

Enhancing Missing Data Imputation through Combined Bipartite Graph and Complete Directed Graph

Zhaoyang Zhang¹, Hongtu Zhu², Ziqi Chen^{1*}, Yingjie Zhang¹,
Hai Shu³

¹School of Statistics, East China Normal University, Shanghai, China.

²Departments of Biostatistics, Statistics, Computer Science and Genetics,
The University of North Carolina at Chapel Hill, Chapel Hill, USA.

³Department of Biostatistics, School of Global Public Health, New York
University, New York, USA.

*Corresponding author(s). E-mail(s): zqchen@fem.ecnu.edu.cn;

Contributing authors: 10215001426@stu.ecnu.edu.cn;

htzhu@email.unc.edu; 2427785647@qq.com; hai.shu@nyu.edu;

Abstract

In this paper, we aim to address a significant challenge in the field of missing data imputation: identifying and leveraging the interdependencies among features to enhance missing data imputation for tabular data. We introduce a novel framework named the Bipartite and Complete Directed Graph Neural Network (BCGNN). Within BCGNN, observations and features are differentiated as two distinct node types, and the values of observed features are converted into attributed edges linking them. The bipartite segment of our framework inductively learns embedding representations for nodes, efficiently utilizing the comprehensive information encapsulated in the attributed edges. In parallel, the complete directed graph segment adeptly outlines and communicates the complex interdependencies among features. When compared to contemporary leading imputation methodologies, BCGNN consistently outperforms them, achieving a noteworthy average reduction of 15% in mean absolute error for feature imputation tasks under different missing mechanisms. Our extensive experimental investigation confirms that an in-depth grasp of the interdependence structure substantially enhances the model's feature embedding ability. We also highlight the model's superior performance in label prediction tasks involving missing data, and its formidable ability to generalize to unseen data points.

Keywords: Graph Neural Network, Missing Data, Missing Mechanism

1 Introduction

Incomplete data is a common occurrence in fields such as clinical research, finance, and economics and survey research [1–5]. Tackling incomplete data for downstream learning tasks requires a comprehensive strategy focused on preserving data integrity and minimizing bias. A prevalent approach is feature imputation, where missing values are estimated using statistical techniques such as regression imputation, as well as more advanced methods, including machine learning algorithms and deep learning models, which capture the underlying relationships in the data [6–10].

When addressing missing data, it is common to classify data structures into two main categories: graph data and tabular data. For graph data, where the connectivity among samples is considered prior knowledge, numerous studies enhance feature imputation based on the known adjacency information among samples, subsequently facilitating downstream tasks such as node classification and link prediction [11–13]. For tabular data, the focus is on the accuracy of the feature imputation and label prediction [9, 14, 15]. Unlike graph data, samples in tabular data are independent of each other. This independence renders most feature imputation methods used in graph data inapplicable to tabular data due to the absence of connecting edges among sample nodes. Nevertheless, some existing methods for feature imputation in tabular data concentrate on capturing the relationships among samples. These methods primarily aim to identify similarities among samples to enhance feature imputation [13, 14]. However, solely focusing on the similarities and connectivity among independent samples in tabular data can be detrimental. For example, Figure 1 (left panel) reveals that IGRM, a variant of GRAPE that considers similarities among samples, performs worse on the Energy dataset compared to the original GRAPE [9, 14]. Moreover, Figure 1 (middle and right panels) illustrates the limitations of IGRM, which only captures similarities among samples when handling missing data.

In this paper, we introduce a cutting-edge framework known as Bipartite and Complete directed Graph Neural Network (BCGNN), which synergizes the structures of bipartite and complete directed graphs to proficiently address the feature imputation task for tabular data. This approach is rooted in inductive learning and a deep comprehension of feature interdependencies.

In BCGNN, the bipartite graph and complete directed graph serve distinct functions. In the bipartite graph, observation nodes and feature nodes aggregate messages from neighboring nodes and corresponding edges to learn their embedding representations. Meanwhile, in the complete directed graph, we innovatively implement an element-wise attention mechanism on the embeddings of feature nodes and introduce the signs of estimated Spearman correlation coefficients to learn and express the interdependence structure among features. The BCGNN framework is constructed by stacking multiple layers of the union of these two subgraphs. At each layer, we perform complete message aggregation and embedding updates on each node, followed

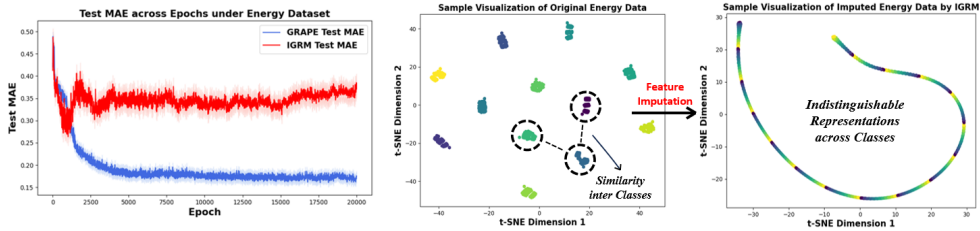


Fig. 1 Comparison of GRAPE and IGRM on the Energy dataset in terms of test MAE in feature imputation (**Left**); Visualizations of the original and the imputed (by IGRM) Energy datasets using t-SNE for dimensionality reduction (**Middle** and **Right**). The poor performance of IGRM can be attributed to erroneously identifying similarities among samples from different classes, leading to indistinguishable sample representations across classes.

by the updates of edge embeddings. Ultimately, the task of feature imputation can be accomplished by integrating the embeddings of corresponding feature nodes and observation nodes at the final layer.

1.1 Related Work

Statistical methods. Classical statistics techniques, such as the Expectation Maximization (EM) algorithm and Multivariate Imputation by Chained Equations (MICE), which rely on distributional assumptions [2, 6, 16–18], have been successful in imputation tasks. Additionally, k-nearest neighbors (KNN) imputation [19] and matrix completion techniques [20–22] have shown effectiveness. However, EM and MICE typically assume a parametric density function for data distribution, which poses limitations when dealing with diverse missing mechanisms and mixed data types. Both MICE and KNN suffer from scalability issues and lack flexibility for downstream tasks. Matrix completion methods also find it challenging to accommodate mixed data types and become computationally intensive for imputing new observations after model training.

Machine and deep learning based methods. Recent deep learning models applied to feature imputation exhibit notable limitations. Generative Adversarial Imputation Nets (GAIN) [23] struggle with missing not at random (MNAR) scenarios. Denoising Autoencoder (DAE) models [24, 25] are limited by their reliance on a single observation, failing to capture complex inter-observation interactions. Optimal Transport (OT) [26] assumes identical distributions for two randomly extracted batches from the same dataset and employs optimal transport distance as the training loss. MIRACLE [10], using causal graphs, faces challenges due to its dependence on restrictive assumptions, limiting its flexibility.

Graph-based framework. Numerous advanced graph-based methods have been proposed for data imputation tasks. Some of these methods are developed for matrix completion and demonstrate excellent performance on large matrices [11, 27–29]. However, their assumptions about feature types limit their applicability to mixed discrete and continuous features. Some Graph Neural Network (GNN)-based models, which employ diversified interaction and information propagation strategies among sample

nodes based on prior known adjacency information, successfully enhance feature imputation in graph data for downstream tasks like node classification and link prediction [11–13]. However, they prove ineffective in tabular data due to the absence of prior topological connectivity among samples. GRAPE [9] showcases success by leveraging a bipartite graph for feature imputation. Following GRAPE, IGRM [14] goes a step further by establishing a friend network among sample nodes to enhance feature imputation by capturing similarities among samples. Meanwhile, GINN [30] utilizes the Euclidean distance of observed data to construct a graph and then reconstructs the tabular data. EGG-GAE [15] learns latent topological structures among sample nodes to assist feature imputation for various downstream tasks.

1.2 Our contributions

It is notable that the current advanced methods for enhancing feature imputation overlook the interdependencies among features. The intricate interdependence among features provides crucial insights for imputing missing data, particularly in scenarios with high missing rates and diverse missing mechanisms. Our work stands as the pioneering effort to capture the interdependencies among features by integrating a complete directed graph into the graph neural network framework. This integration introduces three significant innovations in our BCGNN model:

First, our BCGNN capture the interdependence among features by combining a complete directed graph with a bipartite graph. This integration enables the model to learn node representations inductively and maps out the interdependence between features in a cohesive structure.

Second, our BCGNN employs a unique element-wise additive attention mechanism within the complete directed graph. This allows for a detailed learning of how one feature depends on the attributes of another, a method totally distinct from the usual node-level or edge-level attention in GNNs [31, 32]. By also incorporating the sign of Spearman correlation coefficients, the model gains a refined understanding of feature relationships, enabling more precise imputation.

Lastly, to prevent overfitting, BCGNN uses the DropEdge method [33] across both subgraphs. For enhancing generalization, the framework includes an AttentionDrop feature within the element-wise attention mechanism, preventing undue emphasis on particular segments of the embeddings.

2 Model Architecture and Algorithm

2.1 Problem Definition

Let $\mathbf{D} \in \mathbb{R}^{n \times m}$ be a data matrix composed of n observational samples and m features. The j -th feature of the i -th observation sample is denoted as D_{ij} . Let \mathbf{D}_i be the i -th row of \mathbf{D} . We use a masking matrix $\mathbf{M} = (M_{ij}) \in \{0, 1\}^{n \times m}$ to identify the missing data in data matrix. Specifically, $M_{ij} = 1$ indicates that D_{ij} is observed, while $M_{ij} = 0$ indicates that it is missing. Let $\mathbf{Y} = (Y_i) \in \mathbb{R}^n$ denote the labels for the n samples, we use a corresponding masking vector $\mathbf{M}^y = (M_i^y) \in \{0, 1\}^n$ to identify the missing

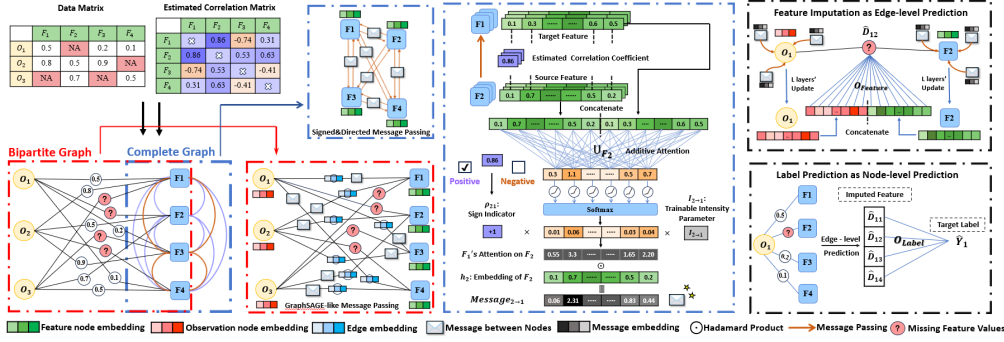


Fig. 2 Flowchart of our BCGNN method. BCGNN consists of a bipartite graph and a complete directed graph constructed from the data matrix and correlation coefficient matrix of features. The working mechanisms of the bipartite graph (**Red Dot-Dashed Box**), the complete directed graph (**Blue Dot-Dashed Box**) and the union of two subgraphs are elaborated in detail in Sections 2.3, 2.4 and 2.5, respectively. With the constructed graph, the feature imputation problem and the label prediction problem are treated as edge-level and node-level prediction tasks, respectively (**Black Dot-Dashed Box**).

values in \mathbf{Y} . We consider feature imputation and label prediction task, i.e, imputing missing feature values D_{ij} at $M_{ij} = 0$ and predicting labels Y_i at $M_i^y = 0$.

2.2 Missing Data Problem as a Union of Two Graphs

Our BCGNN integrates node embeddings, edge embeddings, and a specially designed message-passing mechanism, effectively tackling challenges in missing data imputation. Specifically, BCGNN transforms a data matrix with missing data into a union of two subgraphs as illustrated in Figure 2. One subgraph, denoted as a bipartite graph \mathcal{G}_B , consists of observation nodes, feature nodes, and the edges connecting these two types of nodes. Meanwhile, the other subgraph \mathcal{G}_C is a complete directed graph formed exclusively among the feature nodes. The union graph can then be represented as $\mathcal{G} = \mathcal{G}_B \cup \mathcal{G}_C$. Consequently, the data matrix \mathbf{D} and the masking matrix \mathbf{M} can be naturally represented as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} includes two types of nodes $\mathcal{V} = \mathcal{V}_F \cup \mathcal{V}_O$, with $\mathcal{V}_O = \{u_1, u_2, \dots, u_n\}$ representing observation nodes and $\mathcal{V}_F = \{v_1, v_2, \dots, v_m\}$ representing feature nodes. \mathcal{E} contains two types of edges $\mathcal{E} = \mathcal{E}_{OF} \cup \mathcal{E}_{FF}$, where $\mathcal{E}_{OF} = \{e_{u_i, v_j} | u_i \in \mathcal{V}_O, v_j \in \mathcal{V}_F, M_{ij} = 1\}$ represents edges connecting feature nodes with observation nodes, and $\mathcal{E}_{FF} = \{(e_{v_i \rightarrow v_j}, e_{v_j \rightarrow v_i}) | v_i, v_j \in \mathcal{V}_F, v_i \neq v_j\}$ denotes a set of directed edges connecting each pair of feature nodes. To simplify the notation e_{u_i, v_j} , we use e_{ij} in the context of feature matrix \mathbf{D} , and e_{uv} in the context of graph \mathcal{G} . In \mathcal{E}_{OF} , if D_{ij} is a categorical variable, e_{ij} is initialized with the one-hot vector corresponding to D_{ij} . On the other hand, if D_{ij} is a continuous variable, e_{ij} is initially taken as D_{ij} . In \mathcal{E}_{FF} , the directed edge $e_{v_i \rightarrow v_j}$ is regarded as the message passing from v_i to v_j as detailed in Section 2.4.

Feature imputation as edge prediction Based on the graph representation of the data matrix, the feature imputation problem is naturally transformed into an edge prediction task on the graph. Specifically, the prediction of the j -th feature of the i -th sample can be expressed as predicting e_{ij} in \mathcal{E}_{OF} . Edge prediction is conducted

by minimizing the difference between \widehat{e}_{ij} and e_{ij} . When e_{ij} corresponds to a discrete feature, the difference is measured using cross-entropy. In the case of a continuous feature, the difference is measured using Mean Squared Error (MSE).

Label prediction task in BCGNN When facing label prediction task, we use $\widehat{\mathbf{D}} = (\widehat{\mathbf{D}}_i)$ predicted by the upstream BCGNN as intermediate input. Then \widehat{Y}_i is obtained by using a simple forward propagation through a trainable mapping: $\widehat{Y}_i = f(\widehat{\mathbf{D}}_i)$, $\forall i \in \{1, \dots, n\}$. The mapping $f(\cdot)$ and upstream BCGNN can be learned by minimizing the difference between \widehat{Y}_i and Y_i . When Y_i is a discrete label, we use cross-entropy to measure the difference between \widehat{Y}_i and Y_i . When Y_i is a continuous label, we use MSE to measure it.

2.3 Bipartite Graph: Data-Driven Inductive Learning

During the graph learning, information from data is exclusively conveyed by edges connecting feature nodes and observation nodes, which are established based on the observed data. Inspired by the GraphSAGE architecture [34], we propose to use the bipartite graph $\mathcal{G}_B = (\mathcal{V}, \mathcal{E}_{OF})$ [9]. This framework shares similarities with GraphSAGE but distinguishes itself by the additional incorporation of edge embeddings. In \mathcal{G}_B , each node aggregates information from neighboring nodes and those corresponding edge embeddings for inductive learning.

Specifically, at each GNN layer l , to inductively learn the information from neighboring nodes in \mathcal{G}_B , observation and feature nodes aggregate messages from the embeddings of their neighboring nodes as well as the corresponding edge embeddings as follows: $\forall v \in \mathcal{V}$,

$$m_v^l = \text{AGG}_l \left\{ \sigma(P^l \cdot \text{CONCAT}(h_v^{l-1}, e_{uv}^{l-1}, h_u^{l-1}) + b^l) \mid \forall u \in \mathcal{N}(v, \mathcal{G}_B) \right\}$$

with $h_u^{l-1} \in \mathbb{R}^{d_n}$, $h_v^{l-1} \in \mathbb{R}^{d_n}$ and $e_{uv}^{l-1} \in \mathbb{R}^{d_e}$, where $P^l \in \mathbb{R}^{d_m \times (2d_n + d_e)}$ is a trainable weight matrix, $b^l \in \mathbb{R}^{d_m}$ is a trainable bias term, σ is an activation function, AGG_l is the aggregation function, and d_n , d_e , and d_m respectively denote the dimensions of the node, edge embedding and message. $\mathcal{N}(v, \mathcal{G})$ denotes the set of neighboring nodes of node v in graph \mathcal{G} . Subsequently, the embeddings of observation and feature nodes are updated based on aggregated messages:

$$h_v^l = \sigma(Q^l \cdot \text{CONCAT}(h_v^{l-1}, m_v^l) + c^l), \forall v \in \mathcal{V},$$

where $Q^l \in \mathbb{R}^{d_n \times (d_n + d_m)}$ is a trainable weight matrix, and $c^l \in \mathbb{R}^{d_n}$ is a trainable bias term. After the node embeddings are updated, the edge embeddings will also be updated by

$$e_{uv}^l = \sigma(W^l \cdot \text{CONCAT}(e_{uv}^{l-1}, h_v^l, h_u^l) + d^l), \forall e_{uv} \in \mathcal{E}_{OF},$$

where $W^l \in \mathbb{R}^{d_e \times (2d_n + d_e)}$ is a trainable weight matrix, and $d^l \in \mathbb{R}^{d_e}$ is a trainable bias term.

2.4 Complete Directed Graph: Feature Interdependence Structure Learning

The interdependence structure among features holds valuable information for imputing missing features, yet this aspect is currently neglected by existing methods. The complete directed graph $\mathcal{G}_C = (\mathcal{V}_F, \mathcal{E}_{FF})$ proposed in this subsection is specifically constructed to uncover the interdependence structure among features. Specifically, when dealing with a target feature node v , all $w \in \mathcal{V}_F$ with $w \neq v$ are considered as source nodes under the perspective of message passing. Our approach involves learning the directed dependence structure between h_v and each source feature node embedding h_w . This structure is characterized by message passing from h_w to h_v , with the direction from w to v .

To measure the varying degree of dependency of the target node feature v on different segments of the source node embedding h_w in the complete directed graph \mathcal{G}_C , we employ an element-wise attention weight vector $\alpha_{w \rightarrow v}$. This weight vector is trainable and can be acquired through an additive attention mechanism. In detail, at each GNN layer l , we initially implement an additive attention mechanism $\mathcal{A} : \mathbb{R}^{2d_n} \rightarrow \mathbb{R}^{d_n}$ on $\text{CONCAT}(h_w^{l-1}, h_v^{l-1})$ to derive the attention score vector $\text{Score}_{w \rightarrow v}^l$, where \mathcal{A} represents a single linear layer equipped with a trainable weight matrix $U_w^l \in \mathbb{R}^{d_n \times 2d_n}$, followed by a LeakyReLU activation function. In other words, the vector that represents the attention weight score of the target feature v on the source feature w can be expressed as

$$\text{Score}_{w \rightarrow v}^l = \sigma(U_w^l \cdot \text{CONCAT}(h_w^{l-1}, h_v^{l-1}) + g^l),$$

where σ is the LeakyReLU activation function and $g^l \in \mathbb{R}^{d_n}$ is a trainable bias term. It is noted that, across different target node features v , the weight matrix U_w remains shared for a given source node w . Subsequently, the attention score $\text{Score}_{w \rightarrow v}$ undergoes normalization through the softmax function to yield the element-wise attention weight vector $\alpha_{w \rightarrow v}$. Specifically, the k -th element of the normalized attention weight vector $\alpha_{w \rightarrow v}^l$ is written as

$$\alpha_{w \rightarrow v}^l[k] = \frac{\exp(\text{Score}_{w \rightarrow v}^l[k])}{\sum_{k' \in \{1, 2, \dots, d_n\}} \exp(\text{Score}_{w \rightarrow v}^l[k'])}.$$

Next, we calculate the Spearman correlation coefficient between features w and v . Let ρ_{wv} indicate the sign of the calculated Spearman correlation coefficient, which serves as an indicator of whether the message passing from node w to node v is positive or negative. The details of the calculation are deferred to Appendix A. Furthermore, when considering a specific target feature node v aggregating messages, the strength of messages from different source nodes should vary because the dependence relationships among features are distinct. We thus introduce a trainable scalar parameter $I_{w \rightarrow v}^l$ to quantify the strength of message passing from each source node w to the target node v in \mathcal{G}_C .

Finally, taking into account the attention weight vector $\alpha_{w \rightarrow v}^l$, we represent the dependency of the target feature v on the source feature w as

$$\text{Message}_{w \rightarrow v}^l = (\rho_{wv} \cdot I_{w \rightarrow v}^l \cdot \alpha_{w \rightarrow v}^l) \odot h_w^{l-1},$$

where \odot is the Hadamard product. Message $_{w \rightarrow v}^l$ is the so-called message passing from h_w to h_v at GNN layer l . Exchanging unique messages bidirectionally between each pair of feature nodes effectively captures the directed dependency structure among features.

2.5 Union Graph: Global Learning for Feature Imputation and Label Prediction

After the generation and transmission of messages in the bipartite graph \mathcal{G}_B and the complete directed graph \mathcal{G}_C , we introduce the union graph. In the union graph, feature nodes aggregate messages from both neighboring observation and feature nodes, while observation nodes aggregate messages from neighboring feature nodes. Specifically, at each GNN layer l , we represent the global message aggregation process in BCGNN as follows:

$$m_v^l = \text{AGG}_l \left\{ \begin{array}{l} \sigma(P^l \cdot \text{CONCAT}(h_v^{l-1}, e_{uv}^{l-1}, h_u^{l-1}) \\ + b^l), \rho_{uv} \cdot I_{w \rightarrow v}^l \cdot \alpha_{w \rightarrow v}^l \odot h_w^{l-1} \\ \forall u \in \mathcal{N}(v, \mathcal{G}_B), \forall w \in \mathcal{N}(v, \mathcal{G}_C) \end{array} \right\}, \forall v \in \mathcal{V}_F,$$

$$m_v^l = \text{AGG}_l \left\{ \begin{array}{l} \sigma(P^l \cdot \text{CONCAT}(h_v^{l-1}, e_{uv}^{l-1}, h_u^{l-1}) \\ + b^l) | \forall u \in \mathcal{N}(v, \mathcal{G}_B) \end{array} \right\}, \forall v \in \mathcal{V}_O.$$

After the message passing and aggregation, the embedding of every node $v \in \mathcal{V}$ is updated using h_v^{l-1} and the aggregated message m_v^l as follows:

$$h_v^l = \sigma(Q^l \cdot \text{CONCAT}(h_v^{l-1}, m_v^l) + c^l).$$

Then, BCGNN updates the embedding of edges $e_{uv} \in \mathcal{E}_{OF}$:

$$e_{uv}^l = \sigma(W^l \cdot \text{CONCAT}(e_{uv}^{l-1}, h_v^l, h_u^l) + d^l),$$

which serves as the input for generating messages in the next layer of BCGNN. After BCGNN undergoes forward propagation through L layers of updates, a simple linear readout layer is applied for the final feature imputation:

$$\hat{e}_{uv} = O_{\text{feature}}(h_u^L, h_v^L).$$

This process yields $\hat{\mathbf{D}}$. For a label prediction task related to the i -th observation, a linear readout layer is applied to the combined imputed feature edges of the i -th observation:

$$\hat{Y}_i = O_{\text{label}}(\hat{\mathbf{D}}_i).$$

We use distinct d_n -dimensional one-hot vectors to initialize the embedding h_v for each feature node in \mathcal{V}_F . For observation nodes in \mathcal{V}_O , the embedding h_u is initialized using the d_n -dimensional constant vector $\mathbf{1}_{d_n}$.

2.6 Dropout Schemes in BCGNN

To address the over-fitting problem in BCGNN, we use the DropEdge scheme [33] in the two subgraphs. In training phase, we randomly mask out some edges in both \mathcal{E}_{OF} and \mathcal{E}_{FF} before conducting forward propagation in each epoch:

$$\text{DropB}(\mathcal{E}_{OF}, r_B) = \{e_{ij} | e_{ij} \in \mathcal{E}_{OF}, M_{ij}^B > r_B, M_{ij} = 1\}$$

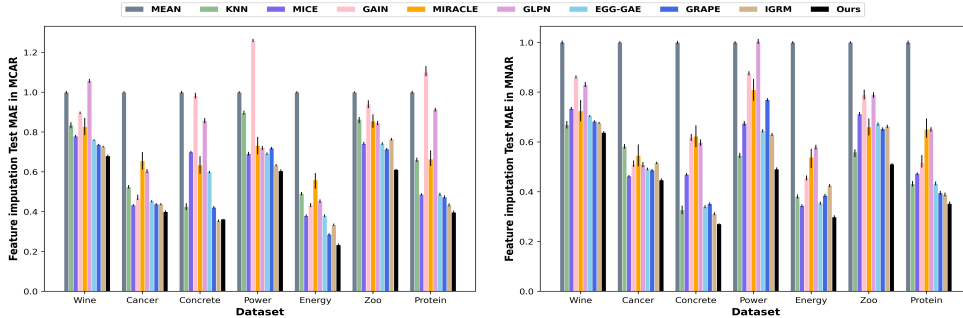


Fig. 3 Average test MAE of feature imputation at a missing rate of 0.3 under MCAR (Left) and MNAR (Right) in UCI datasets over 5 random trials. The results are normalized by the average performance of the Mean imputation.

$$\text{DropC}(\mathcal{E}_{FF}, r_C) = \{e_{v_i \rightarrow v_j} | e_{v_i \rightarrow v_j} \in \mathcal{E}_{FF}, M_{ij}^C > r_C\}$$

where r_B and r_C are dropout rates, the mask matrices $\mathbf{M}^B = (M_{ij}^B) \in \mathbb{R}^{n \times m}$ and $\mathbf{M}^C = (M_{ij}^C) \in \mathbb{R}^{c \times c}$ are both sampled uniformly within the interval $[0, 1]$, with $c = m(m - 1)$. Then, we only input those edges that are not masked into BCGNN for forward computation, and compute the training loss exclusively on the masked edges in \mathcal{E}_{FF} . This approach ensures that the predicted values of edges during training are definitely not functions of themselves, thus avoiding the problem of over-fitting. In the testing phase, the entire graph \mathcal{G} without Dropedge is fed into BCGNN to conduct feature imputation. Furthermore, aim at improving the generalization capacity of BCGNN, we incorporate the AttentionDrop technique in element-wise attention mechanism between feature nodes in \mathcal{G}_C :

$$\text{AttentionDrop}(\alpha_{w \rightarrow v}, a_{drop}) = \{\alpha_{w \rightarrow v}[k] \cdot I(b_k > a_{drop})\},$$

where b_k is uniformly sampled within $[0, 1]$ and a_{drop} is the dropout rate. The detailed forward computation of BCGNN can be found in Algorithm 1.

3 Experiments

3.1 Datasets and Baseline Models

Datasets. We conduct experiments on seven datasets from the UCI Machine Learning Repository. These datasets span diverse domains including medicine (Cancer), civil engineering (Concrete, Energy), industry (Power), business (Wine), and biology (Protein, Zoo). Cancer has the highest feature dimension, consisting of 32 features with 569 observations. Protein has the largest sample size, which contains over 45,000 observations. The Wine, Concrete, Protein, and Energy datasets include a mix of discrete and continuous features. Zoo exclusively comprises discrete features, while the Power and Cancer datasets only have continuous features. The feature values in all dataset are scaled to $[0, 1]$ with a MinMax scaler [35].

Configurations. Since the datasets are fully observed, we simulate missing data by generating a masking matrix \mathbf{M} according to each distinct missing mechanism

Algorithm 1 BCGNN forward computation

Input: Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $\mathcal{V} = \mathcal{V}_F \cup \mathcal{V}_O$ and $\mathcal{E} = \mathcal{E}_{OF} \cup \mathcal{E}_{FF}$; Number of layers L ; DropEdge rate r_B and r_C ; AttentionDrop rate a_{drop} ; $\mathcal{N}(v, \mathcal{G})$ is a mapping providing the set of neighboring nodes of node v within the graph \mathcal{G} ; P^l for message passing; Q^l for node updating; W^l for edge updating; U_w^l for attention weight score calculation; $I_{w \rightarrow v}^l$ quantifying the strength of message passing from node w to node v .

for l **in** $\{1, \dots, L\}$ **do**

$h_v^{(0)} \leftarrow \text{INIT}(v), \forall v \in \mathcal{V}; e_{uv}^{(0)} \leftarrow \text{INIT}(e_{uv}), \forall e_{uv} \in \mathcal{E}$

$\mathcal{E}_{dropB}, \mathcal{E}_{dropC} \leftarrow \text{DropB}(\mathcal{E}_{OF}, r_B), \text{DropC}(\mathcal{E}_{FF}, r_C); \mathcal{G}_{dropB}, \mathcal{G}_{dropC} \leftarrow (\mathcal{V}, \mathcal{E}_{dropB}), (\mathcal{V}_F, \mathcal{E}_{dropC});$

for $v \in \mathcal{V}$ **do**

for $u \in \mathcal{N}(v, \mathcal{G}_{dropB})$ **do**

MessageB $_{uv}^l = \sigma(P^l \cdot \text{CONCAT}(h_v^{l-1}, e_{uv}^{l-1}, h_u^{l-1}) + b^l)$

end for

if $v \in \mathcal{V}_F$ **then**

for $w \in \mathcal{N}(v, \mathcal{G}_{dropC})$ **do**

Score $_{w \rightarrow v}^l = \sigma(U_w^l \cdot \text{CONCAT}(h_w^{l-1}, h_v^{l-1}) + g^l)$

$\alpha_{w \rightarrow v}^l = \text{AttentionDrop}(\text{Softmax}(\text{Score}_{w \rightarrow v}^l), a_{drop})$

MessageC $_{w \rightarrow v}^l = (\rho_{wv} \cdot I_{w \rightarrow v}^l \cdot \alpha_{w \rightarrow v}^l) \odot h_w^{l-1}$

end for

end if

if $v \in \mathcal{V}_O$ **then**

$m_v^l = \text{AGG}_l(\text{MessageB}_{uv}^l | \forall u \in \mathcal{N}(v, \mathcal{G}_{dropB}))$

else if $v \in \mathcal{V}_F$ **then**

$m_v^l = \text{AGG}_l(\text{MessageB}_{uv}^l, \text{MessageC}_{w \rightarrow v}^l | \forall u \in \mathcal{N}(v, \mathcal{G}_{dropB}), \forall w \in \mathcal{N}(v, \mathcal{G}_{dropC}))$

end if

$h_v^l = \sigma(Q^l \cdot \text{CONCAT}(h_v^{l-1}, m_v^l) + c^l)$

end for

for $e_{uv} \in \mathcal{E}_{dropB}$ **do**

$e_{uv}^{(l)} = \sigma(W^l \cdot \text{CONCAT}(h_v^l, e_{uv}^{l-1}, h_u^l) + d^l)$

end for

end for

Output for Feature Imputation: $\hat{e}_{uv} = O_{\text{feature}}(h_u^L, h_v^L), \forall e_{uv} \in \mathcal{E}_{OF}$

Output for Label Prediction: $\hat{Y}_i = O_{\text{label}}(\hat{\mathbf{D}}_i), \forall i \in \{1, \dots, n\}$

and deleting observed values from the data matrix based on \mathbf{M} . We consider three types of missing mechanisms: missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR) [36]. The detailed procedures for generating the masking matrices for various missing mechanisms are provided in Appendix B, and the configurations of BCGNN can be found in Appendix C.

Baseline models. In our experiments, we compare BCGNN with 10 baseline models: (1) Mean: imputation using the feature-wise mean; (2) KNN [19]; (3) MICE [17]; (4) GAIN [23]; (5) Decision Tree: we only use it in the label prediction task [7]; (6) MIRACLE: we only use it in the feature imputation task [10]; (7) GRAPE [9]; (8) GLPN: a method has demonstrated superiority over VGAE and GDN [11], and is only used in the feature imputation task; (9) EGG-GAE [15]; (10) IGRM [14].

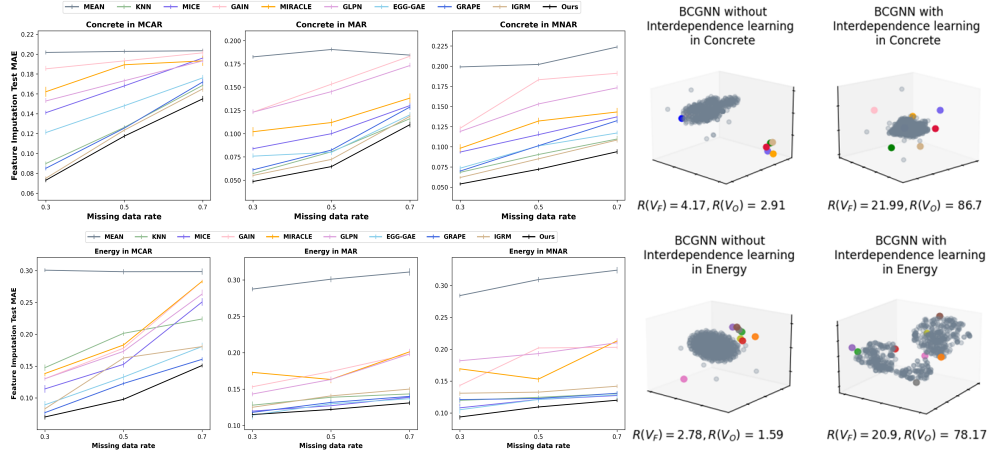


Fig. 4 **Left:** Average test MAE of feature imputation under MCAR, MAR and MNAR with different missing rates in the Concrete and Energy datasets over 5 random trials. **Right:** The embedding spaces V_F and V_O for feature and observation nodes, respectively. They are obtained from the trained BCGNN with/without learning interdependence structure in Concrete and Energy under MAR with a missing rate of 0.3. The **colored** dots represent feature node embeddings and the **grey** dots represent observation node embeddings.

3.2 Feature Imputation

At a missing rate of 0.3, we compare the feature imputation performance of BCGNN with the nine baseline models on the seven datasets under the three missing mechanisms. We record the averaged test error for each method over five random trials. The test error is defined as the mean absolute error (MAE) between D_{ij} and \hat{D}_{ij} over all entries with $M_{ij} = 0$. The results under MCAR and MNAR are presented in Figure 3, and the results for MAR can be found in Figure C1 in Appendix D.1. Compared with the best baseline model GRAPE or IGRM, our BCGNN achieves an average reduction of **11.8%**, **13.2%**, and **16.1%** in the test MAE for MCAR, MAR and MNAR, respectively.

We further consider missing rates of 0.5 and 0.7. At a missing rate of 0.5, compared with the best baseline, GRAPE or IGRM, BCGNN reduces the test MAE by an average of 13.5%, 12.3%, and 16.1% under MCAR, MAR and MNAR respectively. At a missing rate of 0.7, it reduces the test MAE by 10.2%, 11.4% and 12.6%, respectively. Figure 4 (left panel) shows the specific results for feature imputation under various missing rates and missing mechanisms in Concrete and Energy datasets, and the results for all the other datasets are available in Figure C2 in Appendix D.1.

3.3 Feature Embedding Ability of BCGNN

According to Jegelka [37] and Zhang et al. [38], the feature embedding ability of a GNN can be expressed by the size of the embedding space V , which is the embedding output space obtained from the GNN mapping. Here we use lossy coding to measure the size of V with $V = (h_{v_1}, h_{v_2}, \dots, h_{v_m}) \in \mathbb{R}^{m \times d_n}$. The lossy coding scheme maps V to a sequence of binary bits, such that the original embeddings can be recovered up

to an allowable distortion: $E[\|h_{v_i} - \hat{h}_{v_i}\|^2] \leq \epsilon^2$. Then the size of V can be quantified as: $R(V) = 0.5 \log_2 \det(I + mVV^T/(\epsilon^2 d_n))$ [39].

We obtain the output embeddings from the experiments conducted in Section 3.2 for the regular BCGNN and the BCGNN without interdependence structure learning. The latter is achieved by setting all ρ_{wv} between feature nodes to zero, effectively disabling the message passing between feature nodes. Then we evaluate the feature embedding ability of the two types of BCGNN by comparing the sizes of their feature node embedding space $R(V_F)$ and observation node embedding space $R(V_O)$. We report the results for all datasets under different missing mechanisms in Figure 4 (right panel), Figures C2 (right panel), D5 and D6 in the Appendix. Clearly, the BCGNN with learned interdependence structure among features exhibits larger feature and observation embedding spaces than the BCGNN without it. This directly indicates that the learning of interdependence structure indeed improves the feature embedding ability of BCGNN.

3.4 Generalization on New Observations

When addressing the feature imputation task for new observations, we hope to directly use their observed features as the input to impute missing values, without the need for retraining the model. Thus we investigate the generalization ability of BCGNN in handling new observations. In our experiments, we randomly divide the entire dataset $\mathbf{D} \in \mathbb{R}^{n \times m}$ into two groups of equal sample sizes: $\mathbf{D}^{old} \in \mathbb{R}^{n_{old} \times m}$ and $\mathbf{D}^{new} \in \mathbb{R}^{n_{new} \times m}$, representing the initial and new observations, respectively. As detailed in Appendix B, we respectively generate masking matrices \mathbf{M}^{old} and \mathbf{M}^{new} for \mathbf{D}^{old} and \mathbf{D}^{new} under the three missing mechanisms with a missing rate of 0.3. Then we train BCGNN on the initial observations D_{ij}^{old} with $M_{ij}^{old} = 1$. During the test phase, we employ the new observations D_{ij}^{new} with $M_{ij}^{new} = 1$ as the input for the trained BCGNN and directly impute \hat{D}_{ij}^{new} with $M_{ij}^{new} = 0$. We repeat the same process for the baseline models requiring training (MIRACLE, GAIN, GRAPE, IGRM, EGG-GAE and GLPN). For the other baselines without training, we directly employ them on \mathbf{D}^{new} to impute \hat{D}_{ij}^{new} with $M_{ij}^{new} = 0$. As shown in Figure D3 in Appendix D.2, BCGNN demonstrates outstanding generalization ability on new observations, which yields **11.3%**, **12.8%**, and **14.2%** lower MAE compared with the best baseline, GRAPE, under MCAR, MAR, and MNAR, respectively.

3.5 Label Prediction

To evaluate the performance of BCGNN in label prediction task, we randomly split the target label \mathbf{Y} into a training set \mathbf{Y}_{train} and a test set \mathbf{Y}_{test} with a sample size ratio of 7:3, which are respectively used as observed and missing dataset. BCGNN can then be trained by minimizing the difference between \hat{Y}_i and Y_i over all i with $M_i^y = 1$. For all baseline models except the decision tree, since there are no corresponding end-to-end prediction measures, we first impute and obtain the completed data $\hat{\mathbf{D}}$, and then use a fully connected layer on each $\hat{\mathbf{D}}_i$ to predict \hat{Y}_i . As shown in Figure D4 in Appendix D.2, BCGNN yields **8.8%**, **6.8%**, and **7.3%** lower MAE compared with

the best baseline, GRAPE, under MCAR, MAR, and MNAR respectively over the 7 datasets.

3.6 DropEdge and AttentionDrop

We evaluate the performance of BCGNN under three missing mechanisms with a missing rate of 0.3, which utilizes DropEdge on both bipartite graph and complete directed graph, as well as AttentionDrop on the complete directed graph. We evaluate the BCGNN variants without one, two, or three of these schemes. We report the averaged MAE for feature imputation under different scenarios in Table D2 in Appendix D.4. It shows that using DropEdge on the bipartite graph significantly improves BCGNN’s performance by consistently reducing the average test MAE by 35%, employing DropEdge on the complete directed graph also leads to a 5% reduction in the average test MAE, while AttentionDrop on the complete directed graph achieves an average reduction of 3% in the test MAE.

4 Conclusion

In this paper, we propose BCGNN, a union of bipartite graph and complete directed graph, which addresses feature imputation and label prediction tasks by inductively learning node representations and explicitly learning interdependence structure among features. The directed complete graph is designed to learn and express the interdependence structure, thereby enhancing our model’s expressive capability. Compared to SOTA methods, our model demonstrates significant improvements in accuracy, and robustly adapts to various missing ratios, missing mechanisms, and unseen data points.

Appendix A Spearman Correlation Coefficient Estimator and Its Sign Indicator

A positive correlation coefficient between two features indicates a consistent relationship, where both features tend to increase or decrease simultaneously. Conversely, a negative coefficient signifies an inverse relationship, meaning that when one feature increases, the other tends to decrease. Therefore, in BCGNN, it is crucial to estimate the sign of the correlation between each pair of features to facilitate message passing among feature nodes. We opt for using Spearman correlation coefficient to measure this correlation.

The original formula for estimating the Spearman correlation coefficient between two features X and Y based on n observations is:

$$S_{xy} = \frac{\frac{1}{n} \sum_{i=1}^n (R(x_i) - \bar{R}(x)) \cdot (R(y_i) - \bar{R}(y))}{\sqrt{\frac{1}{n} \sum_{i=1}^n (R(x_i) - \bar{R}(x))^2 \cdot \frac{1}{n} \sum_{i=1}^n (R(y_i) - \bar{R}(y))^2}}$$

where $R(x_i)$ represents the rank of x_i among the n observations for feature X . In practice, BCGNN employs a simplified equivalent formula for estimating the Spearman correlation:

$$S_{xy} = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)},$$

where $d_i = R(x_i) - R(y_i)$. When estimating the Spearman correlation coefficient between features X and Y the presence of missing data, we consider only the observations where both X and Y are simultaneously recorded. In scenarios with high missing rates or complex missing mechanisms, the absolute value of the estimated Spearman correlation coefficient often fails to fully and accurately represent the strength of dependency between the two features [40]. Nevertheless, we extract the signs of estimated Spearman correlation coefficients to facilitate message passing among feature nodes in BCGNN. Specifically, for features w and v with an estimated Spearman correlation coefficient S_{wv} , we define the sign indicator ρ_{wv} as follows: when $|S_{wv}| < 0.1$, $\rho_{wv} = 0$; when $|S_{wv}| \geq 0.1$, $\rho_{wv} = 1$ if $S_{wv} > 0$, and $\rho_{wv} = -1$ if $S_{wv} < 0$.

The selection of 0.1 as a threshold is based on a common consensus in numerous studies. It is generally accepted that when the absolute value of the Spearman correlation coefficient falls below 0.1, the correlation between the random variables is considered negligible. [41, 42]. Moreover, when an estimated Spearman correlation coefficient computed from incomplete data is near zero, there is a realistic possibility that its true counterpart has an opposite sign. Thus, we conservatively set its corresponding sign indicator to 0, effectively circumventing potential inaccuracy in the sign estimation under high missingness and complex missing mechanisms. In this setting, only strong correlations are used to facilitate the feature imputation.

Appendix B Generating missingness

Consider the k -th subject and the i -th feature. We assume that M_{ki} follows a Bernoulli distribution with success probability π_{ki} . The following text describes our approach for creating synthetic datasets that adhere to MCAR, MAR and MNAR patterns of missing data.

B.1 Missing Completely At Random (MCAR)

The missing probability π_{ki} remains constant, and is set to the current experiment’s missing rate.

B.2 Missing At Random (MAR)

We assume that π_{ki} is determined by the observed values of the previous $i - 1$ features of the k -th subject (if observed). Specifically,

$$\pi_{ki} = \frac{p(i) \cdot n \cdot \exp(\sum_{j < i} w_j m_j D_{kj} + b_j(1 - m_j))}{\sum_{l=1}^n \exp(\sum_{j < i} w_j m_j D_{lj} + b_j(1 - m_j))},$$

where $p(i)$ is the average missing rate of the i -th feature, w_j and b_j are sampled from a uniform distribution $U(0, 1)$ once for each dataset, and m_j represents the dependency of the missingness of D_{ki} on D_{kj} .

B.3 Missing Not At Random (MNAR)

The missingness of D_{ki} depends on its own value. Specifically,

$$\pi_{ki} = \frac{p(i) \cdot n \cdot \exp(-w_i D_{ki})}{\sum_{l=1}^n \exp(-w_i D_{li})},$$

where $p(i)$ is the average missing rate of the i -th feature, and w_i is sampled from a uniform distribution $U(0, 1)$, representing the dependency of the missingness of D_{ki} on its own value.

Appendix C Configurations of BCGNN

In the experiments across all datasets, we train BCGNN for 20,000 epochs using Adam optimizer with a learning rate of 0.001. The experiments follow the parameter settings in [43] and [44], and use a three-layer BCGNN with both node embeddings and edge embeddings of 64 dimensions. The AGG_l is implemented using a mean pooling function $MEAN(\cdot)$. ReLU is used as the activation function in message generation, node updates and edge updates, while LeakyReLU is used in additive attention. The DropEdge ratio is set to 0.5, and the AttentionDrop ratio is set to 0.3. O_{feature} and O_{label} are both implemented with a single linear layer. For all experiments, we run 5 trials with different random seeds and report the mean and standard deviation of the results.

Appendix D Additional Experimental Results

D.1 Additional Experimental Results for Feature Imputation

Figure C1 reports the average test MAE of BCGNN and baseline models in UCI datasets under MAR with a missing rate of 0.3 over 5 random trials, and Figure C2 respectively reports the average test MAE of BCGNN and baseline models on the Wine, Cancer, Power, Energy and Protein datasets under different missing mechanisms and different missing rate over 5 random trials.

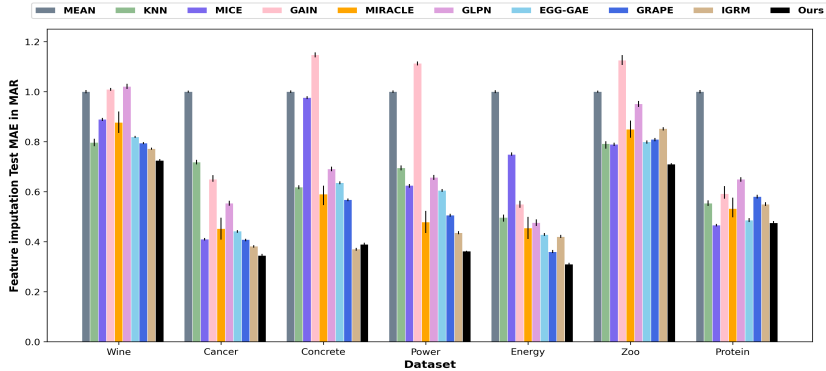


Fig. C1 Average test MAE of feature imputation at a missing rate of 0.3 under MAR in UCI datasets over 5 random trials. The results are normalized by the average performance of the Mean imputation.

D.2 Experimental Results for Generalization on Unseen Data and Label Prediction

At a missing rate of 0.3, we test the performance of BCGNN and baseline models in feature imputation tasks when handling new observations with missing data under different missing mechanisms. We report the average MAE in Figure D3.

At a missing rate of 0.3, we evaluate the performance of BCGNN and baseline models in the label prediction task with missing data under various missing mechanisms. We present the average MAE for label prediction in Figure D4.

D.3 Ablation Study on Aggregation Function and Correlation Estimates

Aggregation Function. In BCGNN, to perform message aggregation on both feature nodes and observation nodes, different aggregation functions such as $\text{MEAN}(\cdot)$, $\text{SUM}(\cdot)$, and $\text{MAX}(\cdot)$ can be used. We evaluate the impact of different aggregation functions on the performance of BCGNN by replicating the experiments in Section 3.2 under MCAR with a missing rate of 0.3.

Table D1 The average test MAE for BCGNN with different aggregation functions and correlation coefficients

	Concrete	Cancer	Energy	Protein	Wine	Zoo	Power
BCGNN with $\text{SUM}(\cdot)$	0.0732	0.0422	0.0701	0.0318	0.0683	0.1483	0.1023
BCGNN with $\text{MEAN}(\cdot)$	0.0731	0.0423	0.0706	0.0310	0.0678	0.1452	0.1012
BCGNN with $\text{MAX}(\cdot)$	0.0753	0.0451	0.0752	0.0351	0.0701	0.1501	0.1085
BCGNN with Spearman	0.0731	0.0423	0.0706	0.0310	0.0678	0.1452	0.1012
BCGNN with Pearson	0.0770	0.0454	0.0724	0.0325	0.0688	0.1480	0.1072
BCGNN with Kendall	0.0755	0.0441	0.0715	0.0321	0.0681	0.1511	0.1041

From Table D1, it is evident that BCGNN performs best with the $\text{SUM}(\cdot)$ function, while it particularly underperforms with $\text{MAX}(\cdot)$. The reason for this is that employing $\text{MAX}(\cdot)$ as the aggregation function in the message aggregation process for feature

nodes may lead to an overemphasis on messages from neighboring feature nodes (due to their tendency to be stronger than those from observation nodes), thereby overshadowing the messages received from observation nodes. This unbalanced choice reduces the expressive capacity of feature nodes.

Correlation Estimators. We investigate the influence of various correlation estimators on the performance of BCGNN. Specifically, we utilize Spearman, Pearson, and Kendall correlation coefficients in the model, respectively. Table D1 indicates that the choice of correlation estimator truly impacts BCGNN’s performance. Spearman and Kendall correlation coefficients generally demonstrate superior performance, whereas the Pearson correlation coefficient shows relatively weaker results. This discrepancy arises from the fact that the Pearson correlation coefficient is primarily suitable for continuous data and is constrained to capturing linear relationships exclusively. Conversely, both Kendall and Spearman correlation coefficients are computed based on ranks, rendering them nonparametric methods that do not require distributional assumptions and are adept at capturing nonlinear relationships.

D.4 Ablation Study on DropEdge and AttentionDrop Scheme

We evaluate the performance of BCGNN under three missing mechanisms with a missing rate of 0.3, which utilizes DropEdge on both bipartite graph and complete directed graph, as well as AttentionDrop on the complete directed graph. We evaluate the BCGNN variants without one, two, or three of these schemes. We report the averaged MAE for feature imputation under different scenarios in Table D2.

Table D2 Average test MAE of feature imputation tasks for BCGNN with or without DropEdge and AttentionDrop schemes

MCAR			Concrete	Cancer	Energy	Zoo	Wine	Power	Protein
Drop in $\mathcal{G}_B \times$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.142	0.072	0.120	0.171	0.109	0.152	0.061
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.078	0.046	0.075	0.149	0.070	0.105	0.032
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \times	0.075	0.044	0.072	0.146	0.066	0.103	0.030
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \checkmark	0.073	0.042	0.070	0.145	0.066	0.101	0.030
MAR			Concrete	Cancer	Energy	Zoo	Wine	Power	Protein
Drop in $\mathcal{G}_B \times$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.073	0.072	0.143	0.169	0.073	0.132	0.071
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.053	0.054	0.119	0.149	0.061	0.102	0.052
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \times	0.051	0.051	0.115	0.147	0.059	0.098	0.049
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \checkmark	0.049	0.050	0.115	0.145	0.057	0.096	0.048
MNAR			Concrete	Cancer	Energy	Zoo	Wine	Power	Protein
Drop in $\mathcal{G}_B \times$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.072	0.073	0.123	0.161	0.081	0.122	0.072
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \times$	AttDrop \times	0.059	0.055	0.100	0.145	0.066	0.102	0.052
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \times	0.056	0.053	0.096	0.142	0.063	0.099	0.048
Drop in $\mathcal{G}_B \checkmark$	Drop in $\mathcal{G}_C \checkmark$	AttDrop \checkmark	0.054	0.052	0.094	0.141	0.061	0.097	0.046

Table D3 Average running time (in seconds) for feature imputation by different methods

Model	Concrete	Energy	Cancer	Zoo	Protein	Wine	Power
Mean	0.0008	0.0010	0.0010	0.0008	0.0135	0.0011	0.0015
KNN	0.2251	0.1343	0.1720	0.0278	636.28	0.5041	11.459
MICE	0.0310	0.0321	0.3720	0.0281	0.2716	0.0531	0.0291
SVD	0.0691	0.0202	0.2710	0.0432	0.5932	0.0564	0.1427
Spectral	0.0710	0.0583	0.0830	0.0341	1.4333	0.0978	0.1814
GAIN	0.1202	0.0132	0.0172	0.0121	0.0501	0.0131	0.0146
MIRACLE	0.0912	0.0928	0.1210	0.8711	2.4213	0.0872	0.0871
GLPN	0.1821	0.1921	0.2121	0.1621	0.5213	0.1732	0.1853
GRAPE	0.0273	0.0151	0.0331	0.0133	0.5683	0.0199	0.0488
BCGNN	0.0403	0.0232	0.1052	0.0231	0.8323	0.0253	0.0572

D.5 Running Time Comparison

We present the running times associated with feature imputation using various methods in our analysis. Specifically, for methods such as Mean, KNN, MICE, SVD, and Spectral, the reported running time corresponds to the duration of a single function call dedicated to imputing missing data. Conversely, for models like GAIN, MIRACLE, GLPN, GRAPE, and our BCGNN, the running time is defined as the duration of one forward pass through the model. Table D3 shows the average running time over 5 random trials for the experiments in Section 3.2.

Appendix E Computational Hardware

All models are trained on a Windows 10 64-bit OS (version 19045) with 16GB of RAM (AMD Ryzen 7 4800H CPU @ 2.9GHz) and 2 NVIDIA GeForce RTX 2060 with Max-Q Design GPUs.

Appendix F Broader Impacts and Limitations

Our method demonstrates superior performance in both feature imputation and label prediction tasks by leveraging the complex interdependence among features for tabular data. It enhances the applicability of missing data imputation tasks to real-world problems, such as clinical research, finance, economics, and microarray analysis. Ethically, we believe that our rather fundamental work has minimal potential for misuse. One limitation is our inability to effectively integrate the causal structure among feature nodes to improve feature imputation accuracy. This remains a significant area of research interest for us in the future.

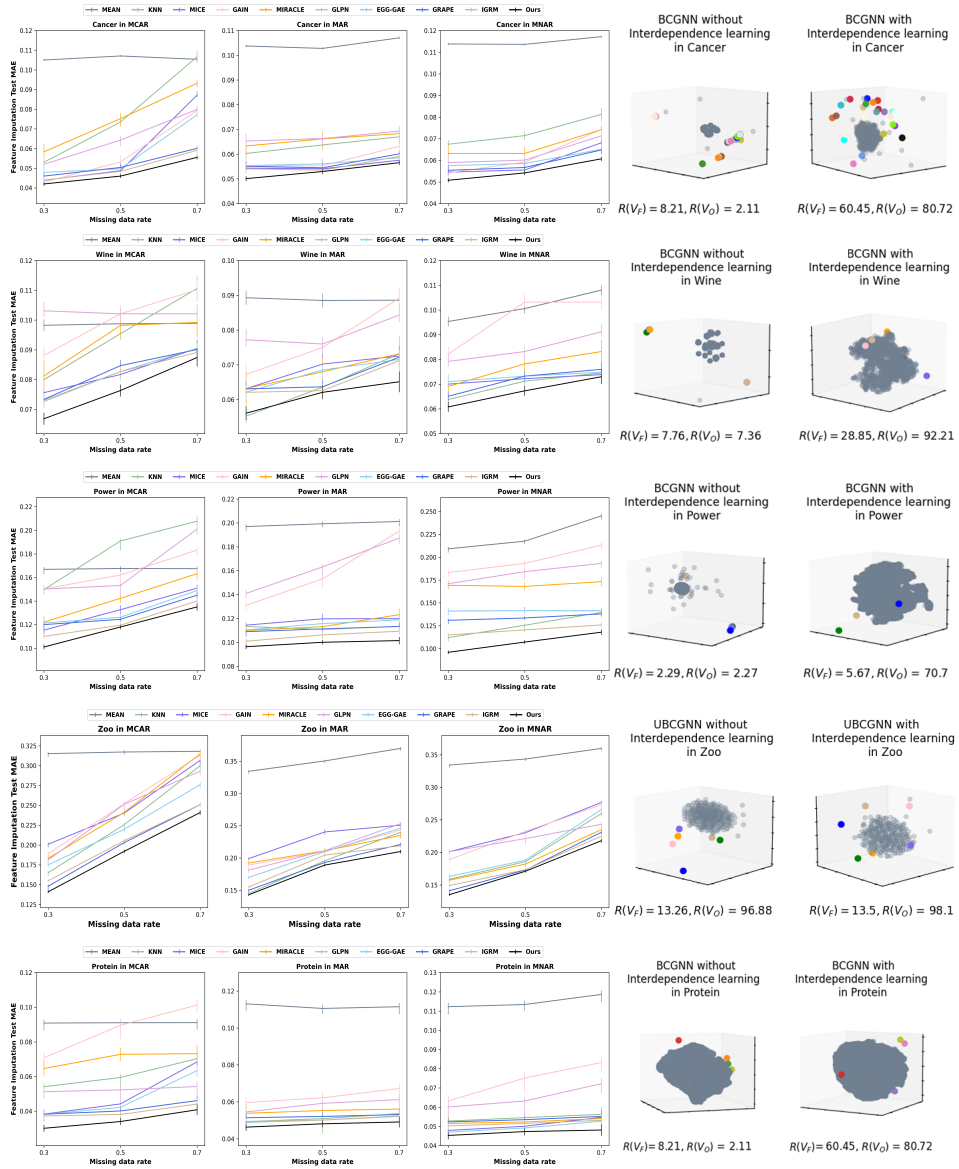


Fig. C2 Left: Average test MAE of feature imputation under MCAR, MAR and MNAR with different missing rates in Cancer, Wine, Power, Energy, Protein over 5 random trials. Right: The embedding spaces V_F and V_O for feature and observation nodes, respectively. They are obtained from the trained BCGNN with/without learning interdependence structure in the above datasets under MAR with a missing rate of 0.3. The colored dots represent feature node embeddings and the grey dots represent observation node embeddings.

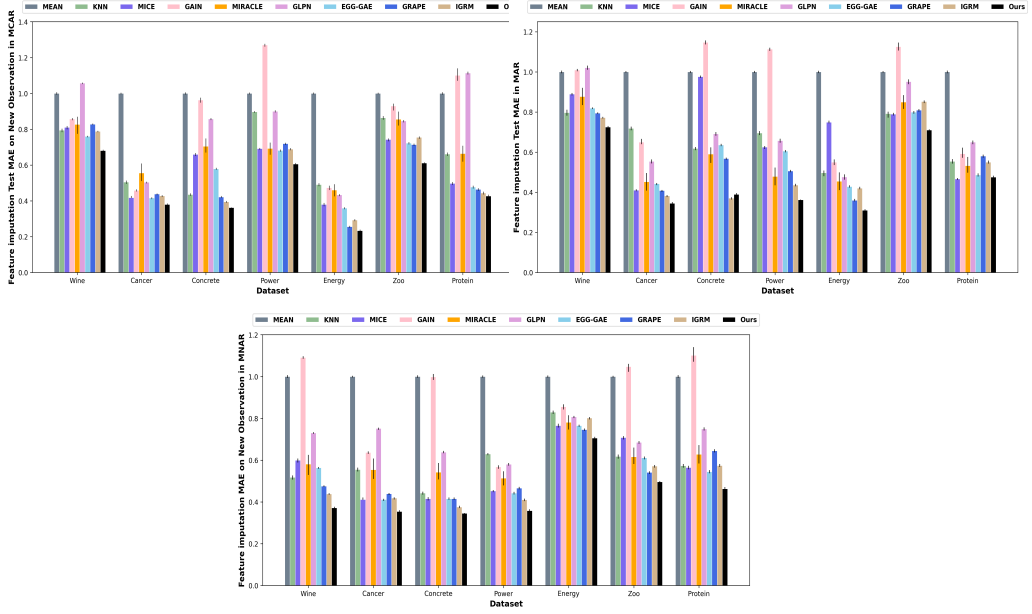


Fig. D3 Average MAE of feature imputation on new observations at a missing rate of 0.3 under MCAR, MAR and MNAR in UCI datasets over 5 random trials. The result is normalized by the average performance of the Mean imputation. Compared with the best baseline (GRAPE), BCGNN demonstrates reductions in the average MAE of 12.3%, 14.8%, and 18.2% under MCAR, MAR and MNAR, respectively.

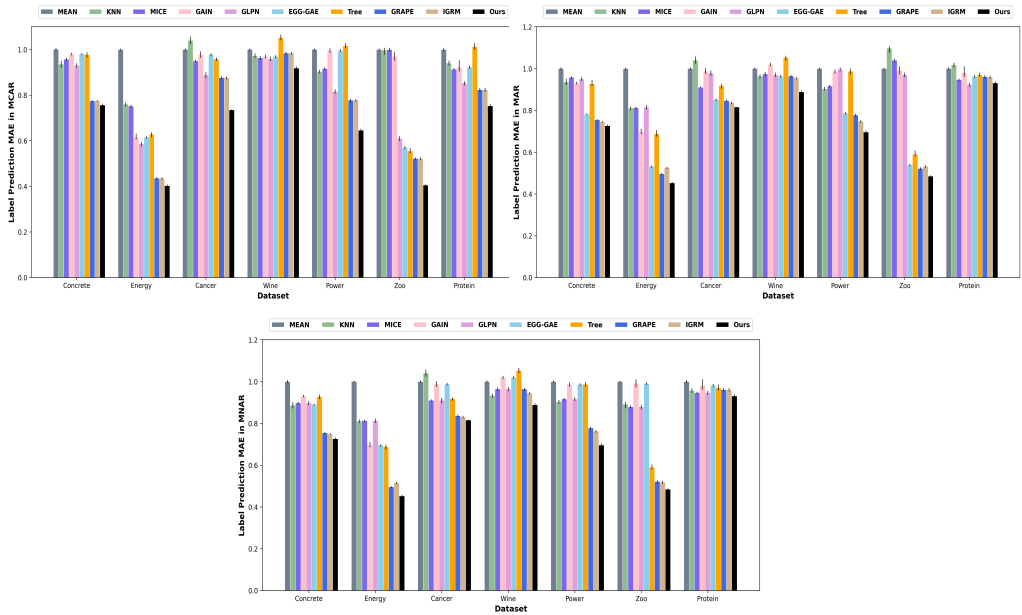


Fig. D4 Average MAE of label prediction at a missing rate of 0.3 under MCAR, MAR and MNAR in UCI datasets over 5 random trials. The results are normalized by the average performance of the Mean imputation.

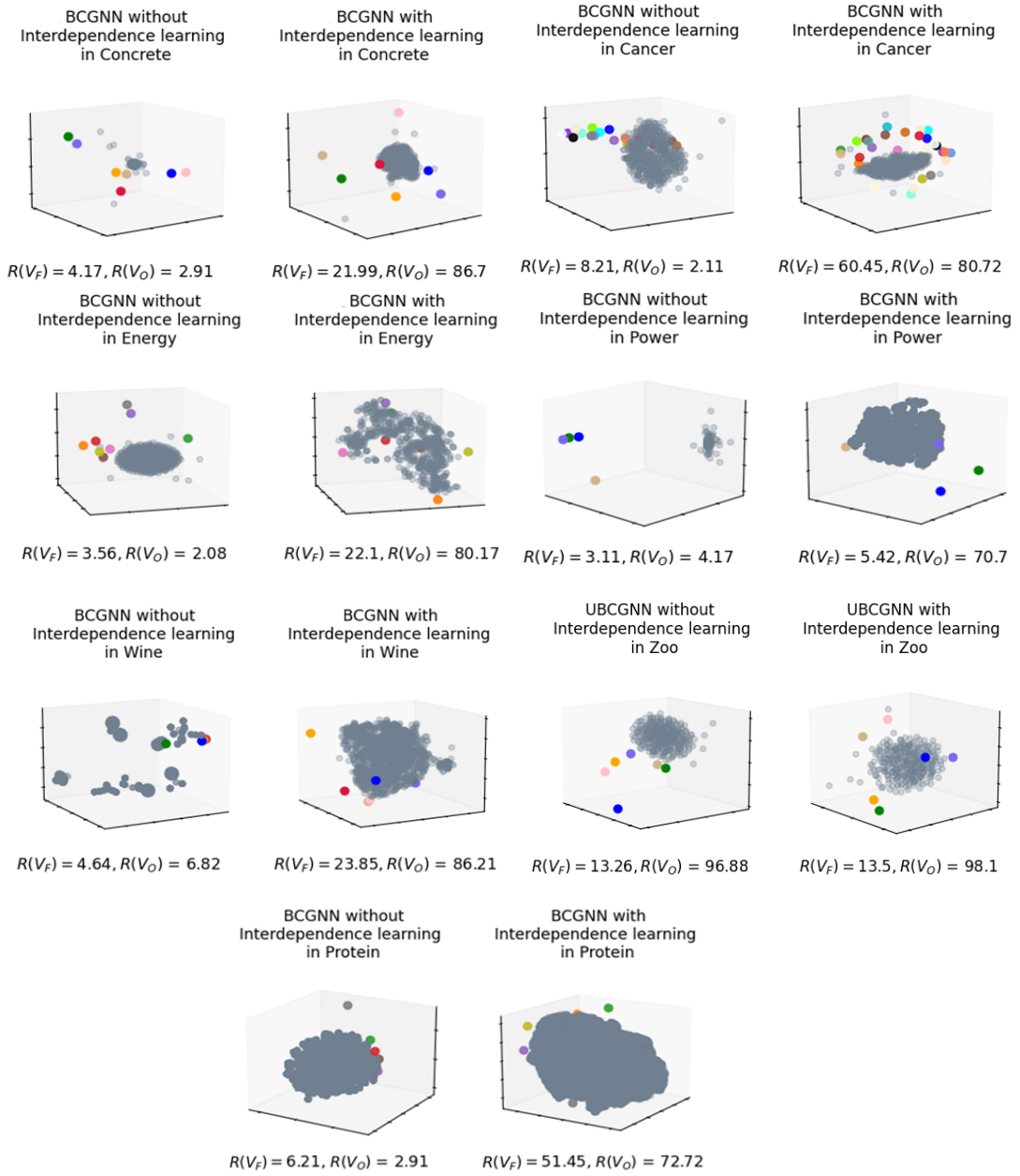


Fig. D5 The sizes of feature embedding space V_F and observation embedding space V_O , which are obtained from the trained BCGNN with/without interdependence structure learning under MCAR by the dimension reduction visualization method t-SNE for all datasets.

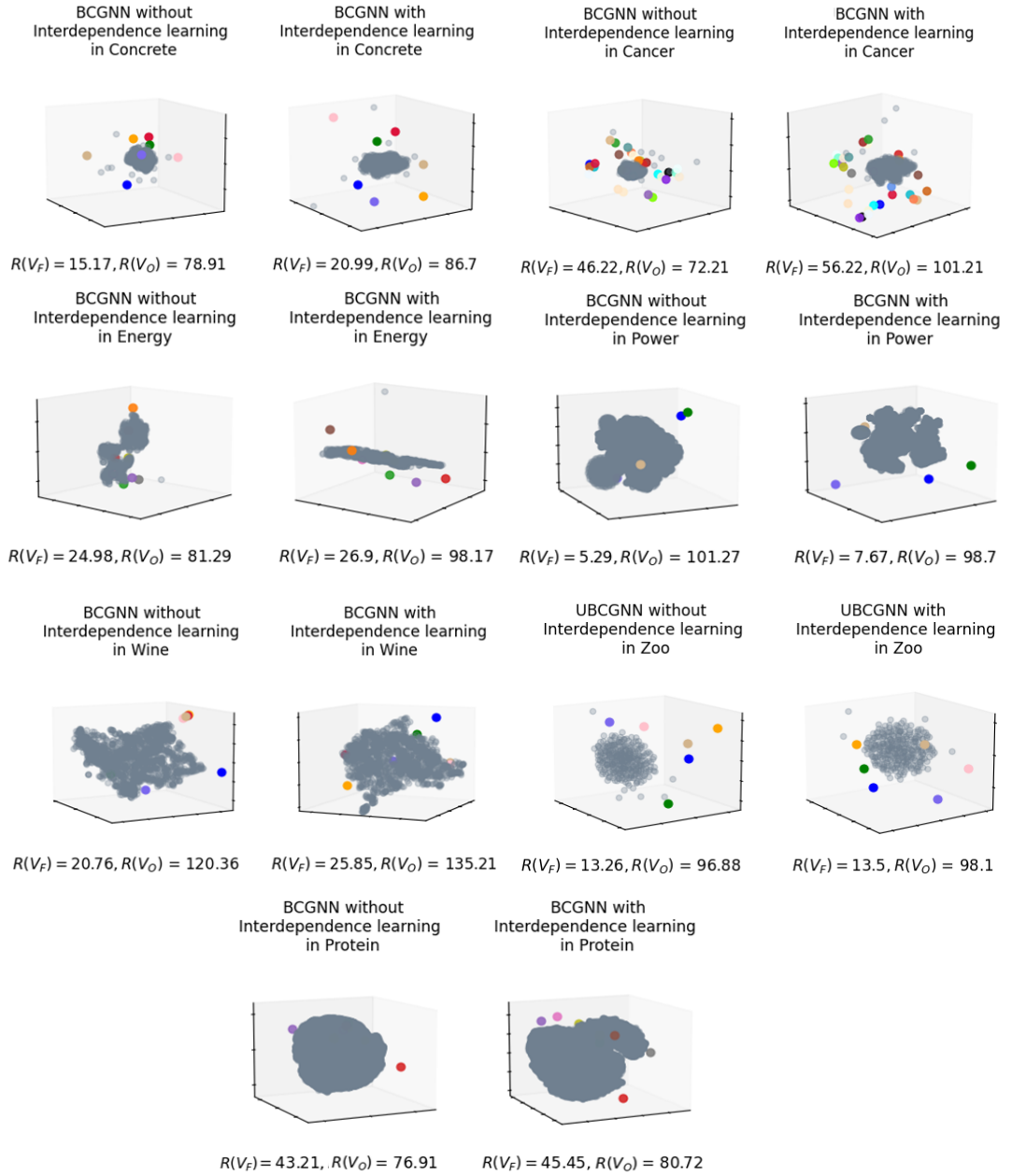


Fig. D6 The sizes of feature embedding space V_F and observation embedding space V_O , which are obtained from the trained BCGNN with/without interdependence structure learning under MNAR by the dimension reduction visualization method t-SNE for all datasets.

References

- [1] Ibrahim, J.G., Molenberghs, G.: Missing data methods in longitudinal studies: a review. *Test* **18**(1), 1–43 (2009)
- [2] Ibrahim, J.G., Chen, M.-H., Lipsitz, S.R., Herring, A.H.: Missing-data methods for generalized linear models: A comparative review. *Journal of the American Statistical Association* **100**(469), 332–346 (2005)
- [3] Brick, J.M., Kalton, G.: Handling missing data in survey research. *Statistical methods in medical research* **5**(3), 215–238 (1996)
- [4] Wooldridge, J.M.: Inverse probability weighted estimation for general missing data problems. *Journal of econometrics* **141**(2), 1281–1301 (2007)
- [5] Sterne, J.A., White, I.R., Carlin, J.B., Spratt, M., Royston, P., Kenward, M.G., Wood, A.M., Carpenter, J.R.: Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *Bmj* **338** (2009)
- [6] Little, R.J., Rubin, D.B.: *Statistical Analysis with Missing Data* vol. 793. John Wiley & Sons, ??? (2019)
- [7] Loh, W.-Y.: Classification and regression trees. *Wiley interdisciplinary reviews: data mining and knowledge discovery* **1**(1), 14–23 (2011)
- [8] Śmieja, M., Struski, L., Tabor, J., Zieliński, B., Spurek, P.: Processing of missing data by neural networks. *Advances in neural information processing systems* **31** (2018)
- [9] You, J., Ma, X., Ding, Y., Kochenderfer, M.J., Leskovec, J.: Handling missing data with graph representation learning. *Advances in Neural Information Processing Systems* **33**, 19075–19087 (2020)
- [10] Kyono, T., Zhang, Y., Bellot, A., Schaar, M.: Miracle: Causally-aware imputation via learning missing data mechanisms. *Advances in Neural Information Processing Systems* **34**, 23806–23817 (2021)
- [11] Zhang, W., Li, G., Tang, J., Li, J., Tsung, F.: Missing data imputation with graph laplacian pyramid network. *arXiv preprint arXiv:2304.04474* (2023)
- [12] Um, D., Park, J., Park, S., Choi, J.Y.: Confidence-based feature imputation for graphs with partially known features. *arXiv preprint arXiv:2305.16618* (2023)
- [13] Gupta, S., Manchanda, S., Ranu, S., Bedathur, S.J.: Grafenne: learning on graphs with heterogeneous and dynamic feature sets. In: *International Conference on Machine Learning*, pp. 12165–12181 (2023). PMLR
- [14] Zhong, J., Gui, N., Ye, W.: Data imputation with iterative graph reconstruction.

- In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 37, pp. 11399–11407 (2023)
- [15] Telyatnikov, L., Scardapane, S.: Egg-gae: scalable graph neural networks for tabular data imputation. In: International Conference on Artificial Intelligence and Statistics, pp. 2661–2676 (2023). PMLR
 - [16] Burgette, L.F., Reiter, J.P.: Multiple imputation for missing data via sequential regression trees. *American journal of epidemiology* **172**(9), 1070–1076 (2010)
 - [17] Van Buuren, S., Groothuis-Oudshoorn, K.: mice: Multivariate imputation by chained equations in r. *Journal of statistical software* **45**, 1–67 (2011)
 - [18] Zhu, H.-T., Lee, S.-Y.: Local influence for incomplete data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **63**(1), 111–126 (2001)
 - [19] Kim, K.-Y., Kim, B.-J., Yi, G.-S.: Reuse of imputed data in microarray analysis increases imputation efficiency. *BMC bioinformatics* **5**, 1–9 (2004)
 - [20] Srebro, N., Rennie, J., Jaakkola, T.: Maximum-margin matrix factorization. *Advances in neural information processing systems* **17** (2004)
 - [21] Cai, J.-F., Candès, E.J., Shen, Z.: A singular value thresholding algorithm for matrix completion. *SIAM Journal on optimization* **20**(4), 1956–1982 (2010)
 - [22] Hastie, T., Mazumder, R., Lee, J.D., Zadeh, R.: Matrix completion and low-rank svd via fast alternating least squares. *The Journal of Machine Learning Research* **16**(1), 3367–3402 (2015)
 - [23] Yoon, J., Jordon, J., Schaar, M.: Gain: Missing data imputation using generative adversarial nets. In: International Conference on Machine Learning, pp. 5689–5698 (2018). PMLR
 - [24] Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.-A.: Extracting and composing robust features with denoising autoencoders. In: Proceedings of the 25th International Conference on Machine Learning, pp. 1096–1103 (2008)
 - [25] Gondara, L., Wang, K.: Multiple imputation using deep denoising autoencoders. *arXiv preprint arXiv:1705.02737* **280** (2017)
 - [26] Muzellec, B., Josse, J., Boyer, C., Cuturi, M.: Missing data imputation using optimal transport. In: International Conference on Machine Learning, pp. 7130–7140 (2020). PMLR
 - [27] Berg, R.v.d., Kipf, T.N., Welling, M.: Graph convolutional matrix completion. *arXiv preprint arXiv:1706.02263* (2017)

- [28] Zhang, M., Chen, Y.: Inductive matrix completion based on graph neural networks. arXiv preprint arXiv:1904.12058 (2019)
- [29] Li, J., Li, J., Liu, Y., Yu, J., Li, Y., Cheng, H.: Deconvolutional networks on graph data. *Advances in Neural Information Processing Systems* **34**, 21019–21030 (2021)
- [30] Spinelli, I., Scardapane, S., Uncini, A.: Missing data imputation with adversarially-trained graph convolutional networks. *Neural Networks* **129**, 249–260 (2020)
- [31] Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y.: Graph attention networks. arXiv preprint arXiv:1710.10903 (2017)
- [32] Chen, J., Chen, H.: Edge-featured graph attention network. arXiv preprint arXiv:2101.07671 (2021)
- [33] Rong, Y., Huang, W., Xu, T., Huang, J.: Dropedge: Towards deep graph convolutional networks on node classification. arXiv preprint arXiv:1907.10903 (2019)
- [34] Hamilton, W.L., Ying, R., Leskovec, J.: *Inductive Representation Learning on Large Graphs* (2018)
- [35] Leskovec, J., Rajaraman, A., Ullman, J.D.: *Mining of Massive Data Sets*. Cambridge university press, ??? (2020)
- [36] Rubin, D.B.: Inference and missing data. *Biometrika* **63**(3), 581–592 (1976)
- [37] Jegelka, S.: Theory of graph neural networks: Representation and learning. arXiv preprint arXiv:2204.07697 (2022)
- [38] Zhang, B., Fan, C., Liu, S., Huang, K., Zhao, X., Huang, J., Liu, Z.: The expressive power of graph neural networks: A survey. arXiv preprint arXiv:2308.08235 (2023)
- [39] Ma, Y., Derksen, H., Hong, W., Wright, J.: Segmentation of multivariate mixed data via lossy data coding and compression. *IEEE transactions on pattern analysis and machine intelligence* **29**(9), 1546–1562 (2007)
- [40] Crysdian, C.: The behaviour of rank correlation coefficients for incomplete data. *Research in Mathematics* **9**(1), 2107793 (2022)
- [41] Mukaka, M.M.: A guide to appropriate use of correlation coefficient in medical research. *Malawi medical journal* **24**(3), 69–71 (2012)
- [42] Schober, P., Boer, C., Schwarte, L.A.: Correlation coefficients: appropriate use and interpretation. *Anesthesia & analgesia* **126**(5), 1763–1768 (2018)

- [43] You, J., Ying, Z., Leskovec, J.: Design space for graph neural networks. *Advances in Neural Information Processing Systems* **33**, 17009–17021 (2020)
- [44] Zhou, F., Li, T., Zhou, H., Zhu, H., Jieping, Y.: Graph-based semi-supervised learning with non-ignorable non-response. *Advances in Neural Information Processing Systems* **32** (2019)