statistic for the CI of latent continuous variables, based solely on discretized observations, is
 65 derived. The key is to build connections between the discretized observations and the parameters
 66 needed for testing the CI of the latent continuous variables. To achieve this, we first develop bridge
 67 equations that allow us to estimate the covariance of the underlying continuous variables with dis-
 68 cretized observations. Then, we leverage a *node-wise regression* [5] to derive appropriate test statistics
 69 for CI relationships from the estimated covariance. By assuming that the continuous variables follow
 70 a Gaussian distribution, we can derive the asymptotic distributions of the test statistics under the null
 71 hypothesis of CI. The major contributions of our paper include that

- 72 • We develop a CI test for ensuring accurate analysis in scenarios where data has been discretized,
 73 which are common due to limitations in data collection or measurement techniques, such as in
 74 financial analysis and healthcare.
- 75 • Our CI test can handle various scenarios including 1). Both variables X_{j_1} and X_{j_2} are discretized
 76 2). Both variables X_{j_1} and X_{j_2} are continuous. 3). One of the variables X_{j_1} or X_{j_2} is discretized.
- 77 • We compare our test with the existing methods on both synthetic and real-world datasets, confirm-
 78 ing that our method can effectively estimate the CI of the underlying continuous variables and
 79 outperform the existing tests applied on the discretized observations.

80 2 DCT: A CI Test in the Presence of Discretization

81 **Problem Setting** Consider a set of independent and identically distributed (i.i.d.) p -dimensional
 82 random vectors, denoted as $\tilde{\mathbf{X}} = (X_1, X_2, \dots, \tilde{X}_j, \dots, \tilde{X}_p)^T$. In this set, some variables, indicated
 83 by a tilde (\sim), such as \tilde{X}_j , follow a discrete distribution. For each such variable, there exists a
 84 corresponding latent Gaussian random variable X_j . The transformation from X_j to \tilde{X}_j is governed
 85 by an unknown monotone nonlinear function g_j . This function, $g_j : \mathcal{X} \rightarrow \tilde{\mathcal{X}}$, maps the continuous
 86 domain of X_j onto the discrete domain of $\tilde{\mathcal{X}}$, such that $\tilde{X}_j = g_j(X_j)$ for each observation. Given n
 87 observations $\{\tilde{\mathbf{x}}^1, \tilde{\mathbf{x}}^2, \dots, \tilde{\mathbf{x}}^n\}$ randomly sampled from $\tilde{\mathbf{X}}$, specifically, for each variable X_j , there

88 exists a constant vector $\mathbf{d} = (d_1, \dots, d_M)$ characterized by strictly increasing elements such that

$$\tilde{x}_j^i = \begin{cases} 1 & 0 < g_j(x_j^i) < d_1 \\ m & d_{m-1} < g_j(x_j^i) < d_m \\ M & g_j(x_j^i) > d_m \end{cases} \quad (1)$$

89 This model is also known as the nonparanormal model [20]. The cardinality of the domain after
 90 discretization is at least 2 and smaller than infinity. Our goal is to assess both conditional and
 91 unconditional independence among the variables of the vector $\mathbf{X} = (X_1, X_2, \dots, X_j, \dots, X_p)^T$.
 92 In our model, we assume $\mathbf{X} \sim N(0, \Sigma)$, Σ only contain 1 among its diagonal, i.e., $\sigma_{jj} = 1$ for all
 93 $j \in [1, \dots, p]$. One should note this assumption is *without loss of generality*. We provide a detailed
 94 discussion of our assumption in App. A.8.

95 **Preliminary Framework of DCT** To develop an independence test, one needs to design a test
 96 statistic that can reflect the dependence relation and be calculated from observations. Next, it is
 97 essential to derive the underlying distribution of this statistic under the null hypothesis that the tested
 98 variables are conditionally (or unconditionally) independent. By calculating the value of the test
 99 statistic from observations and determining if this statistic is likely to be sampled from the derived
 100 distribution (i.e., calculating the *p-value* and comparing it with the significance level α), we can
 101 decide if the null hypothesis should be rejected.

102 Our objective is to deduce the independence and CI relationships within the original multivariate
 103 Gaussian model, based on its discretized observations. In the context of a multivariate Gaussian
 104 model, this challenge is directly equivalent to constructing statistical inferences for its covariance
 105 matrix $\Sigma = (\sigma_{j_1, j_2})$ and its precision matrix $\Omega = (\omega_{j, k}) = \Sigma^{-1}$ [3]. The covariance matrix Σ
 106 captures the pairwise covariances between variables, while the precision matrix Ω (also known as the
 107 concentration matrix) provides information about the CI between variables. Specifically, the entry
 108 $\omega_{j, k}$ in the precision matrix is related to the partial correlation coefficient between variables X_j and
 109 X_k , which can be used to test whether these variables are conditionally independent given some other
 110 variables. Technically, we are interested in two things: (1) the calculation of the covariance $\hat{\sigma}_{j_1, j_2}$
 111 and the precision coefficient (or the partial correlation coefficient) $\hat{\omega}_{j, k}$, serving as the estimation
 112 of σ_{j_1, j_2} and $\omega_{j, k}$ respectively (in this paper, a variable with a hat indicates its estimation); and
 113 (2) the derivation of the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ and $\hat{\omega}_{j, k} - \omega_{j, k}$ under the null hypothesis of
 114 independence and CI.

115 In the subsequent section, 1). we first introduce *bridge equations* to address the estimation challenge
 116 of the covariance σ_{j_1, j_2} ; 2). we proceed to derive the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$, demonstrating it
 117 is *asymptotically normal*; 3). utilizing *nodewise regression*, we establish the relationship between
 118 the covariance matrix Σ and the precision matrix Ω , where the regression parameter $\beta_{j, k}$ acts as an
 119 effective surrogate for $\omega_{j, k}$. Leveraging the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$, we further illustrate that
 120 $\hat{\beta}_{j, k} - \beta_{j, k}$ is also *asymptotically normal*.

121 2.1 Design Bridge Equation for Test Statistics

122 **Estimating Covariance with Bridge Equations** The bridge equation establishes a connection
 123 between the underlying covariance σ_{j_1, j_2} of two continuous variables X_{j_1} and X_{j_2} with the ob-
 124 servations. When in the presence of discretization, the discrete transformations make the sample
 125 covariance matrix based on $\tilde{\mathbf{X}}$ inconsistent with the covariance matrix of \mathbf{X} . To obtain the estimation
 126 $\hat{\sigma}_{j_1, j_2}$ of σ_{j_1, j_2} , the bridge equation is leveraged. In general, its form is as follows.

$$\hat{\tau}_{j_1, j_2} = T(\sigma_{j_1, j_2}; \hat{\Lambda}), \quad (2)$$

127 where σ_{j_1, j_2} is the covariance needed to be estimated, $\hat{\tau}_{j_1, j_2}$ is a statistic that can also be estimated
 128 from observations, and $\hat{\Lambda}$ is a set of additional parameters required by the function $T(\cdot)$. The specific
 129 form of the function $T(\cdot)$ will be derived later. Both $\hat{\tau}_{j_1, j_2}$ and $\hat{\Lambda}$ should be able to be calculated
 130 purely relying on observations. Then, given the calculated $\hat{\tau}_{j_1, j_2}$ and $\hat{\Lambda}$, $\hat{\sigma}_{j_1, j_2}$ can be obtained by
 131 solving the bridge equation $\hat{\tau}_{j_1, j_2} = T(\sigma_{j_1, j_2}; \hat{\Lambda})$. As a result, the covariance matrix Σ of \mathbf{X} can be
 132 estimated, which contains information about both unconditional independence and CI (which can be
 133 derived from its inverse).

134 To estimate the covariance of a latent multivariate Gaussian distribution, we need to design appropriate
 135 $\hat{\tau}_{j_1, j_2}$, $\hat{\Lambda}$, and $T(\cdot)$. Notably, bridge equations have to be designed to handle all three possible cases:

136 C1. both observed variables are discretized; C2. one variable is continuous while the other is
 137 discretized; and C3. both variables remain continuous. We will show that cases C1 and C2 can be
 138 merged into a single form of bridge equation with different parameters and a binarization operation
 139 applied to the observations. Our bridge equations are presented in Def. 2.2, Def. 2.3, and Def. 2.4.

140 **Bridge Equations for Discretized and Mixed Pairs** Let us first address the challenging cases
 141 where both observed variables are discretized or where one variable is continuous while the other
 142 is discretized. In general, different bridge equations would need to be designed to handle each case
 143 individually. *However, in our analysis, we provide a unified bridge equation that is applicable to both*
 144 *cases.* This is achieved by binarizing the observed variables, thereby unifying both cases into a binary
 145 case. As some information may be lost in the binarization process, this unification may require more
 146 examples compared to using tailored bridge functions for each specific case. Developing specific
 147 bridge equations for each case to improve sample efficiency is left in future work.

148 Intuitively, for the original continuous variable X_j , binarization separates it into two parts based on
 149 a boundary h_j : the part for X_j larger than h_j and the part for X_j smaller than h_j . In this case, we can
 150 estimate the boundary by calculating the proportion of X_j that exceeds the boundary. In the scenario
 151 of two variables where the threshold h_{j_1} and h_{j_2} divide the space into four regions, the proportions of
 152 these areas are influenced by the covariance σ_{j_1, j_2} , which connects the relation between the binarized
 153 variables with the latent covariance. This approach allows us to initially estimate the threshold h_{j_1} ,
 154 h_{j_2} of a pair of variables, followed by estimating the covariance σ_{j_1, j_2} .

155 Let $\mathbb{P}_n Z$ denote the average of a random variable Z given n i.i.d. observation of Z and $E[Z]$ as the
 156 true mean of Z , \mathbb{P} as the probability and \hat{P} as the empirical probability. We then define the boundary
 157 h_j as follows: for any single discretized variable \tilde{X}_j , there exists a constant c_j such that:

$$\mathbb{1}\{\tilde{x}_j^i > E[\tilde{X}_j]\} = \mathbb{1}\{g_j(x_j^i) > c_j\} = \mathbb{1}\{x_j^i > h_j\},$$

158 where $h_j = g_j^{-1}(c_j)$. Specifically, h_j is the boundary in the original continuous domain to determine
 159 if the discretized observation \tilde{X}_k is larger than its mean. When the continuous variable X_j follows
 160 a normal distribution, there is a relation $\mathbb{P}(\tilde{X}_j > E[\tilde{X}_j]) = 1 - \Phi(h_j)$, where Φ is the cumulative
 161 distribution function (cdf) of a standard normal distribution. We then provide the following definition:

162 **Definition 2.1.** The estimated boundary can be expressed as $\hat{h}_j = \Phi^{-1}(1 - \hat{\tau}_j)$, where $\hat{\tau}_j =$
 163 $\sum_{i=1}^n \mathbb{1}\{\tilde{x}_j^i > \mathbb{P}_n \tilde{X}_j\} / n$, serving as the estimation of $\mathbb{P}(\tilde{X}_j > E[\tilde{X}_j])$.

164 Let $\bar{\Phi}(z_1, z_2; \rho) = \mathbb{P}(Z_1 > z_1, Z_2 > z_2)$, where $(Z_1, Z_2)^T$ follows a bivariate normal distribution
 165 with mean zero, variance one and covariance ρ . We define

$$\tau_{j_1, j_2} = \mathbb{P}(\tilde{x}_{j_1}^i > E[\tilde{X}_{j_1}], \tilde{x}_{j_2}^i > E[\tilde{X}_{j_2}]) = \bar{\Phi}(h_{j_1}, h_{j_2}; \sigma_{j_1, j_2}). \quad (3)$$

166 That is, the proportion of discretized variables larger than their mean can be expressed as a function
 167 of underlying covariance. This equation serves as the key of estimating latent covariance based on the
 168 discretized observations. Specifically, we can substitute those true parameters with their estimation
 169 and construct the bridge equation to get the estimated covariance:

170 **Definition 2.2** (Bridge Equation for A Discretized-Variable Pair). For discretized variables \tilde{X}_{j_1} and
 171 \tilde{X}_{j_2} , the bridge equation is defined as:

$$\hat{\tau}_{j_1, j_2} = \hat{P}(\tilde{X}_{j_1} > \mathbb{P}_n \tilde{X}_{j_1}, \tilde{X}_{j_2} > \mathbb{P}_n \tilde{X}_{j_2}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\tilde{x}_{j_1}^i > \mathbb{P}_n \tilde{X}_{j_1}, \tilde{x}_{j_2}^i > \mathbb{P}_n \tilde{X}_{j_2}\} = T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\}),$$

$$\text{and the function } T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\}) := \bar{\Phi}(\hat{h}_{j_1}, \hat{h}_{j_2}; \sigma) = \int_{x_1 > \hat{h}_{j_1}} \int_{x_2 > \hat{h}_{j_2}} \phi(x_{j_1}, x_{j_2}; \sigma) dx_{j_1} dx_{j_2},$$

172 where ϕ is the probability density function of a bivariate normal distribution, $\hat{h}_{j_1}, \hat{h}_{j_2}$ can be simply
 173 calculated using Def. 2.1.

174 Following the same intuition, we can directly apply the same bridge equation to estimate the co-
 175 variance of mixed pairs. The only difference is there is no need to estimate the boundary \hat{h}_j for the
 176 continuous variable. Instead, we can incorporate its true mean of zero into the equation.

177 **Definition 2.3** (Bridge Equation for A Continuous-Discretized-Variable Pair). For one continuous
 178 variable X_{j_1} and one discretized variable \tilde{X}_{j_2} , the bridge function is defined as follows:

$$\hat{\tau}_{j_1, j_2} = \hat{P}(X_{j_1} > 0, \tilde{X}_{j_2} > \mathbb{P}_n \tilde{X}_{j_2}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_{j_1}^i > 0, \tilde{x}_{j_2}^i > \mathbb{P}_n \tilde{X}_{j_2}\}} = T(\sigma_{j_1, j_2}; \{0, \hat{h}_{j_2}\}),$$

179 and the function $T(\cdot)$ has the same form of Def. 2.2.

180 **A Bridge Equation for A Continuous-Variable Pair** When there is no discretized transformation,
 181 the sample covariance of X_{j_1} and X_{j_2} provides a consistent estimation. In this context, the function
 182 T acts merely as an identity mapping.

183 **Definition 2.4** (A Bridge Equation for A Continuous-Variable Pair). For two continuous variables
 184 X_{j_1} and X_{j_2} , the bridge equation is defined as:

$$\hat{\tau}_{j_1, j_2} := \hat{\sigma}_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i - \frac{1}{n} \sum_{i=1}^n x_{j_1}^i \frac{1}{n} \sum_{i=1}^n x_{j_2}^i = T(\sigma_{j_1, j_2}; \emptyset).$$

185 For two continuous variables X_{j_1} and X_{j_2} , the analytic solution of the estimated covariance can be
 186 simply obtained using Def. 2.4.

187 **Calculation of Estimated Covariance** For the continuous case, the analytic solution of $\hat{\sigma}_{j_1, j_2}$
 188 can be simply obtained using Def. 2.4. For the cases involving the discretized variable as proposed
 189 in Def. 2.2 and Def. 2.3, we can rely on the property that variance Σ only contains 1 among the
 190 diagonal, which implies the covariance σ_{j_1, j_2} should vary from -1 to 1 . Thus, we can calculate the
 191 estimated covariance by solving the objective

$$\min_{\sigma_{j_1, j_2}} \|\hat{\tau}_{j_1, j_2} - T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\})\|^2 \quad s.t. \quad -1 < \sigma_{j_1, j_2} < 1. \quad (4)$$

192 The $\hat{\tau}_{j_1, j_2}$ is a one-to-one mapping with calculated $\hat{\sigma}_{j_1, j_2}$, \hat{h}_{j_1} and \hat{h}_{j_2} , which is proved in App. A.2

193 2.2 Unconditional Independence Test

194 The estimation of covariance $\hat{\sigma}_{j_1, j_2}$ can be effectively solved using the designed bridge equation.
 195 Now, we focus on deriving the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$. These results is used as an unconditional
 196 independence test in the presence of the discretization. Moreover, Thm. 2.5, Lem. 2.6, Lem. 2.7
 197 and Lem. 2.8 will be leveraged in the derivation process of the CI test in Section 2.3. The detailed
 198 derivation steps for both unconditional test and CI test are relatively intricate, therefore, we will
 199 provide a general intuition. For a complete derivation, please refer to the App. A.3.

200 Assume we are interested in the true parameter θ_0 . We denote $\hat{\theta}$ as its estimation which is close to θ_0 ,
 201 and $f(\theta)$ is a continuous function. By leveraging Taylor expansion, we have

$$f(\hat{\theta}) = f(\theta_0) + f'(\theta_0)(\hat{\theta} - \theta_0), \quad (5)$$

202 which directly constructs the relationship between the estimated parameter with the true one. Re-
 203 arrange the term, we get $\hat{\theta} - \theta_0 = (f(\hat{\theta}) - f(\theta_0))/f'(\theta_0)$. If the denominator is a constant and the
 204 numerator can be expressed as a sum of i.i.d samples, we can see $\hat{\theta} - \theta_0$ will be asymptotically
 205 normal according to the central limit theorem [35].

206 Let $\psi_{\hat{\theta}} = [f_{\hat{\theta}}^1(\cdot), f_{\hat{\theta}}^2(\cdot), f_{\hat{\theta}}^3(\cdot)]^T$ contains a group of functions parameterized by $\hat{\theta}$ (For discretized
 207 pairs, $\hat{\theta} = (\hat{\sigma}_{j_1, j_2}, \hat{h}_{j_1}, \hat{h}_{j_2})$). Define $\mathbb{P}_n \psi_{\hat{\theta}}$ as sample mean of these functions evaluated at n sample
 208 points. Similarly, $\mathbb{P}_n \psi_{\hat{\theta}} \psi_{\hat{\theta}}^T$ is defined as sample mean of the outer product $\psi_{\hat{\theta}} \psi_{\hat{\theta}}^T$. The notation
 209 $P\psi_{\hat{\theta}} := E\mathbb{P}_n \psi_{\hat{\theta}}$ denotes the expectations of the functions in $\psi_{\hat{\theta}}$. Furthermore, let $\psi_{\hat{\theta}}'$ denote the
 210 derivative of the functions contained in $\psi_{\hat{\theta}}$. We now provide the main result of derived distribution
 211 $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ under the hull hypothesis that test pairs are independent.

212 **Theorem 2.5** (Independence Test). *In our settings, under the null hypothesis that two observed*
 213 *variables indexed with j_1 and j_2 are statistically independent under our framework, i.e., $\sigma_{j_1, j_2} = 0$,*
 214 *the independence can be tested using the statistic*

$$\hat{\sigma}_{j_1, j_2} = T^{-1}(\hat{\tau}_{j_1, j_2}; \hat{\theta}).$$

215 This statistic is approximated to follow a normal distribution, as detailed below:

$$\hat{\sigma}_{j_1, j_2} \stackrel{\text{approx}}{\sim} N \left(0, \frac{1}{n} \left((\mathbb{P}_n \psi'_\theta)^{-1} \mathbb{P}_n \psi_\theta \psi_\theta^T (\mathbb{P}_n \psi'_\theta)^{-1} \right)_{1,1} \right), \quad (6)$$

216 where the specific form of ψ_θ are presented in Lem. 2.6, Lem. 2.7 and Lem. 2.8.

217 We now provide the specific forms of ψ_θ . Since the variables being tested for independence can be
 218 both discretized, only one being discretized, or neither being discretized. This results in different
 219 forms of ψ_θ consequently differs across these scenarios. Let Z_{j_1} and Z_{j_2} be any two random
 220 variables indexed by j_1 and j_2 . Let $\hat{\sigma}_{j_1, j_2}^i = z_{j_1}^i \cdot z_{j_2}^i - \mathbb{P}_n Z_{j_1} \cdot \mathbb{P}_n Z_{j_2}$ denote the sample covariance
 221 based on a i -th pairwise observation of the variables Z_{j_1} and Z_{j_2} . Let $\hat{\tau}_{j_1}^i = \mathbb{1}_{\{z_{j_1}^i > \mathbb{P}_n Z_{j_1}\}}$ and
 222 $\hat{\tau}_{j_2}^i = \mathbb{1}_{\{Z_{j_2}^i > \mathbb{P}_n Z_{j_2}\}}$, each calculated based on i -th observations of the variables Z_{j_1} and Z_{j_2} ,
 223 respectively. Let $\hat{\tau}_{j_1, j_2}^i$ be $\hat{\tau}_{j_1}^i \cdot \hat{\tau}_{j_2}^i$. We further denote $\bar{\Phi}(\cdot) = 1 - \Phi(\cdot)$. The different forms of ψ_θ
 224 that arise in different cases are defined as follows:

225 **Lemma 2.6.** (ψ_θ for A Continuous-Variable Pair). For two continuous variables X_{j_1} and X_{j_2} ,

$$\psi_\theta := \hat{\sigma}_{j_1, j_2}^i - \hat{\sigma}_{j_1, j_2}. \quad (7)$$

226 **Lemma 2.7** (ψ_θ for A Discretized-Variable Pair). For discretized variables \tilde{X}_{j_1} and \tilde{X}_{j_2} ,

$$\psi_\theta := \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\hat{\sigma}_{j_1, j_2}^i; \{\hat{h}_{j_1}, \hat{h}_{j_2}\}) \\ \hat{\tau}_{j_1}^i - \bar{\Phi}(\hat{h}_{j_1}) \\ \hat{\tau}_{j_2}^i - \bar{\Phi}(\hat{h}_{j_2}) \end{pmatrix}. \quad (8)$$

227 **Lemma 2.8** (ψ_θ for A Continuous-Discretized-Variable Pair). For one discretized variable \tilde{X}_{j_2} and
 228 one continuous variable X_{j_1} ,

$$\psi_\theta := \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\hat{\sigma}_{j_1, j_2}^i; \{0, \hat{h}_{j_2}\}) \\ \hat{\tau}_{j_1}^i - \bar{\Phi}(\hat{h}_{j_2}) \end{pmatrix}. \quad (9)$$

229 Derivation of forms of ψ_θ for different cases and their corresponding distribution defined in Eq (6)
 230 can be found in App. A.4, App. A.5, App. A.6. Up to this point, our discussion has been confined to
 231 the case of covariance σ_{j_1, j_2} , the indicator of unconditional independence. In the next section, we
 232 will present the results of our CI test.

233 2.3 Conditional Independence (CI) Test

234 To construct a CI test of our model, we are interested at two things: calculation of the estimated
 235 precision coefficient $\hat{\omega}_{j, k}$ and the derivation of the corresponding distribution $\hat{\omega}_{j, k} - \omega_{j, k}$. In the
 236 following, we first build $\beta_{j, k}$, which is obtained using nodewise regression and show it serves as a
 237 surrogate of testing for $\omega_{j, k} = 0$, we then construct the formulation of $\hat{\beta}_{j, k} - \beta_{j, k}$ as the combination
 238 of formulation of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ and show it will also be asymptotically normal.

239 **Nodewise Regression for CI** To utilize covariance for testing CI, it is necessary to establish a
 240 relationship between the estimated covariance and a metric capable of reflecting CI. To achieve this,
 241 we employ the nodewise regression which effectively builds the connection between covariance
 242 and precision matrix. Suppose we can access observations $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$ from latent continuous
 243 variables $\mathbf{X} = (X_1, \dots, X_p) \sim N(0, \Sigma)$, nodewise regression will do regression on every dimension
 244 with all other dimensions as predictors.

$$x_{j_1}^i = \sum_{j_1 \neq j_2} x_{j_2}^i \beta_j + \epsilon_{j_1}^i. \quad (10)$$

245 It can be shown that there are deterministic relationships between the regression coefficients and the
 246 covariance and precision matrices of \mathbf{X} , as illustrated below and proved in App. A.7.1.

$$\beta_j = \Sigma_{-j-j}^{-1} \Sigma_{-jj} \in \mathbb{R}^{p-1}, \quad \beta_{j, k} = -\frac{\omega_{j, k}}{\omega_{j, j}}, \quad j \neq k, \quad (11)$$

247 where Σ_{-j-j} is the submatrix of Σ without j th column and j th row, and the Σ_{-jj} is the vector of j th
 248 column without j th row. $\beta_{j, k} \in \mathbb{R}$ is the surrogate of $\omega_{j, k}$ to capture the independence relationship of
 249 X_j with X_k conditioning on other variables. We can use Def. 2.2, Def. 2.3 and Def. 2.4 to get the
 250 estimation $\hat{\Sigma}_{-j-j}$ and $\hat{\Sigma}_{-jj}$ and thus get the estimation $\hat{\beta}_j$.

251 **Statistical Inference for $\beta_{j,k}$** Nodewise regression offers a robust solution for the estimation
 252 problem. A pertinent inquiry pertains to the construction of the distribution of $\hat{\beta}_j - \beta_j$. It is crucial
 253 to recognize that the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ is already established. Therefore, if we can
 254 conceptualize $\hat{\beta}_j - \beta_j$ as a linear combination of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$, the problem is directly solved, i.e.,
 255 the $\hat{\beta}_j - \beta_j$ is linear combination of dependent Gaussian variables. The underlying relationship
 256 between these variables is as follows:

$$\hat{\beta}_j - \beta_j = -\hat{\Sigma}_{-j-j}^{-1} \left((\hat{\Sigma}_{-j-j} - \Sigma_{-j-j})\beta_j - (\hat{\Sigma}_{-jj} - \Sigma_{-jj}) \right).$$

257 The derivation is provided in App. A.7.2. For ease of notation, we further express the distribution of
 258 the difference between the estimated covariance and the true covariance as

$$\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n \xi_{j_1, j_2}^i. \quad (12)$$

259 The specific form of ξ_{j_1, j_2}^i is given in App. A.4, A.5, A.6 respectively for different cases. For
 260 notational convenience, we express $\hat{\Sigma}_{-j-j} - \Sigma_{-j-j} = \frac{1}{n} \sum_{i=1}^n \Xi_{-j, -j}^i$ and $\hat{\Sigma}_{-jj} - \Sigma_{-jj} =$
 261 $\frac{1}{n} \sum_{i=1}^n \Xi_{-j, j}^i$, where ξ_{j_1, j_2}^i is the element of the matrix Ξ at the position indexed by (j_1, j_2) . We
 262 now propose the statistic and its asymptotic distribution for the CI test in the following theorem.

263 **Theorem 2.9** (Conditional Independence test). *In our settings, under the null hypothesis that X_j and*
 264 *X_k are conditional statistically independent given a set of variables \mathbf{Z} , i.e., $\beta_{j,k} = 0$, the statistic*

$$\hat{\beta}_{j,k} = (\hat{\Sigma}_{-j-j}^{-1} \hat{\Sigma}_{-jj})_{[k]}, \quad (13)$$

265 *where $[k]$ denotes the element corresponding to the variable X_k in $\hat{\Sigma}_{-j-j}^{-1} \hat{\Sigma}_{-jj}$. The statistic $\hat{\beta}_{j,k}$*
 266 *has the asymptotic distribution:*

$$\hat{\beta}_{j,k} \sim N\left(0, a^{[k]T} \frac{1}{n^2} \sum_{i=1}^n \text{vec}(B_{-j}^i) \text{vec}(B_{-j}^i)^T a^{[k]}\right),$$

267

$$\text{where } B^i = \begin{bmatrix} \Xi_{-j, -j}^i \\ \Xi_{-j, j}^i \end{bmatrix}, \quad a_l^{[k]} = \begin{cases} \left(\hat{\Sigma}_{-j-j}^{-1} \right)_{[k], l}, & \text{for } l \in \{1, \dots, p-1\} \\ \sum_{q=1}^n \left(\hat{\Sigma}_{-j-j}^{-1} \right)_{[k], l} \left(\tilde{\beta}_j \right)_q, & \text{for } l \in \{p, \dots, p^2 - p\} \end{cases}$$

268 *and $\tilde{\beta}_j$ is β_j whose $\beta_{j,k} = 0$.*

269 In practice, we can plug in the estimation of regression parameter $\hat{\beta}_j$ and set $\hat{\beta}_{j,k} = 0$ as the
 270 substitution of $\tilde{\beta}_j$ to calculate the variance and do the CI test. Specifically, we can obtain the $\hat{\beta}_{j,k}$
 271 using Eq. (13) where the estimated covariance terms can be calculated by solving the bridge equation
 272 Eq. 2. Under the null hypothesis that $\beta_{j,k} = 0$ (conditional independence), we can take the calculated
 273 $\hat{\beta}_{j,k}$ into the distribution defined in Thm. 2.9 and obtain the p-value. If the p-value is smaller than the
 274 predefined significance level α (normally set at 0.05), we will infer the tested pairs are conditionally
 275 dependent; otherwise, we do not. The detailed derivation of the Thm. 2.9 can be found in App. A.7.2.

276 3 Experiments

277 We applied the proposed method DCT to synthetic data to evaluate its practical performance and
 278 compare it with Fisher-Z test [14] (for all three data types) and Chi-Square test [15] (for discrete data
 279 only) as baselines. Specifically, we investigated its Type I and Type II error and its application in
 280 causal discovery. The experiments investigating its robustness, performance in denser graphs and
 281 effectiveness in a real-world dataset can be found in App. C.

282 3.1 On the Effect of the Cardinality of Conditioning Set and the Sample Size

283 Our experiment investigates the variations in Type I and Type II error (1 minus power) probabilities
 284 under two conditions. In the first scenario, we focus on the effects of modifying the sample size,
 285 denoted as $n = (100, 500, 1000, 2000)$, while conditioning on a single variable. In the second,
 286 the sample size is held constant at 2000, and we vary the cardinality of the conditioning set, represented

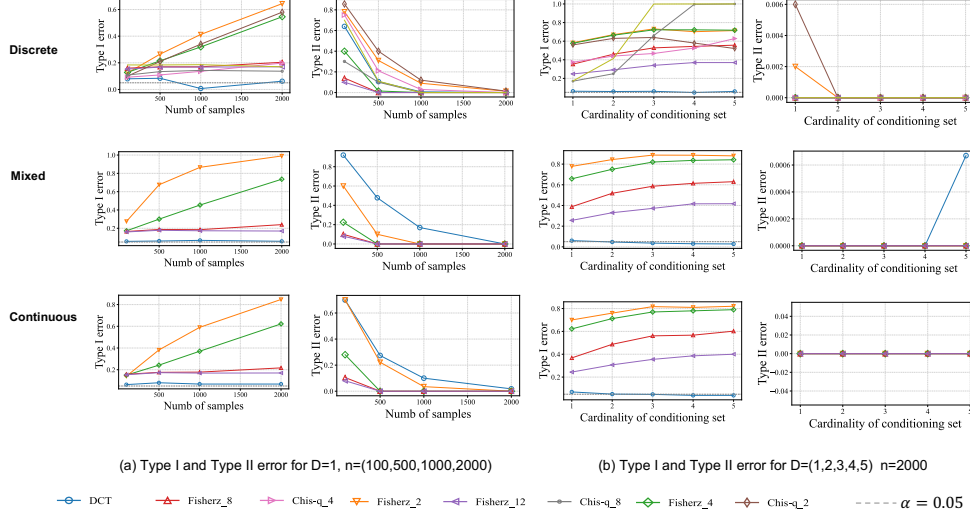


Figure 2: Comparison of results of Type I and Type II error ($1 - \text{power}$) for all three types of tested data (continuous, mixed, discrete) and different number of samples and cardinality of conditioning set. The suffix attached to a test's name denotes the cardinality of discretization; for example, "Fsherz_4" signifies the application of the Fsher-z test to data discretized into four levels. Chi-square test is only applicable for the discrete case.

287 as $D = (1, 2, \dots, 5)$. It is assumed that every variable within this conditioning set is effective, i.e.,
 288 they influence the CI of the tested pairs. We repeat each test 1500 times.

289 We use Y, W to denote the variables being tested and use Z to denote the variables being conditioned
 290 on. The discretized versions of the variables are denoted with a tilde symbol (e.g., \tilde{Z}). For both condi-
 291 tions, we evaluate three distinct types of observations of tested variables: continuous observations
 292 for both variables (Y, W), discrete observations for both variables (\tilde{Y}, \tilde{W}) and a mixed type (\tilde{Y}, W).
 293 The variables in the conditioning set will always be discretized observations (\tilde{Z}).

294 To see how well the derived asymptotic null distribution approximates the true one, we verify if
 295 the probability of Type I error aligns with the significance level α preset in advance. We generate
 296 true continuous multivariate Gaussian data Y, W from Z_i (single $i = 1$ for the first scenario, and
 297 summed over n for the second), structured as $a_i Z_i + E$ and $\sum_{i=1}^n a_i Z_i + E$, where a_i is sampled
 298 from $U(0.5, 1.5)$ and E follows a standard normal distribution, independent of all other variables.
 299 This ensures $Y \perp\!\!\!\perp W | Z$. The data are then discretized into $K = (2, 4, 8, 12)$ levels, with boundaries
 300 randomly set based on the variable range. The first column in Fig. 2 (a) (b) shows the resulting
 301 probability of Type I errors at the significance level $\alpha = 0.05$ compared with other methods.

302 A good test should have as small a probability of Type II error as possible, i.e., a larger power. To
 303 test the power of our DCT, we generate the continuous multivariate Gaussian data Z_i from Y, W ;
 304 constructed as $Z_i = a_i Y + b_i W + E$, where a_i, b_i are sampled from $U(0.5, 1.5)$ and E follows a
 305 standard normal distribution independent with all others, i.e., $Y \perp\!\!\!\perp W | Z$. The same discretization
 306 approach is applied here. The second column in Fig. 2 (a) (b) shows the Type II error with the
 307 changing number of samples and cardinality of the conditioning set compared with other methods.

308 From Fig. 2 (a), we note that the Type I error rates with our derived null distribution are well-
 309 approximated at 0.05 across all three data types in both scenarios. In contrast, other testing methods
 310 show significantly higher Type I error rates, increasing with the number of samples and the size of
 311 the conditioning set. This indicates that such methods are more prone to erroneously concluding
 312 that tested variables are conditionally dependent. Additionally, while alternative tests demonstrate
 313 considerable power with smaller sample sizes, our approach requires a sample size of 2000 to achieve
 314 satisfactory power, particularly in mixed and continuous cases. A possible explanation for this
 315 phenomenon is that our method binarizes discretized data, which may not effectively utilize all
 316 observations. This aspect warrants further investigation in future research. Moreover, our test shows
 317 remarkable stability in response to changes in the number of conditioning sets.

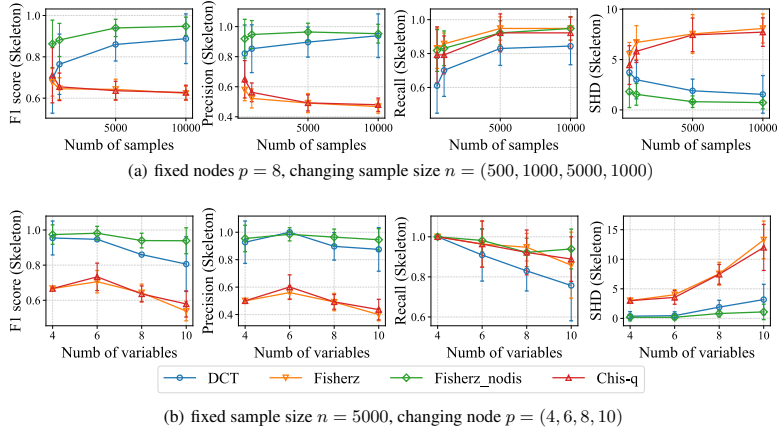


Figure 3: Experiment result of skeleton discovery on synthetic data for changing sample size (a) and changing number of nodes (b). Fisherz_nodis is the Fisher-z test applied to original continuous data. We evaluate F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow).

3.2 Application in Causal Discovery

Causal discovery aims to recover the true causal structure from the data. Constraint-based causal discovery methods like the PC algorithm [30] rely on testing CI from observations to discover causal graphs. However, in the presence of discretization, failures in testing CI leads to false conclusions about causal graphs. To evaluate the efficacy of the DCT, we construct causal graphs utilizing the Bipartite Pairing (BP) model as detailed in [2], with the number of edges being one fewer than the number of nodes. The detailed generation process is provided in App. B due to limited space. Our experiment is divided into two scenarios: (a) fixed data samples $n = 5000$, with changing number of nodes $p = (4, 6, 8, 10)$; (b) fixed number of nodes $p = 8$ and changing sample sizes $n = (500, 1000, 5000, 10000)$.

Comparative analysis is conducted using the PC algorithm integrated with different testing methods. Specifically, we compare DCT against the Fisher-Z test applied to discretized data, the chi-square test, and the Fisher-Z test on original continuous data, the latter serving as a theoretical upper bound for comparison. Since the PC algorithm can only return a completed partially directed acyclic graph (CPDAG), we use the same orientation rules [11] implemented by Causal-DAG [6] to convert a CPDAG into a DAG. We evaluate both the undirected skeleton and the directed graph using criteria such as structural Hamming distance (SHD), F1 score, precision, and recall. For each setting, we run 10 graph instances with different seeds and report the mean and standard deviation of skeleton discovery in Fig. 3, and DAG in Fig. 4 in App B.

According to the result, DCT exhibits performance nearly on par with the theoretical upper bound across metrics such as F1 score, precision, and Structural Hamming Distance (SHD) when the number of variables (p) is small and the sample size (n) is large. Despite a decline in performance as the number of variables increases with a smaller sample size, DCT significantly outperforms both the Fisher-Z test and the Chi-square test. Notably, in almost all settings, the recall of DCT is lower than that of the baseline tests, which is a reasonable outcome *since these tests tend to infer conditional dependencies, thereby retaining all edges given the discretized observations*. For instance, a fully connected graph, would achieve a recall of 1.

4 Conclusion

In this paper, we present a new testing method tailored for scenarios commonly encountered in real-world applications, where variables, though inherently continuous, are only observable in their discretized forms. Our method distinguishes itself from existing CI tests by effectively mitigating the misjudgment introduced by discretization and accurately recovering the CI relationships of latent continuous variables. We substantiate our approach with theoretical results and empirical validation, underscoring the effectiveness of our testing methods.

References

- 352
- 353 [1] Constantin F Aliferis, Alexander Statnikov, Ioannis Tsamardinos, Subramani Mani, and Xenofon D
354 Koutsoukos. Local causal and markov blanket induction for causal discovery and feature selection for
355 classification part i: algorithms and empirical evaluation. *Journal of Machine Learning Research*, 11(1),
356 2010.
- 357 [2] Armen S Asratian, Tristan MJ Denley, and Roland Häggkvist. *Bipartite graphs and their applications*,
358 volume 131. Cambridge university press, 1998.
- 359 [3] Kunihiro Baba, Ritei Shibata, and Masaaki Sibuya. Partial correlation and conditional correlation as
360 measures of conditional independence. *Australian & New Zealand Journal of Statistics*, 46(4):657–664,
361 2004.
- 362 [4] Kunihiro Baba, Ritei Shibata, and Masaaki Sibuya. Partial correlation and conditional correlation as
363 measures of conditional independence. *Australian & New Zealand Journal of Statistics*, 46(4):657–664,
364 2004.
- 365 [5] Laurent Callot, Mehmet Caner, Esra Ulasan, and A. Özlem Önder. A nodewise regression approach to
366 estimating large portfolios, 2019.
- 367 [6] Chandler Squires. *causal.dag: creation, manipulation, and learning of causal models*, 2018.
- 368 [7] Hu Changsheng and Wang Yongfeng. Investor sentiment and assets valuation. *Systems Engineering*
369 *Procedia*, 3:166–171, 2012.
- 370 [8] Aswath Damodaran. *Investment valuation: Tools and techniques for determining the value of any asset*,
371 volume 666. John Wiley & Sons, 2012.
- 372 [9] A Philip Dawid. Conditional independence in statistical theory. *Journal of the Royal Statistical Society*
373 *Series B: Statistical Methodology*, 41(1):1–15, 1979.
- 374 [10] Simon Doods, Toon De Pessemier, and Luc Martens. Movietweetings: a movie rating dataset collected
375 from twitter. In *Workshop on Crowdsourcing and human computation for recommender systems, CrowdRec*
376 *at RecSys*, volume 2013, page 43, 2013.
- 377 [11] Dorit Dor and Michael Tarsi. A simple algorithm to construct a consistent extension of a partially oriented
378 graph. 1992.
- 379 [12] Gary Doran, Krikamol Muandet, Kun Zhang, and Bernhard Schölkopf. A permutation-based kernel
380 conditional independence test. In *UAI*, pages 132–141, 2014.
- 381 [13] Jianqing Fan, Han Liu, Yang Ning, and Hui Zou. High dimensional semiparametric latent graphical model
382 for mixed data. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 79(2):405–421,
383 2017.
- 384 [14] Ronald Aylmer Fisher. On the "Probable Error" of a Coefficient of Correlation Deduced from a Small
385 Sample. *Metron*, 1:3–32, 1921.
- 386 [15] Karl Pearson F.R.S. X. on the criterion that a given system of deviations from the probable in the case of
387 a correlated system of variables is such that it can be reasonably supposed to have arisen from random
388 sampling. *Philosophical Magazine Series 1*, 50:157–175, 2009.
- 389 [16] Kenji Fukumizu, Francis R Bach, and Michael I Jordan. Dimensionality reduction for supervised learning
390 with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5(Jan):73–99, 2004.
- 391 [17] Sverre Urnes Johnson, Pål Gunnar Ulvenes, Tuva Øktedalen, and Asle Hoffart. Psychometric properties
392 of the general anxiety disorder 7-item (gad-7) scale in a heterogeneous psychiatric sample. *Frontiers in*
393 *psychology*, 10:449461, 2019.
- 394 [18] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. Adaptive
395 computation and machine learning. MIT Press, 2009.
- 396 [19] Loka Li, Ignavier Ng, Gongxu Luo, Biwei Huang, Guangyi Chen, Tongliang Liu, Bin Gu, and Kun Zhang.
397 Federated causal discovery from heterogeneous data, 2024.
- 398 [20] Han Liu, John Lafferty, and Larry Wasserman. The nonparanormal: Semiparametric estimation of high
399 dimensional undirected graphs. *Journal of Machine Learning Research*, 10(10), 2009.

- 400 [21] Dimitris Margaritis. Distribution-free learning of bayesian network structure in continuous domains. In
401 AAAI, volume 5, pages 825–830, 2005.
- 402 [22] Karthik Mohan, Mike Chung, Seungyeop Han, Daniela Witten, Su-In Lee, and Maryam Fazel. Structured
403 learning of gaussian graphical models. *Advances in neural information processing systems*, 25, 2012.
- 404 [23] Sarah A Mossman, Marissa J Luft, Heidi K Schroeder, Sara T Varney, David E Fleck, Drew H Barzman,
405 Richard Gilman, Melissa P DelBello, and Jeffrey R Strawn. The generalized anxiety disorder 7-item
406 (gad-7) scale in adolescents with generalized anxiety disorder: signal detection and validation. *Annals of*
407 *clinical psychiatry: official journal of the American Academy of Clinical Psychiatrists*, 29(4):227, 2017.
- 408 [24] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- 409 [25] Christine Peterson, Francesco C Stingo, and Marina Vannucci. Bayesian inference of multiple gaussian
410 graphical models. *Journal of the American Statistical Association*, 110(509):159–174, 2015.
- 411 [26] Zhao Ren, Tingni Sun, Cun-Hui Zhang, and Harrison H Zhou. Asymptotic normality and optimalities in
412 estimation of large gaussian graphical models. 2015.
- 413 [27] Rajat Sen, Ananda Theertha Suresh, Karthikeyan Shanmugam, Alexandros G Dimakis, and Sanjay
414 Shakkottai. Model-powered conditional independence test. *Advances in neural information processing*
415 *systems*, 30, 2017.
- 416 [28] Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara, Takashi
417 Washio, Patrik O. Hoyer, and Kenneth Bollen. Directlingam: A direct method for learning a linear
418 non-gaussian structural equation model, 2011.
- 419 [29] E Isaac Sparling and Shilad Sen. Rating: how difficult is it? In *Proceedings of the fifth ACM conference on*
420 *Recommender systems*, pages 149–156, 2011.
- 421 [30] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT press, 2nd edition, 2000.
- 422 [31] Eric V Strobl, Kun Zhang, and Shyam Visweswaran. Approximate kernel-based conditional independence
423 tests for fast non-parametric causal discovery. *Journal of Causal Inference*, 7(1):20180017, 2019.
- 424 [32] Liangjun Su and Halbert White. A nonparametric hellinger metric test for conditional independence.
425 *Econometric Theory*, 24(4):829–864, 2008.
- 426 [33] A. W. van der Vaart. *M-and Z-Estimators*, page 41–84. Cambridge Series in Statistical and Probabilistic
427 Mathematics. Cambridge University Press, 1998.
- 428 [34] A. W. van der Vaart. *Stochastic Convergence*, page 5–24. Cambridge Series in Statistical and Probabilistic
429 Mathematics. Cambridge University Press, 1998.
- 430 [35] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- 431 [36] Eric P Xing, Michael I Jordan, Richard M Karp, et al. Feature selection for high-dimensional genomic
432 microarray data. In *Icml*, volume 1, pages 601–608. Citeseer, 2001.
- 433 [37] Ming Yuan and Yi Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*,
434 94(1):19–35, 2007.
- 435 [38] Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional indepen-
436 dence test and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- 437 [39] Yishi Zhang, Zigang Zhang, Kaijun Liu, and Gangyi Qian. An improved iamb algorithm for markov
438 blanket discovery. *J. Comput.*, 5(11):1755–1761, 2010.
- 439 [40] Yujia Zheng, Biwei Huang, Wei Chen, Joseph Ramsey, Mingming Gong, Ruichu Cai, Shohei Shimizu,
440 Peter Spirtes, and Kun Zhang. Causal-learn: Causal discovery in python, 2023.

441 *Appendix for*

442

443 **“A Conditional Independence Test in the Presence of Discretization”**

444 Appendix organization:

445

446	A Proof of Things	12
447	A.1 Proof of $\hat{\theta} \xrightarrow{P} \theta_0$	12
448	A.2 Proof of one-to-one mapping between $\hat{\tau}_{j_1, j_2}$ with $\hat{\sigma}_{j_1, j_2}$	13
449	A.3 Proof of Thm. 2.5	13
450	A.4 Derivation of Lem. 2.7	14
451	A.5 Derivation of Lem. 2.8	15
452	A.6 Derivation of Lem. 2.6	16
453	A.7 Proof of Thm. 2.9	16
454	A.7.1 Proof of Relation between Σ, Ω with β	16
455	A.7.2 Detailed derivation of inference for β_j	17
456	A.8 Discussion of assumption of zero mean and identity variance	19
457	B Data Generation and Figure of main experiments: causal discovery	20
458	C Additional experiments	21
459	C.1 Linear non-Gaussian and nonlinear	21
460	C.2 Denser graph	21
461	C.3 multivariate Gaussian with nonzero mean and non-unit variance	21
462	C.4 Real-world dataset	22
463	D Related Work	24
464	E Resource Usage	25
465	F Limiation and Broader Impacts	25
466		
467		

468 **A Proof of Things**

469 **A.1 Proof of $\hat{\theta} \xrightarrow{P} \theta_0$**

470 **Lemma A.1.** *For the estimation $\hat{\theta}$ which is calculated using bridge equation 2.4 2.2 2.3,*
471 *as a zero of Ψ_n defined in Eq. (26),(33), (36) , will converge in probability to $\theta_0 =$*
472 *$(\sigma_{j_1, j_2}, h_{j_1}, h_{j_2}), (\sigma_{j_1, j_2}, h_{j_2}), (\sigma_{j_1, j_2})$ respectively.*

473 *Proof* We first focus on the most challenging one where both variables are discrete. According to
474 the law of large numbers, for the estimated boundary \hat{h}_{j_1} and \hat{h}_{j_2} whose calculations are defined as

475 $\hat{h}_j = \Phi^{-1}(1 - \hat{\tau}_j)$, we should have

$$n \rightarrow \infty, \quad \hat{\tau}_j = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\tilde{x}_j^i > \mathbb{P}_n \tilde{X}_j\}} \xrightarrow{P} \mathbb{P}(\tilde{X}_j > E[\tilde{X}_j]). \quad (14)$$

476 Recall the definition $\mathbb{P}(\tilde{X}_j > E[\tilde{X}_j]) = 1 - \Phi(h_j)$, according to continuous mapping theorem [34],
 477 as long as the function $\Phi^{-1}(1 - \cdot)$ is continuous, we should have $\hat{h}_j \xrightarrow{P} h_j$. And thus $\hat{h}_{j_1} \xrightarrow{P} h_{j_1}$,
 478 $\hat{h}_{j_2} \xrightarrow{P} h_{j_2}$.

479 We have $\hat{\tau}_{j_1, j_2} = \bar{\Phi}(\hat{h}_{j_1}, \hat{h}_{j_2}, \hat{\sigma}_{j_1, j_2})$ and the estimation $\hat{\sigma}_{j_1, j_2}$ can be obtained through solving the
 480 function. Similarly, we also have

$$n \rightarrow \infty, \quad \hat{\tau}_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\tilde{x}_{j_1}^i > \mathbb{P}_n \tilde{X}_{j_1}\}} \mathbb{1}_{\{\tilde{x}_{j_2}^i > \mathbb{P}_n \tilde{X}_{j_2}\}} \xrightarrow{P} \mathbb{P}(\tilde{x}_{j_1}^i > E[\tilde{X}_{j_1}], \tilde{x}_{j_2}^i > E[\tilde{X}_{j_2}]) = \tau_{j_1, j_2}. \quad (15)$$

481 Similarly, according to the continuous mapping theorem, we have $\hat{\sigma}_{j_1, j_2} \xrightarrow{P} \sigma_{j_1, j_2}$. Thus, the
 482 parameter $(\hat{\sigma}_{j_1, j_2}, \hat{h}_{j_1}, \hat{h}_{j_2}) \xrightarrow{P} (\sigma_{j_1, j_2}, h_{j_1}, h_{j_2})$.

483 Apparently, the result above could easily extend to the mixed case where we fix $\hat{h}_1 = h_1 = 0$. Using
 484 the same procedure, we should have $(\hat{\sigma}_{j_1, j_2}, \hat{h}_{j_2}) \xrightarrow{P} (\sigma_{j_1, j_2}, h_{j_2})$.

485 For the continuous case whose estimated variance is calculated as $\hat{\sigma}_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i -$
 486 $\frac{1}{n} \sum_{i=1}^n x_{j_1}^i \frac{1}{n} \sum_{i=1}^n x_{j_2}^i$, according to law of large numbers, we should have

$$n \rightarrow \infty, \quad \hat{\sigma}_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i - \frac{1}{n} \sum_{i=1}^n x_{j_1}^i \frac{1}{n} \sum_{i=1}^n x_{j_2}^i \xrightarrow{P} E(X_{j_1} X_{j_2}) - E(X_{j_1})E(X_{j_2}) = \sigma_{j_1, j_2}. \quad (16)$$

487 A.2 Proof of one-to-one mapping between $\hat{\tau}_{j_1, j_2}$ with $\hat{\sigma}_{j_1, j_2}$

488 **Lemma A.2.** For any fixed \hat{h}_{j_1} and \hat{h}_{j_2} , $T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\}) =$
 489 $\int_{x_1 > \hat{h}_{j_1}} \int_{x_2 > \hat{h}_{j_2}} \phi(x_{j_1}, x_{j_2}; \sigma) dx_{j_1} dx_{j_2}$, is a strictly monotonically increasing function on
 490 $\sigma \in (-1, 1)$.

491 *Proof* To prove the lemma, we just need to show the gradient $\frac{\partial T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\})}{\partial \sigma} > 0$ for $\sigma \in (-1, 1)$.

$$\frac{\partial T(\sigma_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\})}{\partial \sigma} = \frac{1}{2\pi\sqrt{(1-\sigma^2)}} \exp\left(-\frac{(\hat{h}_{j_1}^2 - 2\sigma\hat{h}_{j_1}\hat{h}_{j_2} + \hat{h}_{j_2}^2)}{2(1-\sigma^2)}\right), \quad (17)$$

492 which is obviously positive for $\sigma \in (-1, 1)$. Thus, we have one-to-one mapping between $\hat{\tau}_{j_1, j_2}$ with
 493 the calculated $\hat{\sigma}_{j_1, j_2}$ for fixed \hat{h}_{j_1} and \hat{h}_{j_2} .

494 A.3 Proof of Thm. 2.5

495 In this section, we provide the proof of Thm. 2.5, which utilizes a regular statistical tool: Z-estimator
 496 [33]. Specifically, we are interested in the parameter θ and we have its estimation $\hat{\theta}$. Let $\mathbf{x}_1, \dots, \mathbf{x}_n$
 497 are sampled from some true distribution P , we can construct the function characterized by the
 498 parameter θ related to \mathbf{x} as $\psi_\theta(\mathbf{x})$. As long as we have n observations, we can construct the function
 499 as follows

$$\Psi_n(\theta) = \frac{1}{n} \sum_{i=1}^n \psi_\theta(\mathbf{x}_i) = \mathbb{P}_n \psi_\theta. \quad (18)$$

500 We further specify the form

$$\Psi(\theta) = \int \psi_\theta(\mathbf{x}) d\mathbf{x} = P \psi_\theta. \quad (19)$$

501 Assume the estimator $\hat{\theta}$ is a zero of Ψ_n , i.e., $\Psi_n(\hat{\theta}) = 0$ and will converge in probability to θ_0 , which
 502 is a zero of Ψ , i.e., $\Psi(\theta_0) = 0$. Expand $\Psi_n(\hat{\theta})$ in a Taylor series around θ_0 , we should have

$$0 = \Psi_n(\hat{\theta}) = \Psi_n(\theta_0) + (\hat{\theta} - \theta_0) \Psi_n'(\theta_0) + \frac{1}{2} (\hat{\theta} - \theta_0) \Psi_n''(\theta_0). \quad (20)$$

503 Rearrange the equation above, we have

$$\begin{aligned}\hat{\theta} - \theta_0 &= -\frac{\Psi_n(\theta_0)}{\Psi'_n(\theta_0) + \frac{1}{2}(\hat{\theta} - \theta_0)\Psi''_n(\theta_0)} \\ &= -\frac{\frac{1}{n}\sum_{i=1}^n \psi_{\theta}(\mathbf{x}_i)}{\Psi'_n(\theta_0) + \frac{1}{2}(\hat{\theta} - \theta_0)\Psi''_n(\theta_0)}.\end{aligned}\quad (21)$$

504 According to the central limit theorem, the numerator will be asymptotic normal with variance
505 $P\psi_{\theta_0}^2/n$ as the mean $\Psi(\theta_0) = 0$ is zero. The first term of denominator $\Psi'_n(\theta_0)$ will converge in
506 probability to $\Psi'(\theta_0)$ according to the law of large numbers. The second term $\hat{\theta} - \theta_0 = o_P(1)$.¹
507 As long as the denominator converges in probability and the numerator converges in distribution,
508 according to Slutsky's lemma, we have

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightsquigarrow N\left(0, \frac{P\psi_{\theta_0}^2}{(P\psi'_{\theta_0})^2}\right).\quad (22)$$

509 Extend into the high-dimensional case we should have

$$\hat{\theta} - \theta_0 = -(\Psi'_n(\theta_0))^{-1}\Psi_n(\theta_0),\quad (23)$$

510 where the second order term is omitted, further assume the matrix $P\psi'_{\theta_0}$ is invertible, we have

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightsquigarrow N\left(0, (P\psi'_{\theta_0})^{-1}P\psi_{\theta_0}\psi_{\theta_0}^T(P\psi'_{\theta_0})^{-1}\right),\quad (24)$$

511 Specifically, in our case $\theta_0 = (\sigma_{j_1, j_2}, \mathbf{\Lambda})$, where $\mathbf{\Lambda}$ is another parameter set influencing the estimation
512 of σ_{j_1, j_2} (will discuss case in case in later proof). In the practical scenario, we only have access to
513 the estimated parameter $\hat{\theta}$ and the empirical distribution \mathbb{P}_n , thus we have

$$\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} \overset{\text{approx}}{\rightsquigarrow} N\left(0, ((\mathbb{P}_n\psi'_{\hat{\theta}})^{-1}\mathbb{P}_n\psi_{\hat{\theta}}\psi_{\hat{\theta}}^T(\mathbb{P}_n\psi'_{\hat{\theta}})^{-1})_{1,1}\right).\quad (25)$$

514 Under the null hypothesis of independent, $\sigma_{j_1, j_2}=0$. We provide the proof that $\hat{\theta} \xrightarrow{P} \theta_0$ of our case
515 in App. A.1. Thus, $\mathbb{P}_n\psi_{\hat{\theta}}$, the function parameterized by $\hat{\theta}$, should also converge in $\mathbb{P}_n\psi_{\hat{\theta}_0}$ when
516 $n \rightarrow \infty$. Besides, by the law of large numbers, $\mathbb{P}_n\psi_{\hat{\theta}_0}$ will converge to $P\psi_{\hat{\theta}_0}$. Thus, the equation
517 above will converge to Eq. (24) when $n \rightarrow \infty$.

518 A.4 Derivation of Lem. 2.7

519 Let's first focus on the most challenging case where both variables are discretized observations
520 and our interested parameter will include $\hat{\theta} = (\hat{\sigma}_{j_1, j_2}, \hat{h}_{j_1}, \hat{h}_{j_2})$ (Although we only care about the
521 distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$, the estimation of boundary \hat{h}_{j_1} and \hat{h}_{j_2} will influence the estimation of
522 $\hat{\sigma}_{j_1, j_2}$, thus we need to consider all of them).

523 The next step will be to *construct an appropriate criterion function ψ such that $\Psi_n(\hat{\theta}) = \mathbf{0}$* . Given n
524 observations $\{\tilde{\mathbf{x}}^1, \tilde{\mathbf{x}}^2, \dots, \tilde{\mathbf{x}}^n\}$, which are discretized version of $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$ we should have

$$\Psi_n(\hat{\theta}) = \begin{pmatrix} \Psi_n(\hat{\sigma}_{j_1, j_2}) \\ \Psi_n(\hat{h}_{j_1}) \\ \Psi_n(\hat{h}_{j_2}) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \psi_{\hat{\theta}}(\tilde{\mathbf{x}}^i) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\hat{\sigma}_{j_1, j_2}; \{\hat{h}_{j_1}, \hat{h}_{j_2}\}) \\ \hat{\tau}_{j_1}^i - \bar{\Phi}(\hat{h}_{j_1}) \\ \hat{\tau}_{j_2}^i - \bar{\Phi}(\hat{h}_{j_2}) \end{pmatrix} = \mathbf{0}.\quad (26)$$

525

$$\Psi_n(\theta_0) = \begin{pmatrix} \Psi_n(\sigma_{j_1, j_2}) \\ \Psi_n(h_{j_1}) \\ \Psi_n(h_{j_2}) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \psi_{\theta_0}(\tilde{\mathbf{x}}^i) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\sigma_{j_1, j_2}; \{h_{j_1}, h_{j_2}\}) \\ \hat{\tau}_{j_1}^i - \bar{\Phi}(h_{j_1}) \\ \hat{\tau}_{j_2}^i - \bar{\Phi}(h_{j_2}) \end{pmatrix}.\quad (27)$$

526 The difference between the estimated parameter with the true parameter can be expressed as

$$\begin{aligned}\hat{\theta} - \theta_0 &= \begin{pmatrix} \hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} \\ \hat{h}_{j_1} - h_{j_1} \\ \hat{h}_{j_2} - h_{j_2} \end{pmatrix} = -\frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial \sigma_{j_1, j_2}} & \frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial h_{j_1}} & \frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial h_{j_2}} \\ \frac{\partial \Psi_n(h_{j_1})}{\partial \sigma_{j_1, j_2}} & \frac{\partial \Psi_n(h_{j_1})}{\partial h_{j_1}} & \frac{\partial \Psi_n(h_{j_1})}{\partial h_{j_2}} \\ \frac{\partial \Psi_n(h_{j_2})}{\partial \sigma_{j_1, j_2}} & \frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_1}} & \frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_2}} \end{pmatrix}^{-1} \\ &\quad \cdot \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\sigma_{j_1, j_2}; \{h_{j_1}, h_{j_2}\}) \\ \hat{\tau}_{j_1}^i - \bar{\Phi}(h_{j_1}) \\ \hat{\tau}_{j_2}^i - \bar{\Phi}(h_{j_2}) \end{pmatrix},\end{aligned}\quad (28)$$

¹We will not provide proof of this in this paper; however, interested readers may refer to [33]

527 where the specific form of each entry of the gradient matrix is expressed as

$$\begin{aligned}
\frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial \sigma_{j_1, j_2}} &= -\frac{1}{2\pi\sqrt{1-\sigma_{j_1, j_2}^2}} \exp\left(-\frac{(h_{j_1}^2 - 2\sigma_{j_1, j_2}h_{j_1}h_{j_2} + h_{j_2}^2)}{2(1-\sigma_{j_1, j_2}^2)}\right); \\
\frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial h_{j_1}} &= \int_{h_{j_2}}^{\infty} \frac{1}{2\pi\sqrt{1-\sigma_{j_1, j_2}^2}} \exp\left(-\frac{h_{j_1}^2 - 2\sigma_{j_1, j_2}h_{j_1}x_2 + x_2^2}{2(1-\sigma_{j_1, j_2}^2)}\right) dx_2; \\
\frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial h_{j_2}} &= \int_{h_{j_1}}^{\infty} \frac{1}{2\pi\sqrt{1-\sigma_{j_1, j_2}^2}} \exp\left(-\frac{h_2^2 - 2\sigma_{j_1, j_2}h_{j_2}x_1 + x_1^2}{2(1-\sigma_{j_1, j_2}^2)}\right) dx_1; \\
\frac{\partial \Psi_n(h_{j_1})}{\partial \sigma_{j_1, j_2}} &= 0; \\
\frac{\partial \Psi_n(h_{j_1})}{\partial h_{j_1}} &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{h_{j_1}^2}{2}\right); \\
\frac{\partial \Psi_n(h_{j_1})}{\partial h_{j_2}} &= 0; \\
\frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_2}} &= 0; \\
\frac{\partial \Psi_n(h_{j_2})}{\partial \sigma_{j_1, j_2}} &= 0; \\
\frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_1}} &= 0; \\
\frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_2}} &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{h_{j_2}^2}{2}\right).
\end{aligned} \tag{29}$$

528 For simplicity of notation, we define

$$\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n \xi_{j_1, j_2}^i, \tag{30}$$

529 where the specific form of $\{\xi_{j_1, j_2}^i\}$ is defined in Eq. (28). We should note that $\{\xi_{j_1, j_2}^i\}$ are i.i.d
530 random variables with mean zero (this property will be the key to the derivation of inference of CI).
531 As long as our estimation $\hat{\theta}$ converge in probability to θ_0 as proved in A.1, we have

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightsquigarrow N\left(0, ((P\psi'_{\theta_0})^{-1}P\psi_{\theta_0}\psi_{\theta_0}^T(P\psi_{\theta_0}^T)^{-1})_{1,1}\right), \tag{31}$$

532 where ψ_{θ_0} is defined in Eq. (27). However, in practice, we don't have access to either P or θ_0 . In this
533 scenario, we can plug in the empirical distribution of $\mathbb{P}_n\psi_{\hat{\theta}}$ to get the estimated variance, i.e., the
534 actual variance used in the calculation of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ is

$$\frac{1}{n} \left((\mathbb{P}_n\psi'_{\hat{\theta}})^{-1} \mathbb{P}_n\psi_{\hat{\theta}}\psi_{\hat{\theta}}^T (\mathbb{P}_n\psi_{\hat{\theta}}^T)^{-1} \right)_{1,1}. \tag{32}$$

535 A.5 Derivation of Lem. 2.8

536 Use the same line of procedure as in the derivation of Lem. 2.7, for mixed pair of observations where
537 X_{j_1} is continuous and \tilde{X}_{j_2} is discrete, we can construct the criterion function

$$\Psi_n(\hat{\theta}) = \begin{pmatrix} \Psi_n(\hat{\sigma}_{j_1, j_2}) \\ \Psi_n(\hat{h}_{j_2}) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \psi_{\hat{\theta}}(\tilde{\mathbf{x}}^i) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\hat{\sigma}_{j_1, j_2}; \{0, \hat{h}_{j_2}\}) \\ \hat{\tau}_{j_2}^i - \Phi(\hat{h}_{j_2}) \end{pmatrix} = \mathbf{0}. \tag{33}$$

538

$$\Psi_n(\theta_0) = \begin{pmatrix} \Psi_n(\sigma_{j_1, j_2}) \\ \Psi_n(h_{j_2}) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \psi_{\theta_0}(\tilde{\mathbf{x}}^i) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\sigma_{j_1, j_2}; \{0, h_{j_2}\}) \\ \hat{\tau}_{j_2}^i - \Phi(h_{j_2}) \end{pmatrix}. \tag{34}$$

539 The difference between the estimated parameter with the true parameter can be expressed as

$$\hat{\theta} - \theta_0 = \begin{pmatrix} \hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} \\ \hat{h}_{j_2} - h_{j_2} \end{pmatrix} = -\frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial \sigma_{j_1, j_2}} & \frac{\partial \Psi_n(\sigma_{j_1, j_2})}{\partial h_{j_2}} \\ \frac{\partial \Psi_n(h_{j_2})}{\partial \sigma_{j_1, j_2}} & \frac{\partial \Psi_n(h_{j_2})}{\partial h_{j_2}} \end{pmatrix}^{-1} \begin{pmatrix} \hat{\tau}_{j_1, j_2}^i - T(\sigma_{j_1, j_2}; \{0, h_{j_2}\}) \\ \hat{\tau}_{j_2}^i - \Phi(h_{j_2}) \end{pmatrix}, \quad (35)$$

540 where the specific form of each entry of the gradient matrix can be found in Eq. (29). Using exactly
 541 the same procedure, we should have the same formation of the variance calculated as Eq. (32) with a
 542 different definition of ψ_{θ_0} and $\psi_{\hat{\theta}}$ defined in Eq. (34) (33).

543 A.6 Derivation of Lem. 2.6

544 Use the same line of procedure as in derivation of Lem. 2.7, for a continuous pair of variables, we
 545 can construct the criterion function

$$\Psi_n(\hat{\theta}) = \Psi_n(\hat{\sigma}_{j_1, j_2}) = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i - \frac{1}{n} \sum_{i=1}^n x_{j_1}^i \frac{1}{n} \sum_{i=1}^n x_{j_2}^i - \hat{\sigma}_{j_1, j_2} = 0. \quad (36)$$

546

$$\Psi_n(\theta_0) = \Psi_n(\sigma_{j_1, j_2}) = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i - \frac{1}{n} \sum_{i=1}^n x_{j_1}^i \frac{1}{n} \sum_{i=1}^n x_{j_2}^i - \sigma_{j_1, j_2}. \quad (37)$$

547 Denote $\frac{1}{n} \sum_{i=1}^n x_{j_1}^i$ as \bar{x}_{j_1} and $\frac{1}{n} \sum_{i=1}^n x_{j_2}^i$ as \bar{x}_{j_2} . We should have

$$\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2} = \frac{1}{n} \sum_{i=1}^n x_{j_1}^i x_{j_2}^i - \bar{x}_{j_1} \bar{x}_{j_2} - \sigma_{j_1, j_2}. \quad (38)$$

548 According to Eq. (22), we have

$$\sqrt{n}(\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}) \rightsquigarrow N\left(0, \frac{P\psi_{\theta_0}^2}{(P\psi'_{\theta_0})^2}\right). \quad (39)$$

549 where $(P\psi'_{\theta_0})^2 = 1$. In practical calculation, we have the variance

$$\frac{1}{n} \mathbb{P}_n \psi_{\hat{\theta}}^2 / (\mathbb{P}_n \psi'_{\hat{\theta}})^2 = \frac{1}{n^2} \sum_{i=1}^n (x_{j_1}^i x_{j_2}^i - \bar{x}_{j_1} \bar{x}_{j_2} - \hat{\sigma}_{j_1, j_2})^2. \quad (40)$$

550 A.7 Proof of Thm. 2.9

551 A.7.1 Proof of Relation between Σ, Ω with β

552 Consider our latent continuous variables $\mathbf{X} = (X_1, \dots, X_p) \sim N(0, \Sigma)$ and do nodewise regression

$$X_j = X_{-j}\beta_j + \epsilon_j. \quad (41)$$

553 We can divide its covariance Σ and its precision matrix $\Omega = \Sigma^{-1}$ into X and Y part in our regression:

$$\Sigma = \begin{pmatrix} \Sigma_{jj} & \Sigma_{j-j} \\ \Sigma_{-jj} & \Sigma_{-j-j} \end{pmatrix} \quad \Omega = \begin{pmatrix} \Omega_{jj} & \Omega_{j-j} \\ \Omega_{-jj} & \Omega_{-j-j} \end{pmatrix}. \quad (42)$$

554 Just like regular linear regression, we can get

$$n \rightarrow \infty, \quad \beta_j = \Sigma_{-j-j}^{-1} \Sigma_{-jj}. \quad (43)$$

555 From the invertibility of a block matrix

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} (A - BD^{-1}C)^{-1} & -(A - BD^{-1}C)^{-1}BD^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1} \end{bmatrix}. \quad (44)$$

556 If A and D is invertible, we will have

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} A - BD^{-1}C & 0 \\ 0 & (D - CA^{-1}B)^{-1} \end{bmatrix} \begin{bmatrix} I & -BD^{-1} \\ -CA^{-1} & I \end{bmatrix}. \quad (45)$$

557 Thus, we can get:

$$\begin{aligned}\Omega_{jj} &= \Sigma_{jj} - (\Sigma_{j-j}\Sigma_{-j-j}^{-1}\Sigma_{-jj})^{-1}; \\ \Omega_{j-j} &= -(\Sigma_{jj} - (\Sigma_{j-j}\Sigma_{-j-j}^{-1}\Sigma_{-jj})^{-1})\Sigma_{j-j}(\Sigma_{-j-j})^{-1}.\end{aligned}\quad (46)$$

558 Move one step forward:

$$-\Omega_{jj}^{-1}\Omega_{j-j} = \Sigma_{j-j}(\Sigma_{-j-j})^{-1}. \quad (47)$$

559 Take transpose for both sides, as long as Ω is a symmetric matrix and $\Omega_{-jj} = \Omega_{j-j}^T$, we will have

$$-\Omega_{jj}^{-1}\Omega_{-jj} = \Sigma_{-j-j}^{-1}\Sigma_{-jj} = \beta_j. \quad (48)$$

560 We should note testing $\Omega_{-jj} = 0$ is equivalent to testing $\beta_j = 0$ as the Ω_{jj} will always be nonzero.
561 The variable Ω_{-jj} captures the CI of X_j with other variables. As long as the variable Ω_{jj} is just one
562 scalar, we can get

$$\beta_{j,k} = -\frac{\omega_{j,k}}{\omega_{j,j}} \quad (49)$$

563 capturing the independence relationship between variable X_j with X_k conditioning on all other
564 variables.

565 A.7.2 Detailed derivation of inference for β_j

566 Nodewise regression allows us to use the regression parameter β_j as the surrogate of Ω_{-jj} . The
567 problem now transfers to constructing the inference for β_j , specifically, the derivation of distribution
568 of $\hat{\beta}_j - \beta_j$. The overarching concept is that we are already aware of the distribution of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$
569 and we know that there exists a deterministic relationship between β_j with Σ . Consequently, we can
570 express $\hat{\beta}_j - \beta_j$ as a composite of $\hat{\sigma}_{j_1, j_2} - \sigma_{j_1, j_2}$ to establish such an inference. Specifically, we have

$$\begin{aligned}\hat{\beta}_j - \beta_j &= \hat{\Sigma}_{-j-j}^{-1}\hat{\Sigma}_{-jj} - \Sigma_{-j-j}^{-1}\Sigma_{-jj} \\ &= \hat{\Sigma}_{-j-j}^{-1}\left(\hat{\Sigma}_{-jj} - \hat{\Sigma}_{-j-j}\Sigma_{-j-j}^{-1}\Sigma_{-jj}\right) \\ &= -\hat{\Sigma}_{-j-j}^{-1}\left(\hat{\Sigma}_{-j-j}\beta_j - \Sigma_{-j-j}\beta_j + \Sigma_{-j-j}\beta_j - \hat{\Sigma}_{-jj}\right) \\ &= -\hat{\Sigma}_{-j-j}^{-1}\left((\hat{\Sigma}_{-j-j} - \Sigma_{-j-j})\beta_j - (\hat{\Sigma}_{-jj} - \Sigma_{-jj})\right),\end{aligned}\quad (50)$$

571 where each entry in matrix $(\hat{\Sigma}_{-j-j} - \Sigma_{-j-j})$ and $(\hat{\Sigma}_{-jj} - \Sigma_{-jj})$ denotes the difference between
572 estimated covariance with true covariance. Suppose that we want to test the CI of the variable X_1
573 with other variables, $j = 1$, then

$$\hat{\Sigma}_{-j-j} - \Sigma_{-j-j} = \begin{bmatrix} \hat{\sigma}_{1,1} \dots \hat{\sigma}_{1,j-1}, \hat{\sigma}_{1,j+1} \dots \hat{\sigma}_{1,p} \\ \dots \\ \hat{\sigma}_{j-1,1} \dots \hat{\sigma}_{j-1,j-1}, \hat{\sigma}_{j-1,j+1} \dots \hat{\sigma}_{j-1,p} \\ \dots \\ \hat{\sigma}_{p,1} \dots \hat{\sigma}_{p,j-1}, \hat{\sigma}_{p,j+1} \dots \hat{\sigma}_{p,p} \end{bmatrix} \quad (51)$$

$$- \begin{bmatrix} \sigma_{1,1} \dots \sigma_{1,j-1}, \sigma_{1,j+1} \dots \sigma_{1,p} \\ \dots \\ \sigma_{j-1,1} \dots \sigma_{j-1,j-1}, \sigma_{j-1,j+1} \dots \sigma_{j-1,p} \\ \dots \\ \sigma_{p,1} \dots \sigma_{p,j-1}, \sigma_{p,j+1} \dots \sigma_{p,p} \end{bmatrix}. \quad (52)$$

574 Suppose that we want to test the CI of the variable X_1 with other variables, $j = 1$. then

$$\hat{\Sigma}_{-1-1} - \Sigma_{-1-1} = \begin{bmatrix} \hat{\sigma}_{2,2} \dots \hat{\sigma}_{2,p} \\ \dots \\ \hat{\sigma}_{p,2} \dots \hat{\sigma}_{p,p} \end{bmatrix} - \begin{bmatrix} \sigma_{2,2} \dots \sigma_{2,p} \\ \dots \\ \sigma_{p,2} \dots \sigma_{p,p} \end{bmatrix} \quad (53)$$

$$:= \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} \xi_{2,2}^i \dots \xi_{2,p}^i \\ \dots \\ \xi_{p,2}^i \dots \xi_{p,p}^i \end{bmatrix}, \quad (54)$$

575 where $\{\xi_{j_1, j_2}^i\}$ are i.i.d random variables with specific form defined in Eq. (28) for discrete case,
 576 Eq. (35) for mixed case and Eq. (38) in continuous case. Put them together:

$$\begin{bmatrix} \hat{\beta}_{1,2} - \beta_{1,2} \\ \hat{\beta}_{1,3} - \beta_{1,3} \\ \dots \\ \hat{\beta}_{1,p} - \beta_{1,p} \end{bmatrix} = -\hat{\Sigma}_{-1,-1}^{-1} \frac{1}{n} \sum_{i=1}^n \left(\begin{bmatrix} \xi_{2,2}^i & \xi_{2,3}^i & \dots & \xi_{2,p}^i \\ \xi_{3,2}^i & \xi_{3,3}^i & \dots & \xi_{3,p}^i \\ \dots & \dots & \dots & \dots \\ \xi_{p,2}^i & \xi_{p,3}^i & \dots & \xi_{p,p}^i \end{bmatrix} \begin{bmatrix} \beta_{1,2} \\ \beta_{1,3} \\ \dots \\ \beta_{1,p} \end{bmatrix} - \begin{bmatrix} \xi_{2,1}^i \\ \xi_{3,1}^i \\ \dots \\ \xi_{p,1}^i \end{bmatrix} \right). \quad (55)$$

577 As $\frac{1}{n} \sum_{i=1}^n \xi_{j_1, j_2}^i$ is asymptotically normal, the who vector of $\hat{\beta}_1 - \beta_1$ is a linear combination of
 578 Gaussian distribution. However, We cannot merely engage in a linear combination of its variance as
 579 they are dependent with each other. For example, if Y_1, Y_2 are dependent and we are trying to find
 580 out $Var(aY_1 + bY_2)$, we should have

$$Var(aY_1 + bY_2) = \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} Var(Y_1) & Cov(Y_1, Y_2) \\ Cov(Y_1, Y_2) & Var(Y_2) \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}. \quad (56)$$

581 Now, suppose we are interested in the distribution of $\hat{\beta}_{1,2} - \beta_{1,2}$, we should have

$$\hat{\beta}_{1,2} - \beta_{1,2} = \frac{1}{n} \sum_{i=1}^n (\hat{\Sigma}_{-1,-1}^{-1})_{[2],:} \left(\begin{bmatrix} \xi_{2,2}^i & \xi_{2,3}^i & \dots & \xi_{2,p}^i \\ \xi_{3,2}^i & \xi_{3,3}^i & \dots & \xi_{3,p}^i \\ \dots & \dots & \dots & \dots \\ \xi_{p,2}^i & \xi_{p,3}^i & \dots & \xi_{p,p}^i \end{bmatrix} \begin{bmatrix} \beta_{1,2} \\ \beta_{1,3} \\ \dots \\ \beta_{1,p} \end{bmatrix} - \begin{bmatrix} \xi_{2,1}^i \\ \xi_{3,1}^i \\ \dots \\ \xi_{p,1}^i \end{bmatrix} \right), \quad (57)$$

582 where $(\hat{\Sigma}_{-1,-1}^{-1})_{[2],:}$ is the row of index of X_2 of $\hat{\Sigma}_{-1,-1}^{-1}$ ($[2]$ denotes the index of the variable). For
 583 ease of notation, let

$$\Xi_{-1,-1}^i = \begin{bmatrix} \xi_{2,2}^i & \xi_{2,3}^i & \dots & \xi_{2,p}^i \\ \xi_{3,2}^i & \xi_{3,3}^i & \dots & \xi_{3,p}^i \\ \dots & \dots & \dots & \dots \\ \xi_{p,2}^i & \xi_{p,3}^i & \dots & \xi_{p,p}^i \end{bmatrix}, \quad \Xi_{-1,1}^i = \begin{bmatrix} \xi_{2,1}^i \\ \xi_{3,1}^i \\ \dots \\ \xi_{p,1}^i \end{bmatrix}, \quad (58)$$

584 and let

$$B_{-1}^i = \begin{pmatrix} \xi_{2,1}^i & \xi_{3,1}^i & \dots & \xi_{p,1}^i \\ \xi_{2,2}^i & \xi_{2,3}^i & \dots & \xi_{2,p}^i \\ \xi_{3,2}^i & \xi_{3,3}^i & \dots & \xi_{3,p}^i \\ \dots & \dots & \dots & \dots \\ \xi_{p,2}^i & \xi_{p,3}^i & \dots & \xi_{p,p}^i \end{pmatrix} \quad (59)$$

585 as the concatenation of those two matrices. The variance is calculated as

$$Var\left(\sqrt{n}(\hat{\beta}_{1,2} - \beta_{1,2})\right) = a^{[2]T} \frac{1}{n} \sum_{i=1}^n vec(B_{-1}^i) vec(B_{-1}^i)^T a^{[2]}, \quad (60)$$

586 where

$$a_l^{[2]} = \begin{cases} (\hat{\Sigma}_{-1,-1}^{-1})_{[2],l}, & \text{for } l \in \{1, \dots, p-1\} \\ \sum_{q=1}^n (\hat{\Sigma}_{-1,-1}^{-1})_{[2],l} (\beta_1)_q, & \text{for } l \in \{p, \dots, p^2 - p\} \end{cases} \quad (61)$$

587 $vec(B_{-1}^i)$ is the squeezed vector form of matrix $vec(B_{-1}^i) \in \mathbb{R}^{p \times p-1}$, i.e.,

$$vec(B_{-1}^i) = \begin{pmatrix} \xi_{2,1}^i \\ \xi_{3,1}^i \\ \vdots \\ \xi_{p,p}^i \end{pmatrix}. \quad (62)$$

588 Thus, the distribution of $\hat{\beta}_{j,k} - \beta_{j,k}$ is

$$\hat{\beta}_{j,k} - \beta_{j,k} \sim N\left(0, a^{[k]T} \frac{1}{n^2} \sum_{i=1}^n vec(B_{-j}^i) vec(B_{-j}^i)^T a^{[k]}\right). \quad (63)$$

589 In practice, we can plug in the estimates of β_j to estimate the interested distribution and do the CI
 590 test by hypothesizing $\beta_{j,k} = 0$.

591 **A.8 Discussion of assumption of zero mean and identity variance**

592 In this section, we engage in a more thorough discussion regarding our assumptions about \mathbf{X} .
 593 Specifically, we demonstrate that this assumption of mean and variance does not compromise the
 594 generality. In other words, the true model may possess different mean and variance values, but we
 595 proceed by treating it as having a mean of zero and identity variance.

The key ingredient allowing us to assume such a model is, the discretization function g_j is an unknown nonlinear monotonic function. Suppose the g'_j maps the continuous domain to a binary variable, and we have the "groundtruth" variable, denoted X'_j , with mean a and variance b . Assume the cardinality of the discretized domain is only 2, i.e., our observation \tilde{X}_j can only be 0 or 1. We further have the constant d'_j as the discretization boundary such that we have the observation

$$\tilde{X}_j = \mathbb{1}(g'_j(X'_j) > d'_j) = \mathbb{1}(X'_j > g_j^{-1}(d'_j))$$

596 We can always produce our assumed variable X_j with mean 0 and variance 1, such that $X_j =$
 597 $\frac{1}{\sqrt{b}}X'_j - \frac{a}{\sqrt{b}}$ and the same observation with a different nonlinear transformation g_j and decision
 598 boundary d_j , such that

$$\tilde{X}_j = \mathbb{1}(g_j(X_j) > d_j) = \mathbb{1}(X_j > g_j^{-1}(d_j)) = \mathbb{1}(X'_j > \sqrt{b}g_j^{-1}(d_j) + a)$$

599 As long as the observation \tilde{X}_j is the same, we should have $\sqrt{b}g_j^{-1}(d_j) + a = g'_j^{-1}(d'_j)$. Our assumed
 600 model X_j clearly mimics the "groundtruth" X'_j . Besides, according to Lem. A.2, we have one-to-
 601 one mapping between $\hat{\tau}_{j_1, j_2}$ with the estimated covariance for fixed $\hat{h}_{j_1}, \hat{h}_{j_2}$. Thus, as long as the
 602 observation is the same, the estimation of covariance $\hat{\sigma}_{j_1, j_2}$ remains unaffected by our assumptions
 603 regarding the mean and variance of \mathbf{X} , so do the following inference.

604 We further conduct casual discovery experiments to empirically validate our statement, which is
 605 shown in App. C.3.

606 **B Data Generation and Figure of main experiments: causal discovery**

607 **Data Generation and Code** We construct the true DAG \mathcal{G} using the Bipartite Pairing (BP) model
 608 [2], with the number of edges being one fewer than the number of nodes. The subsequent generation of
 609 true multivariate Gaussian data involves assigning causal weights drawn from a uniform distribution
 610 $U \sim (0.5, 2)$ and incorporating noise via samples from a standard normal distribution for each
 611 variable. Following this, we binarize the data, setting the threshold randomly based on each variable's
 612 range. The code implementation is based on [40].

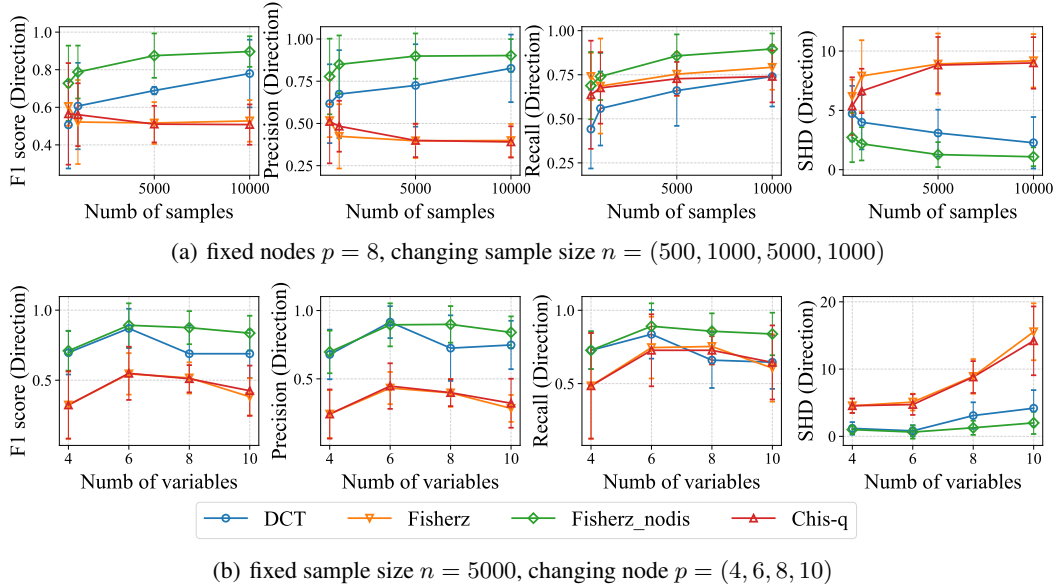


Figure 4: Experiment result of DAG discovery on synthetic data for changing sample size (a) and changing number of nodes (b). Fisherz_nodis is the Fisher-z test applied to original continuous data. We evaluate F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow).

613 C Additional experiments

614 C.1 Linear non-Gaussian and nonlinear

615 Our model requires that the original data must adhere to the hypothesis of following a multivariate
616 normal distribution, which appears to potentially limit the generalizability. Therefore, it is worthwhile
617 to explore its robustness when such assumptions are violated. In this regard, we conducted several
618 experiments, including scenarios involving linear non-Gaussian and nonlinear Gaussian.

619 For both cases, we follow the setting of our experiment where there are $p = 8$ nodes and $p - 1$
620 edges. We explore the effect of changing sample size $n = (100, 500, 2000, 5000)$. Specifically for
621 linear non-Gaussian case, we adhere to some of the settings outlined by [28], conducting experiments
622 where the original continuous data followed: (1) a Student’s t-distribution with 3 degrees of freedom,
623 (2) a uniform distribution, and (3) an exponential distribution. Each variable is generated as $X_i =$
624 $f(PA_i) + noise$, where $noise$ follows the distribution in (1), (2), (3) correspondingly and f is a
625 linear function. The first three rows of Fig. 5 and Fig. 6 show the result of the linear non-Gaussian
626 case.

627 For the nonlinear cases, we follow setting in [19], where every variable X_i is generated as $X_i =$
628 $f(WPA_i + noise)$, $noise \sim N(0, 1)$ and f is a function randomly chosen from (a) $f(x) = \sin(x)$,
629 (b) $f(x) = x^3$, (c) $f(x) = \tanh(x)$, and (d) $f(x) = ReLU(x)$. W is a linear function. Similarly,
630 we set the number of nodes at $p = 8$ and change the number of samples $n = (500, 2000, 5000)$.
631 For both cases, we run 10 graph instances with different seeds and report the result of skeleton
632 discovery in Fig. 5 and DAG in Fig. 6 (The same orientation rules [11] used in the main experiment
633 are employed to convert a CPDAG [6] into a DAG). The last row of Fig. 5 and Fig. 6 shows the result
634 of the nonlinear case.

635 Based on the experimental outcomes, DCT demonstrates marginally superior or comparable efficacy
636 in terms of the F1-score, precision, and SHD relative to both the Fisher-Z test and the Chi-square test
637 when dealing with small sample sizes. Nevertheless, as the sample size increases, DCT’s performance
638 clearly surpasses that of the aforementioned tests across all three evaluated metrics, especially in the
639 linear case. Consistent with observations from the main experiment, DCT exhibits a lower recall in
640 comparison to the baseline tests. This discrepancy can be attributed to the baseline tests being prone
641 to incorrectly infer conditional dependence and connect a large proportion of nodes. According to
642 the results, our test shows notable robustness under the case assumptions are violated, confirming its
643 practical effectiveness.

644 C.2 Denser graph

645 DCT primarily works on cases where CI is mistakenly judged as conditional dependence due
646 to discretization. Consequently, its efficacy is more pronounced in scenarios characterized by a
647 relatively sparse graph, as numerous instances are truly conditionally independent. Nevertheless, the
648 investigation of causal discovery with a dense latent graph is essential for evaluating the power of a
649 test, i.e., its ability to successfully reject the null hypothesis when the tested pairs are conditionally
650 dependent. Thus, we conduct the experiment where $p = 8$, $n = 10000$ and changing edges ($p +$
651 $2, p + 4, p + 6$). Similarly, the latent continuous data follows a multivariate Gaussian model and
652 the true DAG \mathcal{G} is constructed using BP model. We run 10 graph instances with different seeds and
653 report the result of the skeleton discovery and DAG in Fig. 7.

654 According to the experiment results, DCT exhibits better performance in terms of the F1-score,
655 precision, and SHD relative to both the Fisher-Z test and the Chi-square test. As the graph becomes
656 progressively denser, the superiority of the Discrete Causality Test (DCT) correspondingly diminishes
657 as there are few conditional independent cases in the true DAG. Due to the same reason, The recall
658 remains lower than that of other baseline methods.

659 C.3 multivariate Gaussian with nonzero mean and non-unit variance

660 We employed a setting nearly identical to the main experiment, with the only difference being the
661 alteration in data generation: instead of using a standard normal distribution, we used a Gaussian
662 distribution with mean sampled from $U(-2, 2)$ and variance sampled from $U(0, 3)$. We fix the
663 number of variables as $p = 8$ and change the number of samples $n = (100, 500, 2000, 5000)$. The
664 Fig. 8 shows the result and demonstrates the effectiveness of our method.

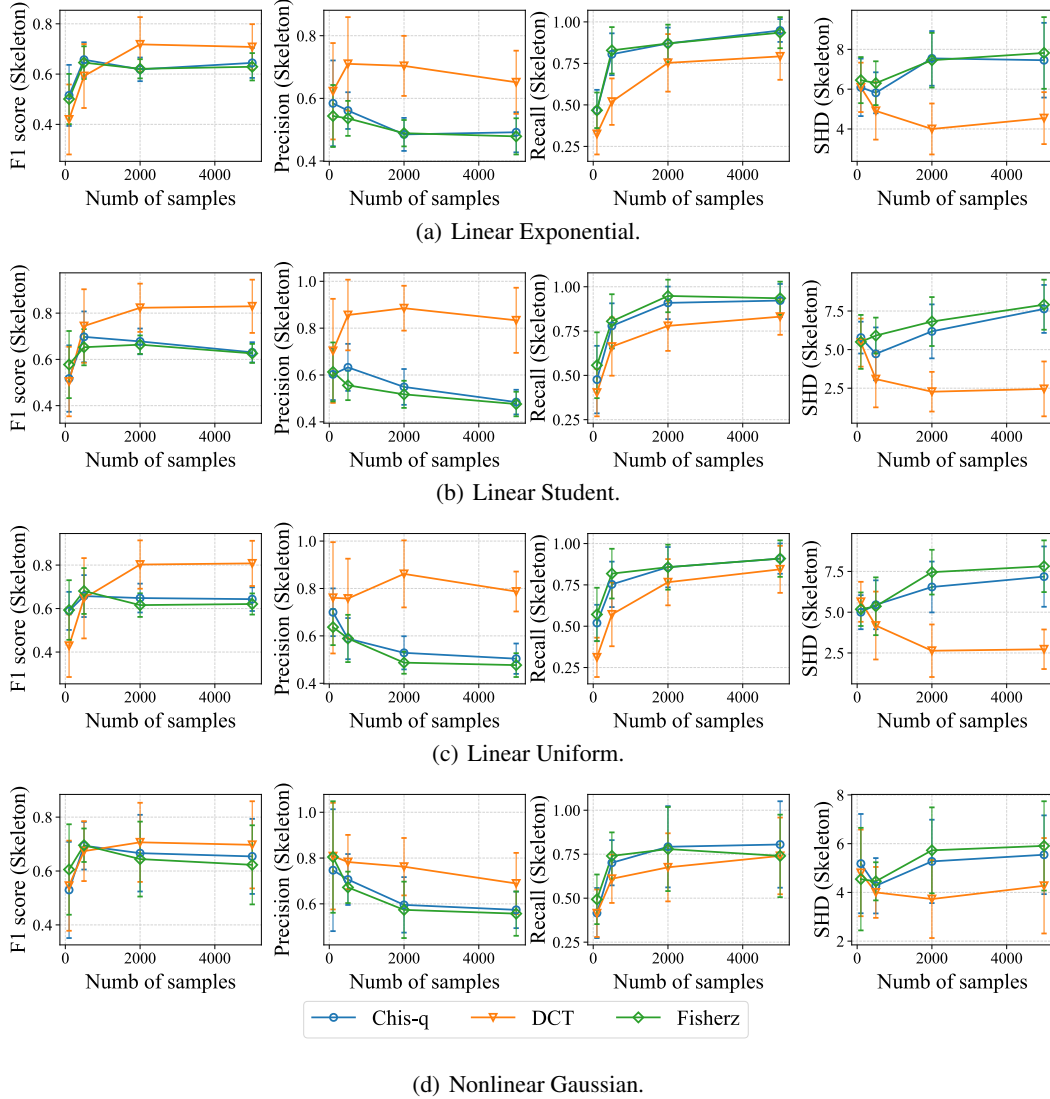


Figure 5: Experiment result of causal discovery on synthetic data with $p = 8$, $n = (100, 500, 2000, 5000)$ where the data generation process violates our assumptions. The data are generated with either nongaussian distributed (a), (b), (c) or the relations are not linear (d). The figure reports F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow) on skeleton.

665 **C.4 Real-world dataset**

666 To further validate DCT, we employ it on a real-world dataset: Big Five Personality
 667 <https://openpsychometrics.org/>, which includes 50 personality indicators and over 19000 data sam-
 668 ples. Each variable contains 5 possible discrete values to represent the scale of the corresponding
 669 questions, where 1=Disagree, 2=Weakly disagree, 3=Neutral, 4=Weakly agree and 5=Agree, e.g.,
 670 "N3=1" means "I agree that I worry about things". This scenario clearly suits DCT, where the degree
 671 of agreement with a certain question must be a continuous variable while we can only observe the
 672 result after categorization. We choose three variables respectively: [N3: I worry about things], [N10:
 673 I often feel blue], [N4: I seldom feel blue]. We then do the casual discovery using PC algorithm with
 674 DCT and compare it with the Chi-square test and Fisher-Z test. The result can be found in Fig. 9.

675 Based on the experimental outcomes, despite the absence of a groundtruth for reference, we observe
 676 that the results obtained via DCT appear more plausible than those derived from Fisher-Z and Chi-
 677 square tests. Specifically, DCT suggests the relationship $N_3 \perp\!\!\!\perp N_4|N_{10}$, which is reasonable as
 678 intuitively, the answer of 'I often feel blue' already captures the information of 'I seldom feel blue'.

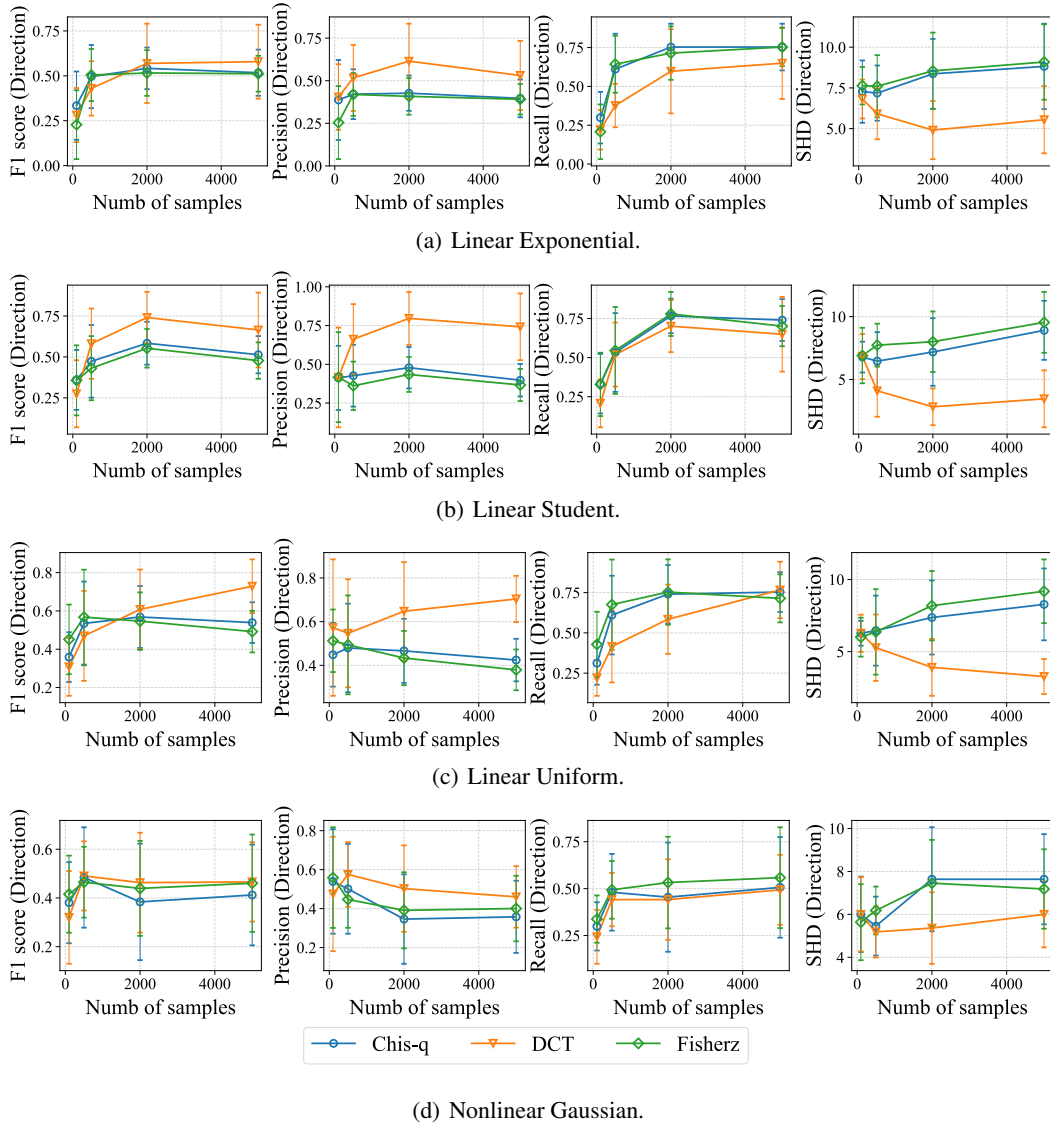


Figure 6: Experiment result of causal discovery on synthetic data with $p = 8$, $n = (100, 500, 2000, 5000)$ where the data generation process violates our assumptions. The data are generated with either nongaussian distributed (a), (b), (c) or the relations are not linear (d). The figure reports F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow) on DAG.

679 As a comparison, both Fisher-Z and Chi-square return a fully connected graph. The results directly
 680 correspond to our illustrative example shown in Fig. 1, substantiating the necessity of our proposed
 681 test.

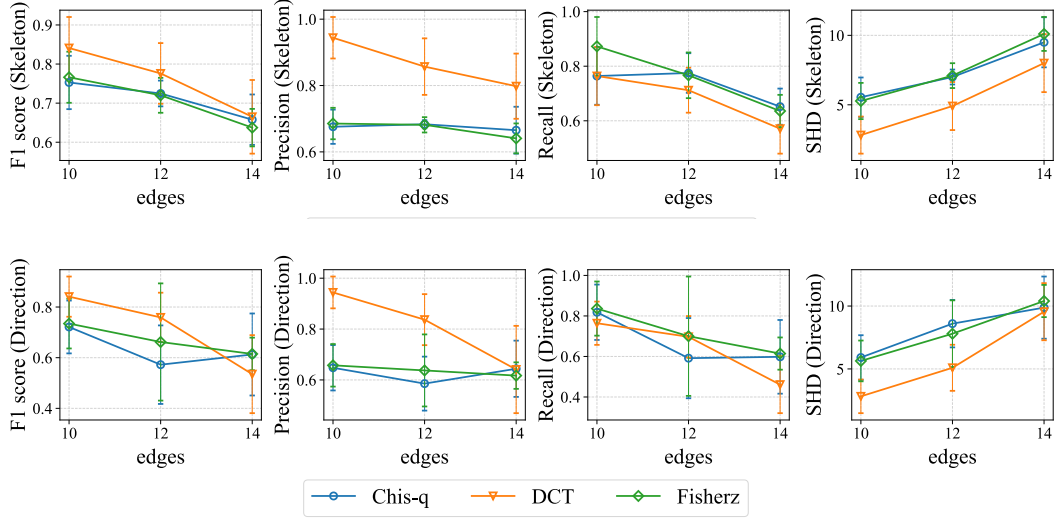


Figure 7: Experimental comparison of causal discovery on synthetic datasets for denser graphs with $p = 8, n = 10000$ and edges varying $p + 2, p + 4, p + 6$. We evaluate F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow) on both skeleton and DAG.

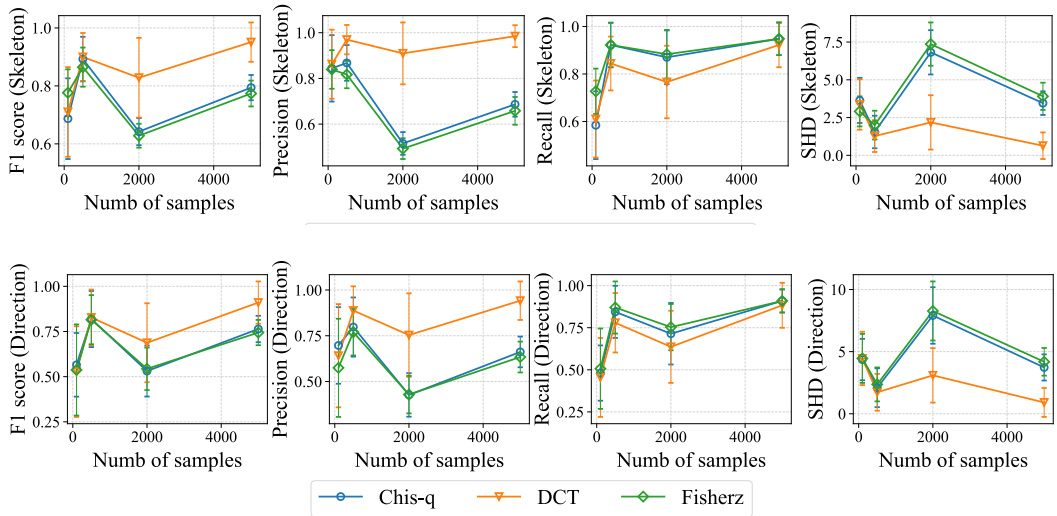


Figure 8: Experimental comparison of causal discovery on synthetic datasets for multivariate Gaussian model with $p = 8, n = (100, 500, 2000, 5000)$ and where mean is not zero. We evaluate F_1 (\uparrow), Precision (\uparrow), Recall (\uparrow) and SHD (\downarrow) on both skeleton and DAG.

682 D Related Work

683 Testing for CI is pivotal in the field of causal discovery [30], and a variety of methods exist for
 684 performing CI tests (CI tests). An important group of CI test methods involves the assumption of
 685 Gaussian variables with linear dependencies. For example, under this assumption, Gaussian graphical
 686 models are extensively studied [37, 25, 22, 26]. To address CI test under Gaussian assumption, partial
 687 correlation serves as a viable method for CI testing [4]. To evaluate the independence of variables
 688 X_1 and X_2 conditional on Z , The technique proposed by [32] determines CI by comparing the
 689 estimations of $p(X_1|X_2, Z)$ and $p(X_1|X_2)$.

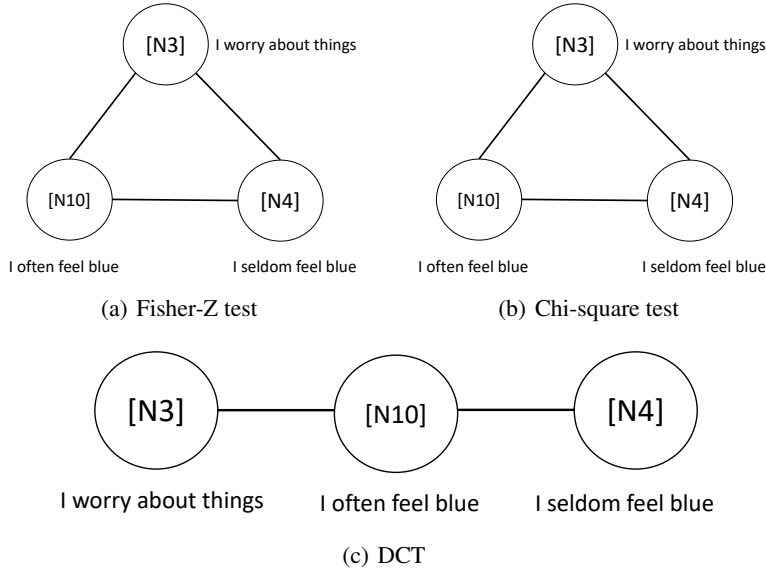


Figure 9: Experimental comparison of causal discovery on the real-world dataset.

690 Another approach involves discretizing Z and performing independent tests within each resulting bin
 691 [21]. Our work, however, diverges from these existing methods in two significant ways. Firstly, we
 692 are equipped to handle data, where partial variables are discretized. Additionally, we postulate that
 693 discrete variables are derived from the transformation of continuous variables in a latent Gaussian
 694 model. With the same assumption, the most closely related study is by [13], where the authors
 695 developed a novel rank-based estimator for the precision matrix of mixed data. However, their work
 696 stops short of providing a CI test for this method. Our research fills this gap, offering the ability to
 697 estimate the precision matrix for both discrete and mixed data and providing a rigorous CI test for
 698 our methodology.

699 Recent advancements in CI testing have utilized kernel methods for continuous variables influenced
 700 by nonlinear relationships. [16] describes non-parametric CI relationships using covariance operators
 701 in reproducing kernel Hilbert spaces (RKHS). KCI test [38] assesses the partial associations of
 702 regression functions linking x , y , and z , while RCI test [31] aims to enhance the KCI test’s efficiency.
 703 In KCIP test [12] employs permutations of samples to emulate CI scenarios. CCI test [27] further
 704 reformulates testing into a process that leverages the capabilities of supervised learning models. For
 705 discrete variable analysis, the G^2 test [1] and conditional mutual information [39] are commonly
 706 employed. However, their method cannot deal with our setting where only discretized version of
 707 latent variables can be observed.

708 E Resource Usage

709 All the experiments are run using Intel(R) Xeon(R) CPU E5-2680 v4 with 55 processors. It costs 4
 710 hours to run experiments in Section 3.1.

711 F Limiation and Broader Impacts

712 **Limitation** So far, the largest limitation of our method is to treat discretized variables as binary,
 713 which wastes the available information. Besides that, the parametric assumption limits its generaliz-
 714 ability. However, we need to point out this is pretty normal in CI test fields.

715 **Broader Impacts** The goal of our proposed method is to test the conditional independence relation-
 716 ship given discretized observation. This task is essential and has broad applications. We are confident
 717 that our method will be beneficial and will not result in negative societal impacts.

718 **NeurIPS Paper Checklist**

719 **1. Claims**

720 Question: Do the main claims made in the abstract and introduction accurately reflect the
721 paper's contributions and scope?

722 Answer: [Yes]

723 Justification: Section1 Introduction and Abstract

724 Guidelines:

- 725 • The answer NA means that the abstract and introduction do not include the claims
726 made in the paper.
- 727 • The abstract and/or introduction should clearly state the claims made, including the
728 contributions made in the paper and important assumptions and limitations. A No or
729 NA answer to this question will not be perceived well by the reviewers.
- 730 • The claims made should match theoretical and experimental results, and reflect how
731 much the results can be expected to generalize to other settings.
- 732 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
733 are not attained by the paper.

734 **2. Limitations**

735 Question: Does the paper discuss the limitations of the work performed by the authors?

736 Answer: [Yes]

737 Justification: Section2.1 line145-line147, Appendix F

738 Guidelines:

- 739 • The answer NA means that the paper has no limitation while the answer No means that
740 the paper has limitations, but those are not discussed in the paper.
- 741 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 742 • The paper should point out any strong assumptions and how robust the results are to
743 violations of these assumptions (e.g., independence assumptions, noiseless settings,
744 model well-specification, asymptotic approximations only holding locally). The authors
745 should reflect on how these assumptions might be violated in practice and what the
746 implications would be.
- 747 • The authors should reflect on the scope of the claims made, e.g., if the approach was
748 only tested on a few datasets or with a few runs. In general, empirical results often
749 depend on implicit assumptions, which should be articulated.
- 750 • The authors should reflect on the factors that influence the performance of the approach.
751 For example, a facial recognition algorithm may perform poorly when image resolution
752 is low or images are taken in low lighting. Or a speech-to-text system might not be
753 used reliably to provide closed captions for online lectures because it fails to handle
754 technical jargon.
- 755 • The authors should discuss the computational efficiency of the proposed algorithms
756 and how they scale with dataset size.
- 757 • If applicable, the authors should discuss possible limitations of their approach to
758 address problems of privacy and fairness.
- 759 • While the authors might fear that complete honesty about limitations might be used by
760 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
761 limitations that aren't acknowledged in the paper. The authors should use their best
762 judgment and recognize that individual actions in favor of transparency play an impor-
763 tant role in developing norms that preserve the integrity of the community. Reviewers
764 will be specifically instructed to not penalize honesty concerning limitations.

765 **3. Theory Assumptions and Proofs**

766 Question: For each theoretical result, does the paper provide the full set of assumptions and
767 a complete (and correct) proof?

768 Answer: [Yes]

769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822

Justification: Assumption: Section2 line81 to line 94, Proof: Appendix A.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section3 and Appendix B,C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874

Answer: [Yes]

Justification: We provide the full code in our supplementary.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 3 and Appendix B, C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Section 3 and Appendix B, C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 875 • It is OK to report 1-sigma error bars, but one should state it. The authors should
876 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
877 of Normality of errors is not verified.
- 878 • For asymmetric distributions, the authors should be careful not to show in tables or
879 figures symmetric error bars that would yield results that are out of range (e.g. negative
880 error rates).
- 881 • If error bars are reported in tables or plots, The authors should explain in the text how
882 they were calculated and reference the corresponding figures or tables in the text.

883 8. Experiments Compute Resources

884 Question: For each experiment, does the paper provide sufficient information on the com-
885 puter resources (type of compute workers, memory, time of execution) needed to reproduce
886 the experiments?

887 Answer: [Yes]

888 Justification: Appendix E.

889 Guidelines:

- 890 • The answer NA means that the paper does not include experiments.
- 891 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
892 or cloud provider, including relevant memory and storage.
- 893 • The paper should provide the amount of compute required for each of the individual
894 experimental runs as well as estimate the total compute.
- 895 • The paper should disclose whether the full research project required more compute
896 than the experiments reported in the paper (e.g., preliminary or failed experiments that
897 didn't make it into the paper).

898 9. Code Of Ethics

899 Question: Does the research conducted in the paper conform, in every respect, with the
900 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

901 Answer: [Yes]

902 Justification: We completely follow NeurIPS Code of Ethics.

903 Guidelines:

- 904 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 905 • If the authors answer No, they should explain the special circumstances that require a
906 deviation from the Code of Ethics.
- 907 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
908 eration due to laws or regulations in their jurisdiction).

909 10. Broader Impacts

910 Question: Does the paper discuss both potential positive societal impacts and negative
911 societal impacts of the work performed?

912 Answer: [Yes]

913 Justification: We propose a new conditional independence test with applications range in
914 multiple fields. Please refer to Appendix F.

915 Guidelines:

- 916 • The answer NA means that there is no societal impact of the work performed.
- 917 • If the authors answer NA or No, they should explain why their work has no societal
918 impact or why the paper does not address societal impact.
- 919 • Examples of negative societal impacts include potential malicious or unintended uses
920 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
921 (e.g., deployment of technologies that could make decisions that unfairly impact specific
922 groups), privacy considerations, and security considerations.

- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

938 11. Safeguards

939 Question: Does the paper describe safeguards that have been put in place for responsible
940 release of data or models that have a high risk for misuse (e.g., pretrained language models,
941 image generators, or scraped datasets)?

942 Answer: [NA]

943 Justification: Method proposed in this paper don't pose such risks.

944 Guidelines:

- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

955 12. Licenses for existing assets

956 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
957 the paper, properly credited and are the license and terms of use explicitly mentioned and
958 properly respected?

959 Answer: [Yes]

960 Justification: We have cited the dataset we use and we provide the code we based in
961 Appendix B.

962 Guidelines:

- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971
- 972
- 973
- 974
- 975
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 976 • If this information is not available online, the authors are encouraged to reach out to
977 the asset’s creators.

978 **13. New Assets**

979 Question: Are new assets introduced in the paper well documented and is the documentation
980 provided alongside the assets?

981 Answer: [Yes]

982 Justification: We have submitted the code.

983 Guidelines:

- 984 • The answer NA means that the paper does not release new assets.
985 • Researchers should communicate the details of the dataset/code/model as part of their
986 submissions via structured templates. This includes details about training, license,
987 limitations, etc.
988 • The paper should discuss whether and how consent was obtained from people whose
989 asset is used.
990 • At submission time, remember to anonymize your assets (if applicable). You can either
991 create an anonymized URL or include an anonymized zip file.

992 **14. Crowdsourcing and Research with Human Subjects**

993 Question: For crowdsourcing experiments and research with human subjects, does the paper
994 include the full text of instructions given to participants and screenshots, if applicable, as
995 well as details about compensation (if any)?

996 Answer: [NA]

997 Justification: We don’t use any crowdsourcing resource.

998 Guidelines:

- 999 • The answer NA means that the paper does not involve crowdsourcing nor research with
1000 human subjects.
1001 • Including this information in the supplemental material is fine, but if the main contribu-
1002 tion of the paper involves human subjects, then as much detail as possible should be
1003 included in the main paper.
1004 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1005 or other labor should be paid at least the minimum wage in the country of the data
1006 collector.

1007 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
1008 **Subjects**

1009 Question: Does the paper describe potential risks incurred by study participants, whether
1010 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1011 approvals (or an equivalent approval/review based on the requirements of your country or
1012 institution) were obtained?

1013 Answer: [NA]

1014 Justification: NA.

1015 Guidelines:

- 1016 • The answer NA means that the paper does not involve crowdsourcing nor research with
1017 human subjects.
1018 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1019 may be required for any human subjects research. If you obtained IRB approval, you
1020 should clearly state this in the paper.
1021 • We recognize that the procedures for this may vary significantly between institutions
1022 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1023 guidelines for their institution.
1024 • For initial submissions, do not include any information that would break anonymity (if
1025 applicable), such as the institution conducting the review.