

小数据集下基于 DRKDE-ICSO 的 BN 结构学习

陈海洋, 刘静, 刘喜庆, 张静

(西安工程大学电子信息学院, 西安, 710048)

摘要 为了解决在小数据集条件下进行数据拓展时产生数据高度相似的问题,提出了基于降维核密度估计的小数据集拓展方法,从而得到较为准确的拓展数据。另外,针对鸡群优化算法求解效率低下和收敛性不足的问题,提出改进的鸡群优化算法进行结构学习:在雄鸡的位置更新公式中引入莱维飞行,使鸡群算法具有更强的跳跃能力;采用指数递减的动态调节惯性权重,以加速局部搜索和提高收敛速度;通过引入最优个体引导策略,增加找到较优位置的概率。实验结果表明,所提算法在小数据集条件下,BIC评分、准确率及汉明距离等指标均优于MCMC算法、BPSO算法、CSO算法、ADLCSO-I算法和SA-ICSO算法。

关键词 鸡群算法;莱维飞行;降维核密度;结构学习

DOI 10.3969/j.issn.2097-1915.2024.02.012

中图分类号 TP181 **文献标志码** A **文章编号** 2097-1915(2024)02-0100-10

A BN Structure Learning Based on DRKDE-ICSO in Small Data Sets

CHEN Haiyang, LIU Jing, LIU Xiqing, ZHANG Jing

(School of Electronics and Information, Xi'an Polytechnic University, Xi'an 710048, China)

Abstract In order to solve the problem of highly similar data in the condition of small data set expansion, the dimensionality reduced kernel density estimation method is utilized for expanding the small data set, obtaining more accurate expanded data. In addition, in order to solve the problems of low efficiency and weak convergence of CSO, an improved ICSO is proposed to learn the structure: Lévy flight is introduced into the position update formula of rooster to make the algorithm jump further; the dynamic adjustment inertia weight with exponential decline is adopted to hasten local search and augmenting convergence speed; by introducing the most advantageous individual guidance approach, the likelihood of discovering the ideal position is increased. The experimental results show that the proposed algorithm is superior to the MCMC algorithm, the BPSO algorithm, the CSO algorithm, the ADLCSO-I algorithm and the SA-ICSO algorithm in terms of BIC score, accuracy and Hamming distance under conditions of small data set.

Key words chicken swarm optimization; Lévy flight; kernel density estimation; structure learn

贝叶斯网络(Bayesian network, BN)的理论研究^[1]包括结构学习、参数学习和推理,其中结构学习是BN理论研究的重要部分。目前,用于BN结构学习的观测数据主要分为充分数据集和小数据集,

收稿日期: 2023-11-29

基金项目: 国家自然科学基金(51905405)

作者简介: 陈海洋(1967—),男,陕西西安人,副教授,博士,研究方向为贝叶斯网络。E-mail: chy_00@163.com

通信作者: 刘静(1998—),女,甘肃金昌人,硕士生,研究方向为贝叶斯网络。E-mail: 2350588223@qq.com

引用格式: 陈海洋,刘静,刘喜庆,等.小数据集下基于DRKDE-ICSO的BN结构学习[J].空军工程大学学报,2024,25(2):100-109. CHEN Haiyang, LIU Jing, LIU Xiqing, et al. A BN Structure Learning Based on DRKDE-ICSO in Small Data Sets[J]. Journal of Air Force Engineering University, 2024, 25(2): 100-109.

在实际应用中,有时会由于客观条件的限制只能获得小数据集,尤其在医疗、军事、航空等领域^[2-5]较为常见,例如,由于数据敏感性和收集难度的原因,军事装备只能获得有限数据。基于充分数据集的 BN 结构学习算法相对成熟^[6-7],而针对小数据集的研究则相对较少。基于小数据集的 BN 结构学习主要有 2 个思路:一是利用约束思想、数据修正等方法对 BN 进行结构学习。戴晶幅等提出了一种双尺度约束模型,用于在缺乏先验信息的情况下学习 BN 结构,显著提高了迭代寻优的收敛速度^[8]。陈海洋等为了提高学习效率,提出基于改进蚁狮优化的 BN 结构学习算法,进而提高了学习性能^[9]。二是对小数据集进行有效拓展,增加可靠的数据信息。Eli-dan 使用了 Bootstrap 抽样来扩充小数据集,并采用 Bagging 方法分析每个数据集得到的 BN 结构,通过实验证明,这种方法能够有效提高不同实际数据集的泛化能力^[10]。王双成等提出了一个结合扩展和修正的小数据集处理机制,成功实现了可靠的基于小数据集的 BN 结构学习^[11]。在现有的数据拓展方法中,只能利用初始数据集进行拓展。因此,如果数据集较小,会出现新生成数据与初始数据过于相似的情况,无法有效改善信息缺失问题。对此,本文提出了基于降维核密度估计的小数据集拓展方法 (dimensionality reduction kernel density estimation, DRKDE),该方法使用降维思想和核密度估计相结合,能够在小数据集下生成与原始数据分布密度相似的拓展数据,从而形成更加充分的数据集,在此基础上需要选择合适的搜索算法进行 BN 结构学习。

目前,在众多搜索算法中,群体智能算法^[12]具备灵活性、鲁棒性、自组织性等特征,在解决寻优问题领域具有显著优势。然而,通过降维核密度拓展数据后,数据集的有效信息增多,加大了 BN 结构学习的计算难度。与其他群体智能算法相比,经典鸡群优化算法 (chicken swarm optimization, CSO) 通过模拟鸡群的层次结构和觅食行为来实现寻优,因此具有相对优秀的全局探索能力和局部优化能力,使得其在处理拓展后的数据时更为可靠有效,但仍存在求解效率低下和收敛性不足的问题。对此,本文提出了改进的鸡群优化算法 (improved chicken swarm optimization, IC-SO),在传统鸡群优化算法的觅食行为中加入动态

调节惯性权重和向最优个体学习行为,增加最优个体产生的概率,减少盲目搜索,同时,在雄鸡的位置更新策略中引入莱维飞行,改进算法的搜索策略,提高算法的收敛精度。

1 相关概念

1.1 贝叶斯网络

1 个 BN 由结构 G 和参数 E 2 部分组成。BN 的基本结构为 1 个有向无环图 $G=(V,E)$,其中, $V=\{X_1,X_2,\dots,X_n\}$ 表示了领域变量的集合,有向边集合 E 用来描述变量之间的相互依赖关系。贝叶斯网络利用条件独立性来分解联合概率分布,它将联合概率分布 P 分解为:

$$P(X_1,X_2,\dots,X_n)=\prod_{i=1}^n P(X_i|\pi(X_i)) \quad (1)$$

式中: X_i 为网络中的第 i 个节点变量; $\pi(X_i)$ 为节点 X_i 的父节点变量集合,节点变量 X_i 可取离散值或者连续值。

1.2 核密度估计法

核密度估计法 (kernel density estimation, KDE) 不需要依赖于整体的分布假设,由已知样本去估计总体的概率分布密度^[13]。Parzen 给出了核密度估计的一般定义,设 X_1,X_2,\dots,X_n 是源自一元连续总体的独立分布样本,在点 $x(x \in R)$ 处总体概率密度函数 $f(x)$ 的核密度估计为:

$$\hat{f}_h(x)=\frac{1}{nh}\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) \quad (2)$$

式中: $K(\cdot)$ 称为核函数; n 为独立同分布随机变量; h 为窗宽或者带宽,常用均方积分误差法来求最佳窗宽。

1.3 鸡群算法

鸡群优化算法^[14]通过模拟鸡群觅食规则与鸡群等级制度来更新个体位置,具有全局搜索能力且收敛精度较高。其遵循以下规则:①根据公鸡的数量划分鸡群为多个子群,每个子群包括 1 只公鸡、若干只母鸡和小鸡;②根据适应性指标的对比,将鸡群划分为公鸡、母鸡和小鸡 3 个群体,适应性指标较高的为公鸡,较低的为小鸡,其余为母鸡;③鸡群内的等级制度是固定的,每隔一定代数会更新 1 次;④所有母鸡均在公鸡的带领下捕食,而小鸡则在母鸡的引导下捕食;⑤在更新位置时,各个等级的鸡会遵循不同的迭代策略。

1.3.1 雄鸡位置更新策略

雄鸡的位置更新公式为:

$$x_{ij}^{t+1} = x_{ij}^t [1 + \text{rand} n(0, \sigma^2)] \quad (3)$$

$$\sigma^2 = \begin{cases} 1, & f_i \geq f_k \\ \exp\left(\frac{f_i - f_k}{|f_i| + \epsilon}\right), & f_i < f_k \end{cases} \quad (4)$$

式中: $\text{rand} n(0, \sigma^2)$ 表示均值为 0、方差为 σ^2 的正态分布随机数; k 表示随机选择雄鸡组里的 1 只雄鸡, 且 $k \neq i$; f_k 、 f_i 分别表示第 k 、 i 只鸡的适应度值; ϵ 表示非常小的常数。

1.3.2 雌鸡位置更新策略

雌鸡的位置更新公式为:

$$x_{ij}^{t+1} = x_{ij}^t + c_1 \text{rand}(x_{ij}^t - x_{ij}^s) + c_2 \text{rand}(x_{ij}^t - x_{ij}^r) \quad (5)$$

式中: rand 为 $[0, 1]$ 区间均匀分布的随机数; r 为第 i 只雌鸡对应的配偶雄鸡; s 为鸡群中任意 1 只雄鸡或雌鸡, 且 $r \neq s$; c_1 和 c_2 分别为向第 r 和 s 只鸡靠拢的学习参数, 其公式分别为:

$$c_1 = \exp((f_i - f_r) / (|f_i| + \epsilon)) \quad (6)$$

$$c_2 = \exp(f_s - f_i) \quad (7)$$

1.3.3 雏鸡位置更新策略

雏鸡的位置更新公式为:

$$x_{ij}^{t+1} = x_{ij}^t + F(x_{mj}^t - x_{ij}^t) \quad (8)$$

式中: m 对应的是第 i 只雏鸡的母亲; x_{mj}^t 表示在第 t 次迭代过程中, 第 m 只雌鸡第 j 维位置; F 为雏鸡跟随母亲觅食的概率。

2 DRKDE-ICSO 算法构建

2.1 算法思想

本文研究的是小数据集条件下的 BN 结构学习问题。首先, 针对在小数据集条件下, 拓展数据集可能存在与初始数据集极其相似的问题, 提出降维核密度估计法来拓展数据集, 并与初始数据集相结合组成充分数据集; 其次, 针对经典 CSO 位置更新方法有效性较低, 导致算法整体搜索能力下降的问题, 提出改进的鸡群优化算法; 最后, 利用充分数据集, 使用改进的鸡群优化算法学习 BN 结构。

2.2 基于 DRKDE 的小数据集拓展

离散贝叶斯网络得到的观测数据为多维数据组, 每个观测节点能得到 1 个观测数据值, 各个节点之间存在关联或独立关系。观测节点越多, 数据维

度越大, 数据之间关系越复杂, 拓展难度也越大, 但离散贝叶斯网络每个节点的状态个数是固定的, 因此, 可以将每组高维数组映射到低维数组中, 降低计算复杂度。降维方法主要思想是将高维数组一对一地对应到低维数组里, 常用的降维方法包括二进制、四进制或八进制转换为十进制, 方法的选取依赖于 BN 节点的状态个数。由于本文所选用的 2 个标准网络中各节点均为 2 种状态, 分别用 0 和 1 表示, 因此选用二进制转换方法。以经典的草坪湿润贝叶斯网络为例, 此网络含有 4 个节点, 每个节点存在 2 个状态, 分别用 0、1 表示, 若某观测数据为 $[0 \ 1 \ 0 \ 0]$, 则可用二进制转换为十进制的方法将其映射为 4。用同样的方法将观测数据都映射到一维数组里, 然后利用概率密度核估计法进行数据拓展。

2.3 基于 ICSO 的 BN 结构学习算法

2.3.1 基于莱维飞行的雄鸡位置更新策略

在整个种群中公鸡是主导者, 他们的地理分布趋向于达到最佳状态。然而, 当前公鸡的位置更新策略可能使得算法的集成进程减慢并具有一定的盲目性。为解决该问题, 在雄鸡位置的更新公式中引入莱维飞行, 用以增加 CSO 算法中公鸡的搜索广度。

莱维飞行^[15]是指一种特定的飞行方式, 其主要优势在于在有限空间内进行高效搜索并且能够完成远距离的移动。莱维飞行的位置更新公式为:

$$x_i^{t+1} = x_i^t + \vartheta \oplus \text{Lévy}(\lambda), i = 1, 2, \dots, N \quad (9)$$

式中: x_i^t 为 x_i 第 t 次迭代的位置; \oplus 为点对点乘法; ϑ 为步长控制量; $\text{Lévy}(\lambda)$ 为随机搜索路径, 并且满足:

$$\text{Lévy}(\lambda) \sim u = t^{-\lambda}, 1 < \lambda \leq 3 \quad (10)$$

式中: u 为正态分布随机数。

在具体操作过程中, Mantegna 算法经常被用于模拟莱维飞行, 相关数学表达式为:

$$s = \frac{u}{|v|^{\frac{1}{\beta}}} \quad (11)$$

$$u \sim N(0, \sigma_u^2) \quad (12)$$

$$v \sim N(0, \sigma_v^2) \quad (13)$$

$$\sigma_u = \left\{ \frac{\Gamma(1+\beta) \sin \frac{\pi\beta}{2}}{\Gamma(\frac{1+\beta}{2}) \beta 2^{\frac{\beta-1}{2}}} \right\}^{\frac{1}{\beta}} \quad (14)$$

式中: s 为随机步长; Γ 为伽玛函数; $\sigma_u = 0.696 \ 6$, $\sigma_v = 1$; β 通常取值为 1.5。

改进的雄鸡位置更新公式为:

$$x_{ij}^{t+1} = x_{ij}^t (1+s) \quad (15)$$

式中: s 为随机步长。

2.3.2 动态调节的惯性权重

CSO 算法迭代初期能在较大范围内快速搜索,使个体位置发生显著变化。然而,在迭代后期,由于搜索范围逐渐变小,适应度值也随之缓慢变化,易陷入局部极值。因此,受非线性递减权重策略及动态权重策略的启发,提出一种指数递减的动态调节惯性权重策略。

动态调节的惯性权重公式为:

$$\omega = \omega_{\text{end}} (\omega_{\text{start}} / \omega_{\text{end}})^{(1/(1+at/T))} \text{normrnd}(1,0.1) \quad (16)$$

式中: ω_{start} 和 ω_{end} 分别是迭代开始时和迭代结束时的 ω 取值,针对不同的问题,其取值也会不同,一般在 $0.9 \sim 0.4$ 内取值; t 为当前迭代次数; T 为最大迭代次数; $\text{normrnd}(1,0.1)$ 表示服从均值为 1、标准差为 0.1 的正态分布随机数; a 为影响调节递减效果的参数,本文 a 取 $15^{[16]}$ 。

2.3.3 最优个体引导策略

在基本 CSO 算法中,公鸡自主觅食,母鸡围绕公鸡觅食,小鸡围绕母鸡觅食。公鸡搜索到局部最优解时,母鸡和小鸡也会紧随其后陷入局部最优。在鸡群位置更新策略中加入最优个体引导策略,可提高找到较优位置的概率,有效避免个体陷入局部最优。

改进的雄鸡位置更新公式为:

$$x_{ij}^{t+1} = x_{ij}^t (\omega + s) + d_1 \text{rand}(x_{\text{best}j}^t - x_{ij}^t) \quad (17)$$

式中:各参数定义与式(3)和式(4)相同; ω 为动态惯性权重; $x_{\text{best}j}^t$ 为当前代数的最优个体; d_1 为雄鸡向最优个体学习的参数。

改进的雌鸡位置更新公式为:

$$x_{ij}^{t+1} = \omega x_{ij}^t + c_1 \text{rand}(x_{ij}^t - x_{ij}^t) + c_2 \text{rand}(x_{sj}^t - x_{ij}^t) + d_2 \text{rand}(x_{\text{best}j}^t - x_{ij}^t) \quad (18)$$

式中:各参数定义与式(5)相同; d_2 为雌鸡向最优个体学习的参数。

改进的雏鸡位置更新公式为:

$$x_{ij}^{t+1} = \omega x_{ij}^t + F(x_{mj}^t - x_{ij}^t) + d_3 \text{rand}(x_{\text{best}j}^t - x_{ij}^t) \quad (19)$$

式中:各参数定义与式(8)相同; d_3 为雏鸡向最优个体学习的参数。

为了确定参数 d_1 、 d_2 、 d_3 的取值,利用控制变量法进行了对比实验,随机选择其中 2 个参数为定值,改变另 1 个参数的值。通过迭代寻优的方式评估算法性能,如图 1 所示。当 d_1 、 d_2 、 d_3 的取值分别为 2、

1.8、1.6 时,算法能够在最少的迭代次数内找到最优解,因此本文 d_1 、 d_2 、 d_3 分别取 2、1.8、1.6。

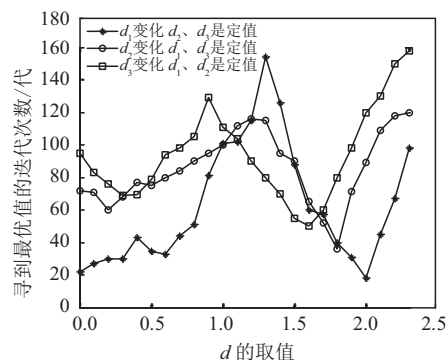


图 1 最优个体引导策略对比图

3 算法实现步骤

本文算法实现的基本步骤为:

步骤 1 将高维初始数据映射到低维空间上,确定核函数与最优窗宽 h ,利用核密度估计法求出总体密度函数 $f(x)$ 。

步骤 2 随机生成概率密度函数为 $f(x)$ 的拓展数据集,与初始数据集组成充分数据集。

步骤 3 初始化改进鸡群优化算法各参数,鸡群种群数量设为 N ,随机生成 N 个初始贝叶斯网络结构,当前进化代数 $t=0$,最大迭代次数为 M 。

步骤 4 结合充分数据集,将 BIC 评分作为适应度指标,确定初代个体最佳位置和全局最优位置。

步骤 5 根据 BIC 评分值进行排序,建立鸡群等级制度。

步骤 6 利用鸡群位置更新策略更新位置信息,处理随机生成的无效结构,更新局部最优位置和全局最优位置。

步骤 7 判断是否满足迭代次数,若满足,输出最优结构,若不满足,则转入步骤 5。

4 实验结果分析

实验分 4 部分,仿真环境为:Windows 10,64 bit,CPU 为 2.5 GHz,程序实现采用 Matlab R2014a 以及软件工具包 FullBNT-1.0.4。

4.1 小数据集拓展效果验证

为验证概率密度核估计效果,使用验证 BN 结构学习领域的经典网络五节点过劳死网络和八节点胸部疾病诊断网络进行仿真实验。

利用 2 个贝叶斯网络随机生成 50 组初始观测数据,映射到一维数组后,再生成容量为初始数据集 10 倍的拓展数据集,其初始给定与拓展后的核密度估计结果对比图分别如图 2 和图 3 所示。可以看出,在五节点和八节点网络中都存在 1 个数据量点,当数据量小于该点时,拓展误差较大,大于该点时,核密度估计曲线平滑完整,总体密度相似且误差较小,该点在图 2 和图 3 中分别为 1.384 和 5.269。因此,总体核密度估计拓展效果理想,但初始数据量不易过小。

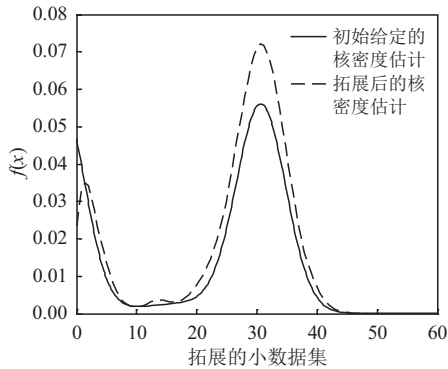


图 2 五节点核密度估计对比图

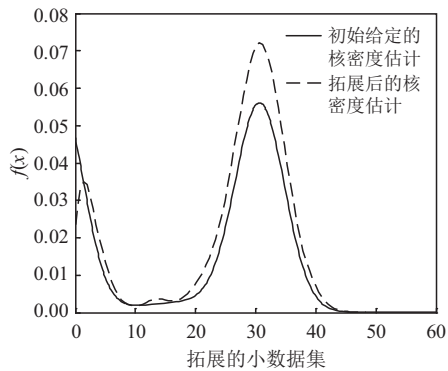


图 3 八节点核密度估计对比图

为了验证拓展数据是否可以有效建立 BN 结构,利用 2 个网络分别产生容量为 50、100、200、500 的初始数据,利用降维核密度估计法分别扩展 10 倍后作为实验数据,再分别利用 CSO 算法和 BPSO 算法进行 BN 结构学习,在初始小数据集和拓展小数据集下进行汉明距离的比较,图 4 和图 5 分别展示了在 2 个网络中各算法的比较结果。

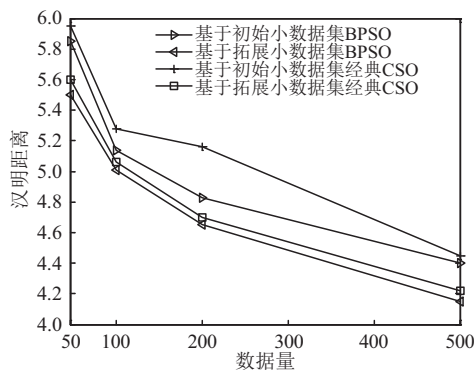


图 4 五节点网络各算法汉明距离

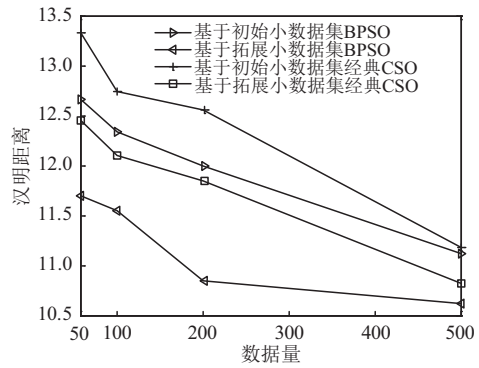


图 5 八节点网络各算法汉明距离

由图 4 和图 5 可知,在过劳死贝叶斯网络与胸部疾病诊断网络中,采用拓展小数据集后,各算法的汉明距离均优于小数据集下的汉明距离。这是因为利用降维核密度估计法对小数据集进行拓展,有效利用了有限数据,降低了汉明距离,从而提高了结构学习精度。

4.2 测试函数及仿真实验

为了验证 ICSO 算法的有效性,本文采用 5 个不同的测试函数进行测试和验证,并将结果与 PSO 算法、DE 算法、CSO 算法、ADLCSO-I 算法^[17]和 SA-ICSO 算法^[18]进行对比。测试函数如下:

1) Ackley 函数

$$f_1(x) = 20 + e - 20 \exp\left(-\frac{1}{5} \sqrt{\frac{1}{n} \sum_{i=1}^n x_i}\right) - \exp\left(\frac{1}{n} \sum_{i=1}^n \cos(2\pi x_i)\right), x_i \in [-32, 32] \quad (20)$$

2) Griewank 函数

$$f_2(x) = \frac{1}{4000} \sum_{i=1}^n (x_i)^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right) + 1, x_i \in [-600, 600] \quad (21)$$

3) Rastrigin 函数

$$f_3(x) = \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i) + 10), x_i \in [-5.12, 5.12] \quad (22)$$

4) Sphere 函数

$$f_4(x) = \sum_{i=1}^n x_i^2, x_i \in [-5.12, 5.12] \quad (23)$$

5) Zakharov 函数

$$f_5(x) = \sum_{i=1}^n x_i^2 + \left(\sum_{i=1}^n 0.5 i x_i\right)^2 + \left(\sum_{i=1}^n 0.5 i x_i\right)^4, x_i \in [-5, 10] \quad (24)$$

各算法的种群规模设定为 100,维数设定为 20,最高迭代次数设定为 1 000。每组测试仿真 50 次取平均值,测试结果如表 1 所示。

表 1 函数测试结果

函数	算法	最好值	均值	最差值	标准差
$f_1(x)$	PSO 算法	2.184 5e-06	8.786 7e-06	1.886 9e-05	2.413 6e-05
	DE 算法	3.498 1e-08	4.014 4e-07	4.719 1e-07	3.613 3e-07
	CSO 算法	4.440 9e-15	4.440 9e-15	4.440 9e-15	0
	ICSO 算法	8.881 8e-16	8.881 8e-16	8.881 8e-16	0
	ADLCSO-I 算法	8.881 8e-16	8.881 8e-16	8.881 8e-16	0
	SA-ICSO 算法	8.881 8e-16	8.881 8e-16	8.881 8e-16	0
$f_2(x)$	PSO 算法	2.760 3e-15	2.729 1e-11	3.999 6e-11	2.259 3e-11
	DE 算法	3.685 9e-14	2.790 0e-11	1.248 9e-10	1.292 2e-10
	CSO 算法	0	0	0	0
	ICSO 算法	0	0	0	0
	ADLCSO-I 算法	0	0	0	0
	SA-ICSO 算法	0	0	0	0
$f_3(x)$	PSO 算法	2.771 1e-13	2.323 5e-10	1.020 9e-08	1.524 3e-08
	DE 算法	2.812 3e-08	1.786 6e-07	1.571 0e-06	1.248 5e-10
	CSO 算法	0	0	0	0
	ICSO 算法	0	0	0	0
	ADLCSO-I 算法	0	0	0	0
	SA-ICSO 算法	0	0	0	0
$f_4(x)$	PSO 算法	1.343 2e-17	3.920 0e-11	2.025 3e-10	1.875 4e-10
	DE 算法	7.454 2e-16	1.028 0e-15	2.551 0e-15	2.099 2e-15
	CSO 算法	4.457 6e-80	3.493 4e-76	7.871 1e-76	7.402 9e-76
	ICSO 算法	0	0	0	0
	ADLCSO-I 算法	4.743 5e-80	3.756 4e-78	3.675 3e-78	2.958e-79
	SA-ICSO 算法	2.572 5e-100	3.246 54e-96	2.477 5e-95	8.121 14e-96
$f_5(x)$	PSO 算法	5.213 9e-15	3.902 8e-12	4.923 2e-11	5.428 0e-11
	DE 算法	36.163 3	59.004 3	73.004 7	27.682 4
	CSO 算法	4.808 6e-11	2.034 9e-07	7.089 1e-07	7.691 2e-07
	ICSO 算法	0	0	0	0
	ADLCSO-I 算法	1.783 94e-112	4.676 54e-112	4.453 65e-102	4.654 6e-112
	SA-ICSO 算法	2.375 9e-41	7.407 07e-39	5.580 4e-38	1.825 45e-38

(注:加粗字体为每组最优值)

由表 1 可知,在函数 $f_1(x)$ 中,ADLCSO-I 算法、SA-ICSO 算法和 ICSO 算法都找到了最优值;在函数 $f_2(x)$ 和 $f_3(x)$ 中,除了 PSO 算法和 DE 算法外,其他算法均能在规定的迭代次数内找到理论最优值;在函数 $f_4(x)$ 和 $f_5(x)$ 中,只有 ICSO 算法能够找到理论最优值。总体来看,ICSO 算法在测试函数中表现最优。

为了更准确评估各算法性能,从统计学角度出

发,采用 Friedman 检验^[19]。实验中,对八节点网络采样提取了 5 个数据集,分别应用 PSO 算法、MCMC 算法、CSO 算法、ADLCSO-I 算法、SA-ICSO 算法和 ICSO 算法对目标进行寻优。通过这种方法,得到各算法排名,所得数据是基于 10 次实验的平均值,以获得一个更稳定和可靠的评价结果,在数值后面用括号标记该算法在所有算法中的排名,具体评估结果如表 2 所示。

表 2 Friedman 检验结果

数据集	PSO 算法	MCMC 算法	CSO 算法	ADLCSO-I 算法	SA-ICSO 算法	ICSO 算法
50	0.591 5(3)	0.372 9(6)	0.538 0(5)	0.543 0(4)	0.667 1(1)	0.599 8(2)
100	0.562 7(3)	0.421 7(6)	0.474 1(5)	0.526 9(4)	0.564 1(2)	0.565 2(1)
200	0.528 1(3)	0.469 2(6)	0.496 9(5)	0.513 4(4)	0.533 1(2)	0.661 2(1)
500	0.579 7(2)	0.490 6(5)	0.516 2(4)	0.447 1(6)	0.559 3(3)	0.602 1(1)
1 000	0.597 1(4)	0.386 0(6)	0.660 9(3)	0.542 6(5)	0.688 2(2)	0.723 9(1)

由表 2 可知,在算法寻优性能方面,ICSO 算法在多数数据集上表现出色,且随着数据量的增加,优势越发明显,相较于其他算法具有更显著的优势。

4.3 算法复杂度分析

在对算法的复杂度进行分析时,主要关注耗时和占用空间最多的步骤。鉴于本文算法与其他 5 种算法相比,在空间复杂度上没有明显的差异,因此主要关注各算法的时间复杂度。

表 3 3 种算法的时间复杂度分析

PSO 算法步骤	复杂性	DE 算法步骤	复杂性	ICSO 算法步骤	复杂性
初始化粒子位置	ND	初始化位置	ND	初始化位置	ND
初始化粒子速度	ND	计算适应度值	ND	计算适应度值	ND
计算适应度值	N	更新最好位置	N	更新最好位置	N
更新最好位置	N			位置更新	ND
更新速度	ND			计算适应度值	ND
计算新位置	ND				

由表 3 可以看出,PSO 算法、DE 算法和 ICSO 算法的时间复杂度分别为 $2N+4ND$ 、 $N+2ND$ 和 $N+4ND$ 。显然,ICSO 算法与 PSO 算法的时间复杂度相当接近,DE 算法的时间复杂度优于其他 2 种算法。但仅考虑时间复杂度无法全面评价算法的优越性,总体来看 ICSO 算法表现更优越。

4.4 BN 结构学习仿真实验验证

为验证本文算法结构学习的寻优准确性和全局收敛性,利用 4.1 节提到的标准过劳死网络和胸部

群体智能优化算法的时间复杂度与问题优化的规模 D 及群体内个体的数量 N 成正比关系。由于 ICSO 算法、CSO 算法、ADLCSO-I 算法和 SA-ICSO 算法在一次迭代所需要的主要步骤和相应时间复杂度是相同的,因此这里只列出了 PSO 算法、DE 算法和 ICSO 算法在一次迭代中所需的主要步骤及对应的时间复杂度,如表 3 所示。

疾病诊断网络进行仿真实验,2 个网络分别产生容量为 50、100、200、500 的初始数据,利用降维核密度估计法分别扩展 10 倍后作为实验数据。将本文算法与 BPSO 算法^[20]、CSO 算法、ADLCSO-I 算法、SA-ICSO 算法和 MCMC 算法进行比较,为降低数据随机性的影响,每组仿真结果为 50 次的平均值。

针对不同小数据集拓展的充分数据集,分别在五节点和八节点网络上运行不同结构学习算法的仿真结果如表 4 和表 5 所示。

表 4 五节点网络 6 种算法仿真结果对比

初始数据组	算法	BIC 评分	相同边	运行时间/s	准确率/%	汉明距离
50	MCMC 算法	-1 753.38	18.19	1.39	72.76	6.81
	BPSO 算法	-1 753.57	19.12	5.11	77.66	5.50
	CSO 算法	-1 756.38	18.34	10.00	76.61	5.60
	ADLCSO-I 算法	-1 755.67	19.35	11.55	77.40	5.65
	SA-ICSO 算法	-1 745.16	19.86	9.56	79.44	5.14
	ICSO 算法	-1 752.00	21.22	9.90	79.71	5.40
100	MCMC 算法	-3 417.50	18.99	1.52	75.96	6.00
	BPSO 优化算法	-3 419.15	19.86	5.44	79.85	5.01
	CSO 算法	-3 422.32	19.06	10.79	79.02	5.06
	ADLCSO-I 算法	-3 421.87	19.80	11.89	79.20	5.20
	SA-ICSO 算法	-3 414.73	20.21	10.27	80.84	4.79
	ICSO 算法	-3 418.06	20.38	10.81	81.26	4.70
200	MCMC 算法	-6 674.96	19.72	1.77	78.88	5.28
	BPSO 算法	-6 676.07	20.53	5.79	81.53	4.65
	CSO 算法	-6 679.41	19.62	11.67	80.67	4.70
	ADLCSO-I 算法	-6 678.09	20.15	12.13	80.60	4.85
	SA-ICSO 算法	-6 670.09	20.40	10.81	81.60	4.60
	ICSO 算法	-6 674.60	20.73	11.81	82.92	4.55
500	MCMC 算法	-16 411.24	19.92	2.06	79.68	5.08
	BPSO 算法	-16 416.67	20.85	7.59	83.40	4.15
	CSO 算法	-16 406.21	20.78	15.12	83.12	4.22
	ADLCSO-I 算法	-16 415.34	20.95	15.96	83.80	4.05
	SA-ICSO 算法	-16 404.67	21.01	14.80	84.04	4.00
	ICSO 算法	-16 403.79	21.21	14.78	84.20	3.98

(注:加粗字体为每组最优值)

表 5 八节点网络 6 种算法仿真结果对比

初始数据集	算法	BIC 评分	相同边	运行时间/s	准确率/%	汉明距离
50	MCMC 算法	-2 166.02	51.21	2.65	80.29	12.57
	BPSO 算法	-2 168.01	52.15	9.01	81.68	11.70
	CSO 算法	-2 179.41	51.55	17.48	80.55	12.45
	ADLCSO-I 算法	-2 175.30	53.20	17.96	83.13	10.80
	SA-ICSO 算法	-2 166.91	52.89	18.71	82.64	11.11
	ICSO 算法	-2 160.19	53.44	17.83	84.24	10.00
100	MCMC 算法	-4 311.60	51.40	2.71	80.87	12.16
	BPSO 算法	-4 316.20	53.15	9.95	82.15	11.55
	CSO 算法	-4 336.84	52.95	18.23	81.40	12.10
	ADLCSO-I 算法	-4 324.88	53.97	18.98	84.33	10.03
	SA-ICSO 算法	-4 306.79	53.78	20.35	84.03	10.22
	ICSO 算法	-4 305.48	54.21	18.60	84.89	9.65
200	MCMC 算法	-8 525.50	52.75	3.14	81.50	11.97
	BPSO 算法	-8 527.46	54.20	10.53	83.32	10.85
	CSO 算法	-8 548.87	53.67	20.22	81.91	11.85
	ADLCSO-I 算法	-8 544.19	54.40	21.13	85.00	9.60
	SA-ICSO 算法	-8 523.24	54.22	24.56	84.72	9.78
	ICSO 算法	-8 511.43	54.84	20.61	85.43	9.35
500	MCMC 算法	-21 836.26	53.18	4.21	83.09	10.82
	BPSO 算法	-21 840.44	55.68	14.22	83.98	10.62
	CSO 算法	-21 872.67	55.14	28.10	83.59	10.82
	ADLCSO-I 算法	-21 801.45	54.60	29.54	85.31	9.40
	SA-ICSO 算法	-21 786.43	54.67	31.45	85.42	9.33
	ICSO 算法	-21 790.24	56.21	28.41	85.79	9.31

(注:加粗字体为每组最优值)

分析表 4 和表 5 数据可知,随着数据量的增加,各算法所共同拥有的边数量也会增加,这表明有效数据信息的增多可以提高所学到的贝叶斯网络的准确性,从而进一步验证了拓展小数据集的有效性。从 BIC 评分来看,ICSO 算法的 BIC 评分明显优于其他 4 种算法,由此可以说明,在拓展小数据集下,ICSO 算法对于结构学习寻优的效果最佳。为了更直观地比较不同数据集下不同结构学习算法仿真结果的变化趋势,对各算法按照运行时间、准确率和汉明距离进行详细对比,对比结果分别如图 6~图 11 所示。

由图 6 和图 7 可知,ICSO 算法运行时间高于 BPSO 算法和 MCMC 算法且与 CSO 算法、ADLCSO-I 算法、SA-ICSO 算法处于相同数量级。但由图 8 和图 9 可知,ICSO 算法在 2 个网络中的准确率均高于其他 5 种算法,这归因于 ICSO 算法采用了一系列优化措施。由图 10 和图 11 可知,除了在五节点网络中,当初始数据量小于 100 时,ICSO 算法的汉明距离与 SA-ICSO 算法相比略差外,仍优于其他算法。然而,随着数据量的增加,ICSO 算法的汉明距离逐渐优于 SA-ICSO 算法。此外,在五节点网络中,当初始数据量小于 500 时,ADLCSO-I 算法的准确率和汉明距离略低于 BPSO 算法和 CSO 算法,

由此说明,在节点数目较小的网络中 ADLCSO-I 算法没有太大优势。

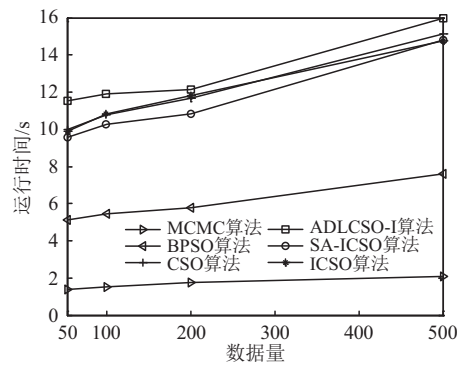


图 6 五节点网络中各算法时间对比

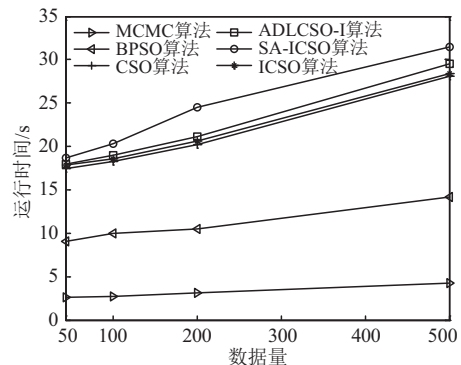


图 7 八节点网络中各算法时间对比图

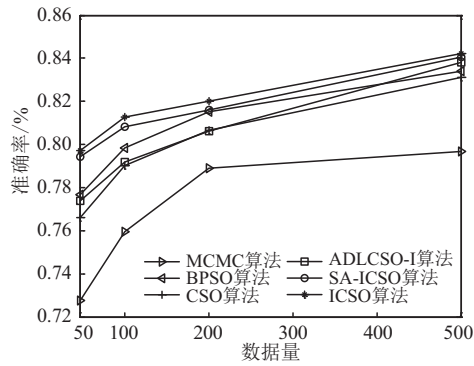


图8 五节点网络中各算法的准确率

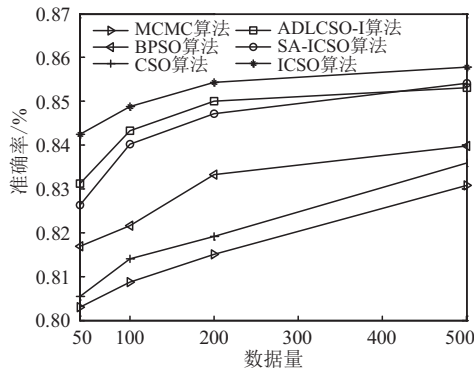


图9 八节点网络中各算法的准确率

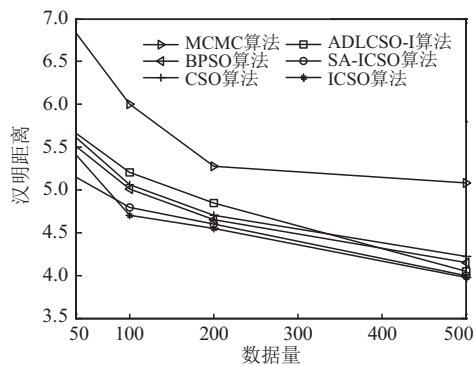


图10 五节点网络中各算法的汉明距离

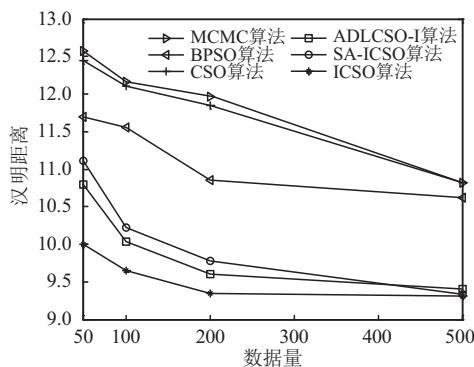


图11 八节点网络中各算法的汉明距离

综上所述,在拓展小数据集条件下,ICSO 算法在寻优过程中,无论在寻优效率还是寻优精度方面,均表现出优于其他算法。这表明该算法在解决小数据集优化问题时,具有较高准确性。

5 结语

本文首先提出了基于降维核密度估计的小数据集拓展方法,其次提出了改进的鸡群优化算法进行贝叶斯网络结构寻优。仿真实验结果表明,降维核密度估计法拓展小数据集效果理想;改进的鸡群优化算法具有较高的寻优能力,优于 MCMC 算法、BPSO 算法、CSO 算法、ADLCSO-I 算法和 SA-ICSO 算法。实验中,本文提出的动态惯性权重、最优个体引导策略以及莱维飞行的引入,均对鸡群优化算法寻优能力的提升起到了重要作用,体现了本文算法的优越性。

参考文献

- [1] WANG J Y, LIU S Y. Novel Binary Encoding Water Cycle Algorithm for Solving Bayesian Network Structures Learning Problem [J]. Knowledge-Based Systems, 2018, 150: 95-110.
- [2] REFAI A, MEROUANI H F, AOURAS H. Maintenance of a Bayesian Network: Application Using Medical Diagnosis [J]. Evolving Systems, 2016, 7 (3): 187-196.
- [3] MCLACHLAN S, DUBE K, HITMAN G A, et al. Bayesian Networks in Healthcare: Distribution by Medical Condition [J]. Artificial Intelligence in Medicine, 2020, 107: 101912.
- [4] 翟贵敏,董龙明,邱瑞波,等.基于贝叶斯网络的空中目标威胁估计算法[J].火力与指挥控制,2016,41(11):90-93,97.
- [5] 高天祥,王刚,岳韶华,等.基于贝叶斯决策理论的NSHV分段建模威胁评估[J].空军工程大学学报(自然科学版),2019,20(1):60-66.
- [6] 王守会,覃枫.基于集成学习和反馈策略的贝叶斯网络结构学习[J].计算机学报,2021,44(6):1051-1063.
- [7] SCANAGATTA M, CORANI G, CAMPOS C P D, et al. Approximate Structure Learning for Large Bayesian Networks [J]. Machine Learning, 2018, 107 (8-10):1209-1227.
- [8] 戴晶帆,任佳,董超,等.基于双尺度约束模型的BN结构自适应学习算法[J].自动化学报,2021,47(8):1988-2001.
- [9] 陈海洋,尚珊珊,任智芳,等.基于改进蚁狮优化的贝叶斯网络结构学习算法[J].空军工程大学学报,2023,24(2):104-111.
- [10] ELIDAN G. Bagged Structure Learning of Bayesian Network [J]. Journal of Machine Learning Research, 2011,15:251-258.

