

BClean: A Bayesian Data Cleaning System

Jianbin Qin^{1,2}, Sifan Huang¹, Yaoshu Wang², Jing Zhu¹, Yifan Zhang¹,
Yukai Miao³, Rui Mao^{1,2}, Makoto Onizuka⁴, Chuan Xiao^{4,5}

¹Shenzhen University

²Shenzhen Institute of Computing Sciences

³Zhongguancun Laboratory

⁴Osaka University

⁵Nagoya University

{qinjianbin,mao}@szu.edu.cn,{huangsifan2020,2300271063,zhangyifan2021}@email.szu.edu.cn

yaoshuw@sics.ac.cn,miaoyk@zgcclab.edu.cn,{onizuka,chuanx}@ist.osaka-u.ac.jp

ABSTRACT

There is a considerable body of work on data cleaning which employs various principles to rectify erroneous data and transform a dirty dataset into a cleaner one. One of prevalent approaches is probabilistic methods, including Bayesian methods. However, existing probabilistic methods often assume a simplistic distribution (e.g., Gaussian distribution), which is frequently underfitted in practice, or they necessitate experts to provide a complex prior distribution (e.g., via a programming language). This requirement is both labor-intensive and costly, rendering these methods less suitable for real-world applications. In this paper, we propose BClean, a Bayesian Cleaning system that features automatic Bayesian network construction and user interaction. We recast the data cleaning problem as a Bayesian inference that fully exploits the relationships between attributes in the observed dataset and any prior information provided by users. To this end, we present an automatic Bayesian network construction method that extends a structure learning-based functional dependency discovery method with similarity functions to capture the relationships between attributes. Furthermore, our system allows users to modify the generated Bayesian network in order to specify prior information or correct inaccuracies identified by the automatic generation process. We also design an effective scoring model (called the compensatory scoring model) necessary for the Bayesian inference. To enhance the efficiency of data cleaning, we propose several approximation strategies for the Bayesian inference, including graph partitioning, domain pruning, and pre-detection. By evaluating on both real-world and synthetic datasets, we demonstrate that BClean is capable of achieving an F-measure of up to 0.9 in data cleaning, outperforming existing Bayesian methods by 2% and other data cleaning methods by 15%. Our source code is available at <https://github.com/yysl88/BClean>.

1 INTRODUCTION

Data cleaning is an essential step in preprocessing data for downstream applications such as data analysis and machine learning (ML) model construction. Existing data cleaning solutions utilize user-defined rules, outlier detection techniques, crowdsourcing, or knowledge bases to identify errors in the data [1, 28, 42], and subsequently carry out data repairing [35, 41, 50]. The majority of research on data cleaning has been dedicated to error correction using discriminative models based on integrity constraints [15, 31, 62], external data [14, 25, 51, 58, 67], statistical approaches [43, 45], ML techniques [61], or hybrid methods [7, 41]. Despite these approaches, certain challenges persist. For instance, some solutions based on integrity constraints often encounter difficulties in real-world applications, as they require expert input for rule creation.

Solutions that rely on external data incur substantial costs associated with expert information collection. Moreover, ML methods face potential confusion of learning models when generating feature vectors from noisy data.

To address these shortcomings, probabilistic methods [34, 37] have been proposed for data cleaning in a data-driven manner. These methods employ probabilistic inference [46] for data cleaning, such as using probabilistic graphical models (PGM). Furthermore, probabilistic programming languages (PPL) [8, 10, 44] have been developed as powerful tools for describing and executing probabilistic models. Another typical method is HoloClean [50], in which probabilistic inference acts as a feature generator to generate the degrees of validity of denial constraints (DCs), one type of dependency rules. Configuring rules in HoloClean is simpler than in PPL. Despite achieving robust cleaning performance in a few empirical evaluations, HoloClean is a semi-supervised method that requires manually annotated data. Other semi-supervised methods, such as Raha [42] and Baran [41], feature automatic rule discovery. They need labels for around 20 tuples to perform error detection and correction. However, error propagation may occur from detection to correction, despite ease of deployment. Currently, probabilistic inference-based data cleaning remains an active area of study, as PGMs offer natural advantages when modeling dirty data as generative models [29].

Recently, Bayesian methods, a subset of probabilistic methods, have been explored and have demonstrated promising results [35]. For instance, existing studies [16, 29, 66] employ Bayesian networks to correct errors. These methods encompass Bayesian network construction and inference. The network construction involves structure learning and parameter estimation from the observed data, whereas the inference aims to determine the most probable value to fill in each data cell, given other attributes of the same tuple and the learned data distribution. If the inferred value differs from the original value in the cell, the inferred value replaces the original one, serving as a correction for the cell.

Motivation. Existing Bayesian methods face three key challenges: (1) a high dependence on accurate domain knowledge from experts, which is required for PPL [35], and necessitates user interaction and domain knowledge editing, (2) inefficiency of inference using Bayesian networks [16, 29], in which variable elimination following the topological order of the nodes incurs significant computational cost, and (3) the presence of errors in the constructed Bayesian network originating from noisy data, which may propagate into subsequent inferences, resulting in inferior inference performance.

Example 1: Table 1 shows a Customer table that includes 6 tuples with 9 attributes. Most existing methods can accurately identify and

Table 1: Customer table.

Tid	Name	Department	Jobid	City	State	ZipCode	InsuranceCode	InsuranceType
1	Johnny.R	315 w hickory st	25676000	sylacauga	CA	35150	2567600035150	
2	Johnny.R	400 northwood dr	25676x00	sylacauga	KT	35150	2567600035150	Normal
3	Johnny.R	315 w hicky st	25676000	sylacauga	CA	35150	2567600035150	Normal
4	Henry.P	400 northwood dr	25600180	centre	KT		2560018035960	Low
5	Henry.P	400 nprthwood dr	25600180	centre	NY	3960	25600v5960	High
6	Henry.P		25600180	centre	KT	35960		Low

correct errors in Tuples 1 – 3, given that the context is relatively clean and complete. Functional dependency (FD) learning [65] can be employed to discover the FD $InsuranceCode \rightarrow InsuranceType$ to correct the missing value “NULL \rightarrow Normal” in the first line, and another FD $ZipCode \rightarrow State$ to correct the error “KT \rightarrow CA” in the second line. Repairs such as “25676x00 \rightarrow 25676000” and “315 w hicky st \rightarrow 315 w hickory st” can be identified using the context [41]. In contrast, Tuples 4 – 6 in Table 1 present more of a challenge due to the presence of numerous critical errors. The values for Department, InsuranceCode, InsuranceType, and State cannot be definitively inferred due to little evidence from Name, Jobid, and ZipCode. Due to the lack of observations, the errors “400 nprthwood dr”, “NY”, “3960”, “25600v5960”, and “High” may be perceived as correct values by probabilistic models.

To address these problems, many existing data cleaning methods incorporate external data [14, 50]. In Bayesian methods, domain knowledge plays a crucial role in encoding data distribution and Bayesian network structure. PClean [35] is a state-of-the-art Bayesian method, whose network construction and the prior knowledge are hand-crafted. This requires specification of data types, compliant distributions, and possible noises (e.g., uncommon symbols or Gaussian noise). However, authoring the PPL code in PClean introduces a steep learning curve and high user cost. For example, PClean requires users to precisely partition Table 1 into four parts: $P_1 = \{Name \sim Dist(\theta_1), Name \sim Dist(\theta_2), Jobid \sim Dist(\theta_3)\}$; $P_2 = \{City \sim Dist(\theta_4), State \sim Dist(\theta_5), ZipCode \sim Dist(\theta_6)\}$; $P_3 = \{InsuranceCode \sim Dist(\theta_7), InsuranceType \sim Dist(\theta_8)\}$; and $P_4 = \{P_1, P_2, P_3, error_dist(\theta_9)\}$, where $A \sim Dist(\theta)$ in each part denotes an attribute A that follows a distribution with parameter θ . It is hard for users to specify such partition, and the specification itself is error-prone.

Our approach. In this paper, we introduce BClean, an unsupervised Bayesian data cleaning system. Given a noisy relational dataset, BClean performs error correction in two stages: the construction stage and the inference stage. In the construction stage, BClean automatically constructs a Bayesian network from the dataset and provides a user interaction function that optionally collects lightweight domain knowledge from users to fine-tune the Bayesian network. In the inference stage, BClean leverages the dataset and any constraints specified by users to develop a compensatory scoring model, which acts as a supplementary model providing additional statistical information for inference. Notably, it approximates the value of a term in the formula of Bayesian inference, thereby addressing the problem of error amplification caused by inaccurate Bayesian network construction from dirty data. BClean processes each cell in the dataset, using the constructed Bayesian network and the compensatory scoring model to compute the probability of candidate repair options, based on the evidence from other attributes of the

same tuple. Furthermore, BClean incorporates a series of optimization techniques to boost efficiency, including graph partitioning, domain pruning, and pre-detection. An experimental evaluation on real-world and synthetic datasets shows that BClean can achieve an F-measure of up to 0.9 in data cleaning, outperforming existing Bayesian methods by 2% and other data cleaning methods by 15%.

In BClean, akin to PClean but more straightforward, users only need to define constraints by specifying simple expressions (e.g., categories of values, minimum/maximum lengths of attributes, minimum/maximum values of numerical attributes, whether null values are allowed, and regular expressions) instead of adhering to a specific distribution. As a result, BClean incurs substantially lower user costs than existing solutions (e.g., PClean), which require users to learn PPL to inject prior knowledge. Importantly, the user constraints supported in BClean are not limited to the aforementioned expressions but can be any function that returns a binary output, thus covering a wide range of constraint forms such as dependency rules (FDs and DCs), arithmetic expressions, and even deep neural networks. Nevertheless, our experiments show that using simple expressions is sufficient to deliver satisfactory data cleaning results.

Our contributions are summarized as follows:

- We propose a Bayesian-inference-based data cleaning system that repairs errors in datasets using a Bayesian network.
- We adapt an existing FD discovery method to automatically generate a Bayesian network, and provide users with an interface to optionally modify the Bayesian network to specify prior knowledge of the dataset.
- We design a compensatory scoring model for Bayesian inference, leveraging user-specified constraints to filter obviously incorrect values and rank candidate repair options.
- We introduce a series of approximation techniques, including graph partitioning, domain pruning, and pre-detection, to enhance the efficiency of data cleaning.
- We conduct experiments on real-world and synthetic datasets. The results demonstrate the effectiveness and the efficiency of our system and its competitiveness with alternative solutions, especially those methods based on Bayesian inference.

The remainder of this paper is organized as follows. Section 2 introduces preliminaries. Section 3 offers an overview of BClean and its data-driven modeling. Sections 4 – 6 present the technical details of BClean. Section 7 reports the experimental results. Section 8 reviews related work. Section 9 concludes the paper.

2 PRELIMINARIES

The aim of BClean is to repair erroneous data in a structured dataset D by employing a Bayesian network (BN, a.k.a. Bayes network or belief network), thereby producing a cleaner version of the dataset.

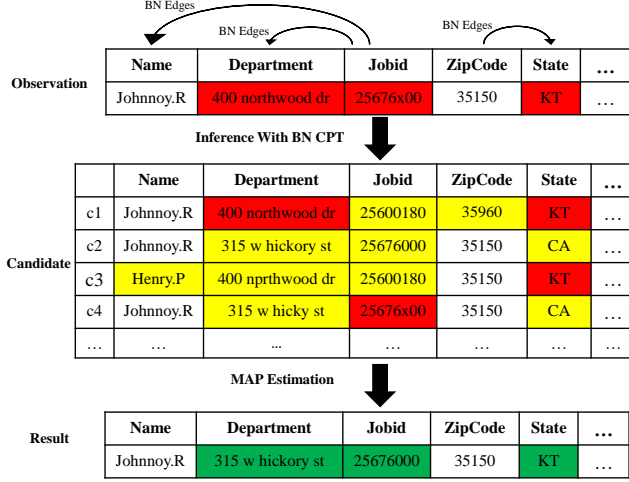


Figure 1: A running example of BClean. Red, yellow, and green represent erroneous, candidate, and clean cell values.

The observed dataset D comprises n tuples, $\{T_1, \dots, T_n\}$, and is characterized by m attributes: $\mathcal{A} = \{A_1, \dots, A_m\}$. The domain of attribute A_j is denoted as $dom(A_j)$. $T_i[A_j]$ represents the observation of the j -th attribute in the i -th tuple. To repair the dataset, BClean considers candidate values $c \in dom(A_j)$ and infers the most probable candidate value c^* to replace $T_i[A_j]$ for all $i \in [1, n]$ and $j \in [1, m]$. Moreover, candidate values must satisfy the constraints specified by the user; e.g., the ZIP code in the US must be five digits. We employ a user constraint (UC) to verify if the input (a cell, tuple, or a dataset) adheres to the user-specified constraints. A UC is a function $UC(\cdot)$ that returns 1 if the input satisfies the user constraints, or 0 otherwise. It can be any function that returns a binary output, such as rules (e.g., DCs and FDs), arithmetic expressions, regular expressions, or even deep neural networks. As the interaction should be easy to use, we primarily focus on the following UCs in this paper: (1) minimum/maximum attribute lengths (or minimum/maximum values for numerical attributes), (2) non-null constraints, and (3) simple regular expressions for digits and dates.

The above UCs are straightforward. The first pertains to attribute statistics, the second mandates value assignment, and the third specifies the attribute’s format. Even for users unfamiliar with regular expressions, numerous online tools exist for generating them from examples, such as the one detailed in [6], with its online demonstration available at [5]. The underlying rationale is that these UCs, serving as prior knowledge about the dataset, do not necessitate database expertise. Therefore, it is unnecessary to presume user familiarity with databases, particularly advanced techniques used in HoloClean, including conditional FDs [9] and metric FDs [32]. Additionally, users are not required to label data for many tuples (approximately 40, as in Raha+Baran). This approach makes our solution more broadly applicable. For instance, improving data preprocessing is a relatively relevant challenge in economics and social science research [3]. While model accuracy is crucial in these fields, we cannot expect researchers to be versed in concepts like functional dependencies, to devote substantial time to acquiring

these database skills, or to drill down to the dataset to label the correctness of certain values in a record-by-record manner.

By inferring the value of each cell in D , we obtain the cleaned dataset D^* . Figure 1 depicts an example of applying BClean to a dataset. Although a dataset typically contains many tuples, we illustrate only one tuple in this example. The nodes in the BN represent attributes. The edges in the BN represent the dependency relationships between attributes. Candidate values are generated for each cell, and their probabilities are computed by maximum a posteriori (MAP) estimation using the BN.

A BN is a type of Bayesian method and is presently one of the most effective theoretical models in the field of uncertain knowledge representation and inference [46]. A BN is a directed acyclic graph (DAG) (N, E, θ) , composed of a set of nodes N representing random variables, a set of directed edges E indicating the conditional dependencies between random variables, and a set of conditional probability tables (CPTs) θ that weight the edges to express the strength of conditional dependency. For each random variable, its probability can be computed using the probabilities of its node’s parents. In case of no parent node, a prior probability is expressed by the prior information. Given an observed dataset D , the structure of a BN and its CPTs can be learned, with nodes representing attributes and edges representing dependency relationships. Figure 2 illustrates the structure of a BN learned from the dataset in Table 1. The probability of a tuple $T = (t_1, \dots, t_m)$ is

$$\Pr[t_1, \dots, t_m] = \prod_{A_i \in N} \Pr[T[A_i] = t_i | T[A_j] = t_j, j \in Ans(i)],$$

where $Ans(i)$ represents the set of ancestor nodes influencing A_i . For nodes without any ancestor nodes, $\Pr[T[A_i] = t_i | T[A_j] = t_j, j \in Ans(i)]$ is determined by the prior probability of A_i , which can be inferred from D .

3 MODELING IN BCLEAN

An overview of BClean is depicted in Figure 2. BClean accepts as input a dataset D as observation, along with UCs provided by users. BClean executes data cleaning in two steps:

Probabilistic modeling. We employ graphical lasso [60] to generate a covariance matrix and a directed graph structure learning method [49] to generate a BN. Inference on an erroneous dataset D can exacerbate deviation in the results. To counter this issue, we leverage a compensatory score. This score considers value frequency (i.e., for each attribute A_i , we collect the observations for each value in $dom(A_i)$ as its frequency) and pairwise attribute correlation.

Inference grounding and pruning. We generate candidate values for data cleaning by iterating through the domain values of each attribute and compute the probabilities of the candidate values. Considering the number of candidate values can be vast, we curtail unnecessary candidate generation through three techniques. First, we partition the BN by capitalizing on the Markov property [46]. Second, based on the partitioned BN, we prune the domain of each attribute and discard values that obviously do not conform to the domain semantics. Third, we use the compensatory score to perform preliminary detection and inference for the variables likely to be correct.

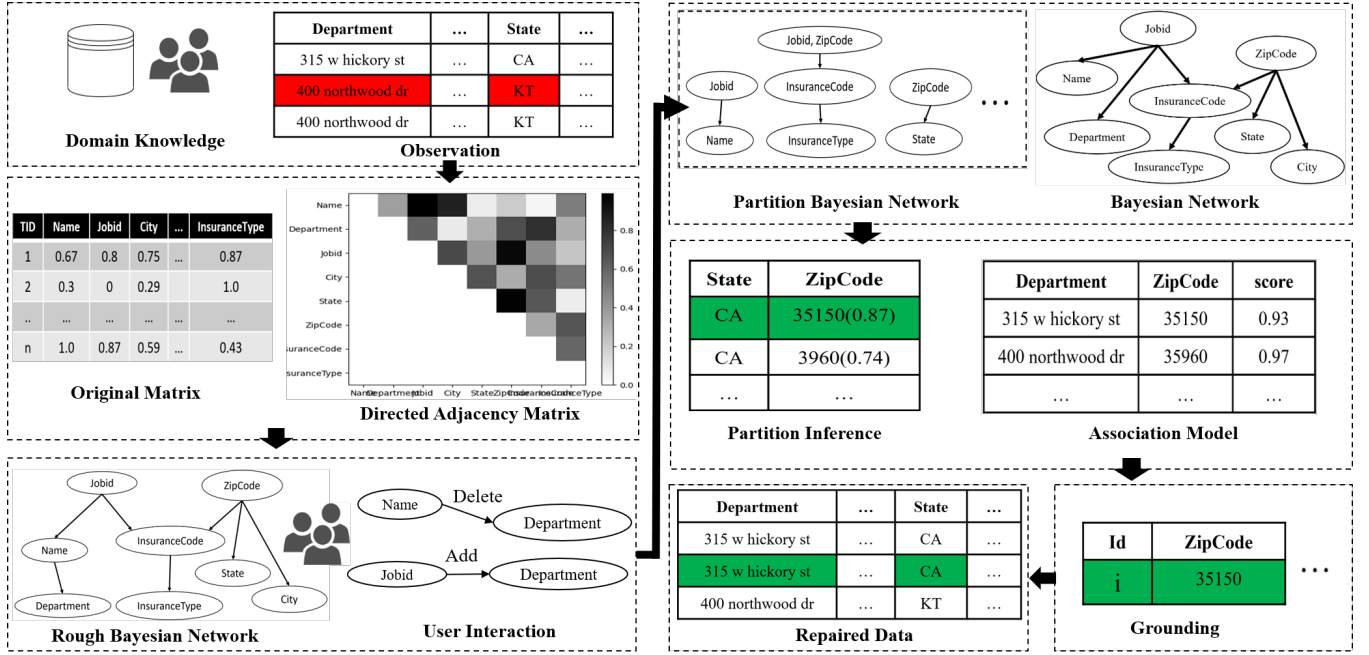


Figure 2: An overview of the BClean framework.

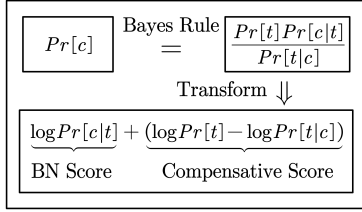


Figure 3: BClean's data-driven model.

We represent the data cleaning task as an MAP problem to determine the optimal data distribution. Figure 3 illustrates the modeling procedure of BClean. Given D , our objective is to solve an MAP problem to derive a cleaned dataset D^* . For each T_i in D and each attribute A_j , we infer the most probable candidate value c^* in $dom(A_j)$. To calculate the probability of each candidate value, we leverage the other attributes of T_i and apply the following transformation via Bayes' rule:

$$c^* = \arg \max_{c \in dom(A_j)} \frac{Pr[t]Pr[c|t]}{Pr[t|c]}, \text{ subject to } UC(c) = 1,$$

where t denotes the observed values of attributes other than A_j in $T[i]$, i.e., $t = T_i[A_1, \dots, A_{j-1}, A_{j+1}, \dots, A_m]$. By transforming it to logarithmic form, we obtain

$$c^* = \arg \max_{c \in dom(A_j)} (\log Pr[c|t] + \log Pr[t] - \log Pr[t|c]), \quad (1)$$

subject to $UC(c) = 1$.

The terms in the above equation can be divided into two parts: $\log Pr[c|t]$ and $\log Pr[t] - \log Pr[t|c]$. The first term can be obtained

Algorithm 1: BClean

Input : Observed dataset D with n rows and m columns, Bayesian network BN .
Output : Cleaned dataset D^* .

```

1 for  $i = 1$  to  $n$  do
2   for  $j = 1$  to  $m$  do
3      $c^* \leftarrow T_i[A_j]$ ;
4      $p^* \leftarrow \log(BN[A_j](c^*)) + \log(CS[A_j](c^*))$ ;
5     for  $c \in dom(A_j)$  such that  $UC(c) = 1$  do
6        $p \leftarrow \log(BN[A_j](c)) + \log(CS[A_j](c))$ ;
7       if  $p > p^*$  then
8          $c^* \leftarrow c, p^* \leftarrow p$ ;
9      $D_i^*[A_j] \leftarrow c^*$ ; /* clean ( $i$ -th row,  $j$ -th column) */;
10 return  $D^*$ ;
    
```

from the BN using the nodes having an edge to A_j . While the second term cannot be directly computed from the BN, we compute it with a compensatory scoring model, whose details will be presented in the next section.

Algorithm 1 outlines the procedure of BClean. It iterates through all rows and columns in the observed dataset D and infers the most probable value for each cell. For each cell, c^* is initialized as the original value $T_i[A_j]$. Then, for each candidate value that satisfies UCs, we compute its probability using the BN and the compensatory scoring model (Equation 1), and update c^* if the probability is higher. All the c^* values constitute the cleaned dataset D^* , which is returned by the algorithm.

4 BAYESIAN NETWORK CONSTRUCTION

We introduce our technique for constructing a BN. The BN partitioning technique will be described in Section 6.

The problem of constructing a BN from a given dataset is NP-hard [12]. Therefore, an exhaustive search algorithm that traverses every structure is not practical for large datasets. Many existing approaches resort to heuristics to find a feasible solution. For instance, hill-climbing-based approaches, such as MMHC [55] provided in the Pgmmy toolkit [4], add one edge at a time and evaluate its score to determine each edge and its direction. However, such approaches often converge to a local optimum, limiting their effectiveness. Other typical approaches include tree search, which necessitates specifying the root state, and the PC algorithm [53], which requires a conditional independence hypothesis given by the user.

In BClean, we construct a BN from the dataset D , as depicted in Figure 2. Automatic BN construction is challenging because it either requires accurate data or prior knowledge of the network (e.g., as specified using PPL [35]). However, D is often unclean and sparse. Our objective is to create a supportive system that obviates the need for users to initiate BN construction from scratch. This makes the aforementioned heuristic approaches ill-suited to our use case. To mitigate this, we leverage a structure learning algorithm that supports domain knowledge to construct an approximate BN. Our strategy involves extending the FDX method [65] with error tolerance to generate an *inverse covariance matrix* using *graphical lasso* [60], and then employing the inverse covariance matrix to create a BN skeleton. Given that the BN skeleton could be noisy and incomplete, we offer a user-interaction feature enabling users to view and manually adjust the BN. It is worth noting that there are alternative methods to construct a BN from D . However, we have chosen the solution outlined above due to its ability to tolerate errors in D and its superior empirical performance. Next, we describe the details.

Our method, based on in statistical modeling, operates directly on D and models the distribution by which clean data is generated. To model the relationships between attributes, we utilize the FDs within the dataset. An FD $X \rightarrow Y$ states that $T_i[Y] = T_j[Y]$, if $T_i[X] = T_j[X]$. Given the presence of errors in the dataset, the application of strict FDs lacks flexibility. Therefore, we propose to soften the FDs by introducing a similarity measure: for any pair of tuples $T_i, T_j \in D$, we denote $Sim(T_i[X], T_j[X])$ as the similarity between two values. The similarity ranges between 0 and 1, and serves as probability in our modeling. For numerical data, we use $\frac{|T_i[X] - T_j[X]|}{(|T_i[x]| + |T_j[x]|)/2}$, and for string data, we use unit-cost edit distance normalized by string lengths as follows.

$$Sim(T_i[X], T_j[X]) = 1 - \frac{2 \cdot ED(T_i[X], T_j[X])}{len(T_i[X]) + len(T_j[X])},$$

where ED denotes unit-cost edit distance (i.e., Levenshtein distance) and len denotes string length. For example, as illustrated in the Department attributes of Tuples 1 and 3 in Figure 1, the features obtained here are represented by 1 because they refer to the same entity. If an FD is applied, it would be 0, whereas our softened method reports a similarity of 0.86, thereby tolerating errors in the dataset.

We compute pairwise similarity for the attributes in each tuple, considering these similarity values as observations from a multivariate Gaussian distribution. The graphical lasso [60] is subsequently utilized to calculate the covariance matrix Σ of the underlying distribution. We decompose the inverse covariance matrix $\Theta = \Sigma^{-1}$, adhering to the method typically used in learning the structure of linear models [40, 49] and employed in FDX [65]:

$$\Theta = \Sigma^{-1} = (I - B)\Omega(I - B)^T,$$

where Θ denotes the inverse covariance matrix for pairwise attributes in \mathcal{A} , I is the identity matrix, and B is the autoregression matrix of the model, which is exactly the adjacency matrix that stores the weights of the edges of the BN skeleton. Moreover, we set a weight threshold, thereby retaining only edges with weights exceeding the threshold.

For example, for dataset D in Table 1, we compute pairwise similarities, as depicted in the second part of Figure 2 labeled original matrix. Following this, we compute the inverse covariance matrix and decompose it to obtain the adjacency matrix of the BN skeleton. (*ZipCode*, *InsuranceCode*) becomes a directed edge since its weight in the adjacency matrix surpasses the threshold.

Upon constructing the BN skeleton, we conduct parameter learning to estimate the joint probability of attributes and derive the CPTs concerning the observations. It is important to note that during the construction of the BN skeleton, we are oblivious to any errors in the observations, and our BN construction models errors as part of the distribution.

Given that the BN skeleton may be noisy, we offer an interface through which users can view and manually modify the automatically constructed skeleton to generate the BN employed for inference. Users can add or remove edges in the BN, as depicted in Figures 2 (f) – (g), or merge nodes, as shown in Figures 2 (g) – (h). For nodes to be merged, if all of them have an edge from/to node A_j , the edges from/to A_j will be merged into one, as displayed in the incoming edge to *InsuranceCode* in Figure 2 (h). Other edges from/to the nodes to be merged will be removed after merging. Since edges are associated with CPTs, the CPTs need to be updated if users modify the BN. For efficiency, we only recalculate the CPTs for the attributes involved in the modification instead of all attributes.

Remarks. To use the FDX method, we first sort tuples according to each attribute, and only compute similarities and check equality within two adjacent tuples. As such, we do not need to compute each tuple pair in D . We construct the Bayesian network by using the FDX method and extend it to support fuzzy matching, e.g., edit similarity. The time complexity is $(nm \log n + nm + r|B|)$, where r is the number of iterations of graphical lasso learning, and $|B|$ is the size of the autoregression matrix.

5 BAYESIAN INFERENCE WITH COMPENSATORY SCORE

Upon construction of the BN, we design an inference method incorporating a compensatory strategy to identify the most probable candidate value for each attribute of a tuple T . Traditional approaches to BN inference involve the iterative cleaning of data via MAP estimation. However, such a method is often ineffective as most real-world datasets are unclean and the BNs constructed

from them lack accuracy. If a repair is incorrect, the error propagates to the next step, rendering subsequent repairs based on it likely invalid. We refer to this problem as error amplification. Recall that we partition the logarithmic form of $\log \Pr[c]$ into two terms, $\log \Pr[c|t]$ and $\log \Pr[t] - \log \Pr[t|c]$, where $\log \Pr[c|t]$ can be computed through BN inference via CPTs. To minimize the errors produced by $\log \Pr[c|t]$, we consider the second term as a compensatory score $\text{Score}_{\text{comp}}$, i.e., $\text{Score}_{\text{comp}} = \log \Pr[t] - \log \Pr[t|c]$.

Example 2: Consider the BN of the relational data in Table 1, where we aim to populate the missing value for the attribute *Department* of t_6 . Owing to the unclean data, we derive $\Pr[\text{"400 nrpthwood dr"}|t_6] = 0.54$ and $\Pr[\text{"400 northwood dr"}|t_6] = 0.51$. Without $\text{Score}_{\text{comp}}$, "400 nrpthwood dr" would be erroneously filled into the cell for *Department* of t_6 .

Computing $\text{Score}_{\text{comp}}$ directly is difficult, as it necessitates the knowledge of $\Pr[c^*]$, thereby inverting the dependency between c^* and t . Given that $\Pr[t]$ is a constant value once the BN is constructed and fixed, a lower value of $\Pr[t|c]$ suggests $\Pr[c]$ is closer to $\Pr[c^*] = 1$. We approximate the computation of the compensatory score as the correlation between c and t :

$$\text{Score}_{\text{comp}}(c, t) \approx \text{Score}_{\text{corr}}(c, t).$$

To explain this approximation, intuitively, if a candidate value $c \in \text{dom}(A_i)$ of a tuple T is c^* , c is likely to exhibit a strong dependency to other attributes, i.e., coexist with other values frequently. We assume that erroneous values in real-world datasets are rare. For example, consider an error c' which only appears in one tuple T . In this case, $\Pr[t|c] = 1$ and $-\log \Pr[t|c]$ attains the minimum value, i.e., 0. Thus, modeling the correlation between c and t , i.e., $\text{Score}_{\text{corr}}$, aligns with the same objective as $\text{Score}_{\text{comp}}$.

Inspired by Bayeswipe [16], we adopt a straightforward yet effective idea to compute $\text{Score}_{\text{corr}}$ by counting the occurrences of (c, t) . Accordingly, $\text{Score}_{\text{corr}} = \frac{\text{count}(c, t)}{|D|}$. Since this scoring suffers from sparsity and tends to be inaccurate for infrequent observations, to mitigate this, we separately consider the correlation between each attribute value of t and c and accumulate them. If more values in t have a high correlation with c , then $\text{Score}_{\text{corr}}(c, t)$ is likely to be significant. Since $c = T_i[A_j]$ and $t = T_i[A_1, \dots, A_j - 1, A_{j+1}, \dots, A_m]$, we define $\text{Score}_{\text{corr}}$ as follows:

$$\text{Score}_{\text{corr}}(c, t, A_j) = \sum_{A_k \in \mathcal{A} \setminus \{A_j\} \wedge e=t[A_k]} \text{corr}(c, e, A_j, A_k), \quad (2)$$

where $\text{corr}(c, e, A_j, A_k)$ denotes the correlation between value c of attribute A_j and value e of attribute A_k . Due to the presence of noise, instead of directly counting the occurrence of (c, e) as their correlation such that $\text{corr}(c, e, A_j, A_k) = \frac{\text{count}(c, e, A_j, A_k)}{|D|}$, we incorporate the UCs provided by users and introducing the concept of confidence (conf) for each tuple T_i . We assign different weights (i.e., confidences) to tuples. A tuple T_i with a higher confidence value suggests that it is more accurate, and the correlations among its values are more reliable. We define $\text{conf}(T_i)$ as follows, where $UC(\cdot)$ is a function to check each attribute value against user constraints

Algorithm 2: compensatory score computation

Input: Observed dataset D , user constraints $UC(\cdot)$, confidence threshold τ , penalty parameter β ;
Output: compensatory scores corr .

```

1 corr = {};
2 for  $T \in D$  do
3    $p := 0$ ;
4   Compute  $\text{conf}(T)$  using Equation 3;
5   for  $A_i \in A$  do
6      $c := T[A_i]$ ;
7     for  $A_j \in A/A_i$  do
8        $v := T[A_j]$ ;
9       if  $\text{conf} \geq \tau$  then
10        |  $\text{corr}[\langle c, v, A_i, A_j \rangle] = \text{corr}[\langle c, v, A_i, A_j \rangle] + 1$ ;
11       else
12        |  $\text{corr}[\langle c, v, A_i, A_j \rangle] = \text{corr}[\langle c, v, A_i, A_j \rangle] - \beta$ ;
13 return corr;
```

and output either 0 or 1:

$$\text{conf}(T_i) = \max \left\{ 0, \frac{\sum_{e \in T_i} \mathbf{1}_{\{UC(e)=1\}} - \lambda \cdot \sum_{e \in T_i} \mathbf{1}_{\{UC(e)=0\}}}{|T_i|} \right\}, \quad (3)$$

where λ is a non-negative parameter used to penalize the erroneous values that violate the UCs, and $\mathbf{1}$ is the indicator function. Then, we define $\text{corr}(c, e, A_j, A_k)$ as

$$\text{corr}(c, e, A_j, A_k) = \frac{\sum_{T \in \Omega} (\mathbf{1}_{\text{conf}(T) \geq \tau} - \beta \mathbf{1}_{\text{conf}(T) < \tau})}{|D|},$$

$$\Omega = \{T' | T' \in D \wedge (c, e) = (T'[A_j], T'[A_k])\}.$$

In the above equation, if the confidence of a tuple $T \in D$ is no less than a predefined constant τ , T is considered reliable and the combinations of its attribute values (i.e., $(T[A_j], T[A_k])$) contribute to the computation of corr . Otherwise, T is considered unreliable, and we impose a penalty β such that corr would decrease once T is taken into account.

Since we incorporate UCs in Equation 3, more violations of UCs will reduce $\text{Score}_{\text{corr}}$. It is important to note that the inference process is a competition among the candidate values in the domain, comparing their scores. In other words, the relative order is significant, not the scores themselves. Therefore, when using $\text{Score}_{\text{corr}}$, it is not necessary to require $\text{Score}_{\text{corr}}$ to be equal to $\text{Score}_{\text{comp}}$.

To efficiently compute $\text{Score}_{\text{corr}}$ for all attributes, we design a compensatory score computation algorithm (Algorithm 2). We scan each attribute value in D and compute the correlation with other attribute values (lines 2 – 12). Specifically, we check each tuple $T \in D$ (line 2) and compute its confidence $\text{conf}(T)$ using Equation 3. Then, we enumerate all pairs of attributes (A_i, A_j) with two cases: if $\text{conf}(T)$ is no less than the predefined threshold τ , we update corr by accumulating the counts of (A_i, A_j) ; otherwise, we impose a penalty β to decrease corr . Finally, we compute and store the correlation values of all pairs of values in D , and $\text{Score}_{\text{corr}}$ could be readily computed using Equation 2.

Example 3: Continuing with Example 2, we set $\lambda = 0.25$, $\beta = 2$, and $\tau = 0.75$. Consider a $UC(\cdot)$ as a spell checker that returns 1 if the

spelling is correct and 0 otherwise. $c = 400$ northwood dr in t_5 violates the UC, and its $\text{conf}(t_5) = 0 < \tau$ and $\text{Score}_{\text{corr}}(c, t_5, \text{Department}) = -0.31$. $c = 400$ northwood dr in t_6 does not violate the UC, and its $\text{conf}(t_6) = 0.79 \geq \tau$. $c = 400$ northwood dr in t_4 does not violate the UC, and its $\text{Score}_{\text{corr}}(c', t_4, \text{Department}) = 0.13$. Eventually, we have $\Pr[400 \text{ northwood dr} | t_6] = 0.64$ and $\Pr[400 \text{ nprthwood dr} | t_6] = 0.23$. Upon executing Bayesian inference with a compensatory score, the Department of t_6 is filled with 400 northwood dr.

Remarks. While our idea is inspired by Bayeswipe [16], its realization is completely different. Instead of focusing on total loss as Bayeswipe does, our correlation-based score $\text{Score}_{\text{corr}}$ is derived from the co-occurrence of attribute-value pairs in D and the confidences of tuples. In real-world data, clean data constitute the majority and exhibit dependency and correlation. If a tuple is clean, its correlation score is high as it is likely to share common attribute values with a few other tuples. Conversely, if a tuple contains errors, it is less likely to share common values with other tuples, resulting in a very low correlation score for t observed. Thus, we approach this as a missing value prediction task in the correlation-based score. For a tuple $T_i \in T$, the value of every cell $T_i[A_j] | 0 \leq i \leq m$ is more likely to be clean for higher co-occurrence times with the remaining observations $T_j[A_j] | 0 \leq j \leq m \wedge j \neq i$ of the column. We first create a co-occurrence dictionary where all candidate values c are stored along with the count of their co-occurrence with every possible value from all other columns. For a specific tuple t , $\text{Score}_{\text{corr}}(c, t)$ will collect the count of cell observations from different columns shown in T_i as features. The norm of the feature is computed as a weighted score, where the weights are the counts of co-occurrence of (c, t) . The distance between an observation and a candidate value is matched with the weighted score. Candidate values with higher scores have a greater likelihood of being the final ground truth. The time complexity of Bayesian inference is dominated by compensatory score computation, which is $O(nm^2)$.

6 BAYESIAN INFERENCE OPTIMIZATION

When the dataset is large, the number of variables in BN increases, rendering the BN inference cost prohibitive. To address this issue, a partition inference mechanism that conducts inference on nodes in BN, only interacting with those nodes within one hop in the Markov blanket. In addition, propose two pruning strategies: a tuple pruning strategy that identifies the cells in the dataset requiring repair and bypasses those that are largely clean, and a domain pruning strategy that eliminates values in domains that are clearly not candidate values for repair.

6.1 Partitioned Bayesian Inference

Bayesian inference methods can be either exact inference such as variable elimination and belief propagation, or approximate inference such as Gibbs sampling. Unlike existing techniques primarily concentrate on the global distribution of BN, according to the Markov property, we can treat the inferred node A_j as a state in BN, while disregarding the set of nodes that are not directly connected to A_j . This process essentially transforms a BN into several sub-networks. Each sub-network computes the partitioned distribution pertaining to the dividing node A_j , with all other nodes considered as observations. During Bayesian inference on each

node, only nodes and edges within its sub-network participate in the calculation.

BN partitioning. The variable elimination strategy for exact inference necessitates the inclusion of all precursor states of A_j in the topology, according to the inference rules of BN. However, the states of precursor nodes might have been altered in previous repairing steps. As such, the prior of A_j may not be the original one, and inferring from the first node could lead to higher costs when inferring every A_j . If the modified states are erroneous, these errors could propagate to the inference of A_j . To avoid this erroneous propagation and to expedite the inference, we partition BN into l sub-networks (we abuse the notation of a set of attributes to denote a sub-network): $\{\mathcal{A}_{\text{joint}}^{(1)}, \mathcal{A}_{\text{joint}}^{(2)}, \dots, \mathcal{A}_{\text{joint}}^{(l)}\}$, following the Markov blanket. $\mathcal{A}_{\text{joint}}^{(i)}$ ($1 \leq i \leq l$) denotes the i -th sub-network containing the inferred node A_j , alongside its one-hop parent nodes $\mathcal{A}^{(i)}$ parent and child nodes $\mathcal{A}^{(i)}$ child. Formally,

$$\mathcal{A}_{\text{joint}}^{(i)} = \mathcal{A}_{\text{parent}}^{(i)} \cup \{A_j\} \cup \mathcal{A}_{\text{child}}^{(i)}.$$

Multiple sub-networks might intersect at a node A_k , but $A_k \in \mathcal{A}_{\text{joint}}^{(i)}$ does not affect other sub-networks.

Bayesian inference. In the partitioned BN, nodes can be categorized into two types. One is isolated nodes, \mathcal{A}_{iso} , which do not connect to other nodes. The other is joint nodes, $\mathcal{A}_{\text{joint}} \in \text{BNsets}$, which connect to other nodes. For each isolated node, the CPT is modeled as uniform distribution, under the assumption that candidate values are uniformly distributed in the domain. For each joint node, all observations are known; thus, the probability distribution of candidate values is directly obtained from the CPT. The set of connected nodes, $\mathcal{A}_{\text{connected}}$, includes parent nodes $\mathcal{A}_{\text{parent}}$ and child nodes $\mathcal{A}_{\text{child}}$, with CPTs given by $\Pr[\mathcal{A}_j | \mathcal{A}_{\text{parent}}]$ and $\Pr[\mathcal{A}_{\text{child}} | A_j]$. In accordance with the properties of BN, we treat A_j as known, denoting every value in the domain as a candidate. As such, $\mathcal{A}_{\text{parent}}$ and $\mathcal{A}_{\text{child}}$ are independent. Formally,

$$\Pr[A_j | \mathcal{A}_{\text{connected}}] = \Pr[A_j | \mathcal{A}_{\text{parent}}] \Pr[\mathcal{A}_{\text{child}} | A_j].$$

6.2 Pruning Strategies

Data cleaning on D with values drawn from the domain of their corresponding attribute, and the inference cost is capped by the number of cells and the number of candidate values. Consequently, we introduce tuple and domain pruning strategies to respectively prune cells that do not need inspection and candidate values that are not correct answers.

Tuple pruning. Given a tuple T and an attribute A_i , we design a filtering mechanism to determine whether $T[A_i]$ should be inferred. The underlying intuition is that if $T[A_i]$ co-occurs with other values of T more frequently, indicating a stronger correlation, then $T[A_i]$ is more likely to be correct and $T[A_i]$ has a lower priority to be inferred. We define a function $\text{Filter}(T, A_i)$ and a threshold τ_{clean} to tell whether $T[A_i]$ needs inference. Formally,

$$\text{Filter}(T, A_i) = \frac{1}{m-1} \sum_{A_j \in \mathcal{A} \setminus \{A_i\}} \frac{\text{count}(T[A_i], T[A_j])}{\text{count}(T[A_j])}.$$

In the data cleaning process, we first compute $\text{Filter}(T, A_i)$. If the result is not less than τ_{clean} , we consider $T[A_i]$ as relatively reliable and omit the BN inference in the current iteration; otherwise, $T[A_i]$

is deemed obscure and requires repair. It is important to note that tuple pruning differs from standard error detection. Even though both can identify erroneous values, tuple pruning prioritizes cells with the fewest conflicts.

Domain pruning. We treat each sub-network as an independent semantic space and perform domain pruning from the semantic perspective, akin to the cloze test task in natural language processing.

In this semantic space, we take the dividing variable as the fill-in-the-blank variable, with all variables excluding the dividing variable acting as the context, and compute each answer under these semantics. For every value $v \in \text{dom}(A_i)$, we assign a weight using TF-IDF. Formally,

$$\begin{aligned} \text{score}(v) &= \text{TF}(v, \text{context}) \text{IDF}(v, D) \\ &= \text{context}(v) \cdot \log\left(\frac{|D|}{1 + \text{count}(v, D)}\right), \end{aligned}$$

where $\text{context}(v)$ denotes the number of sub-networks filled with v . These semantics treat different sub-networks as distinct contexts. The more frequently v appears in a specific semantic and the less frequently v appears in other semantics, the more likely v is to become the ground truth in that semantic. For every non-semantic $v \in \text{dom}(A_i)$, its score is significantly lower, possibly even zero. Each sub-network only computes a few candidates likely to be the ground truth during inference.

7 EXPERIMENTS

7.1 Experiment Setup

Datasets. We use six benchmark datasets of varying sizes and error types as summarized in Table 2.

- **Hospital:** This is a dataset derived from [13, 50], primarily comprising details of various hospitals such as HospitalName and Address. Approximately 5% of the data contains errors, with ground truth available for all cells. This dataset exhibits significant duplication across cells and possesses substantial causality between attributes.
- **Flights:** Used in [50] and collected from [38], this dataset contains departure and landing times of different airlines recorded on various websites. Ground truth information is available for part of the cells.
- **Soccer:** This dataset contains profiles of football players, including attributes like Name, Date of Birth, City, etc. The clean version of Soccer is collected from [48]. Errors account for approximately 5% of all data.
- **Beers:** This dataset used in [41, 42] contains two numerical attributes, ounces and abv. We acquire both the clean and dirty versions from the same source.
- **Inpatient:** This dataset of inpatient profiles was collected from the CMS [11].
- **Facilities:** This dataset contains information on medical enterprises collected from the CMS [11].

Error Injection. In line with the error injection approach utilized in the benchmarks (e.g., Hospital) by Raha+Baran [41, 42] and HoloClean [50], by default, we categorize errors into three types: typos (T), missing values (M), inconsistencies (I). For T, we modify the original value by randomly adding, deleting, or replacing a

character. For M, we randomly replace a value with NULL. For I, we interchange two values from the domains of two columns or a specific column. Their frequencies do not exhibit a significant difference, e.g., in the Inpatient dataset, we observe 1210, 1800, and 1480 instances of T, M, and I errors, respectively. In addition to these default error settings, we also evaluate swapping value errors (S), which are generated by swapping values within the same attribute, i.e., the same domain.

Methods. We compare the following methods.

- **BClean:** This is the basic version of BClean without any optimizations as detailed in Section 6.
- **BClean_{UC}:** This variant of BClean does not employ UCs.
- **BClean_{PI}:** This variant of BClean incorporates Partition Inference to our approach, which reduces unnecessary calculations by partitioning the network.
- **BClean_{PIP}:** This variant of BClean enables both Partition Inference and Pruning based on sub-network semantics.
- **PClean [35]:** This method uses PPL to detect and repair errors. PClean requires the user to modify the relations and distributions of the given datasets and outperforms other Bayesian data cleaning systems due to its superior ability to encode expert domain knowledge.
- **Raha+Baran [41, 42]:** This system combines two methods, Raha and Baran, to construct a machine learning-based data cleaning system. Raha is responsible for error detection using multiple error detection mechanisms, while Baran repairs detected errors using various error correction rules. We use the default settings, such as the number of labels, ML models, and so forth.
- **HoloClean [50]:** This semi-supervised inference system for data cleaning detects errors using signals and repairs them by compiling integrity rules, matching dependencies, and statistical signals into features in a Factor Graph. We run HoloClean with the signals either provided by the dataset owners or collected by ourselves.
- **Garf [47]:** This rule-based method employs sequence generative adversarial networks. It automatically learns rules from dirty data without any prior knowledge and uses them to repair the data. We use the default settings as provided in the paper.

Because DCs and labels used in HoloClean [50] and Raha+Baran [41, 42] papers have not been released, we configure them by ourselves.

Prior Knowledge. Table 2 also reports the statistical numbers of user inputs for all the competitors. For BClean, we inject 12 regular expressions (detailed in Table 3), and all attributes in these datasets adhere to UCs. We avoid using specific values in regular expressions to prevent the leakage of ground truth. We assign data-quality experts to each dataset to formulate the UCs based on the attribute formats. Furthermore, these experts are also trained to label 20 tuples for Raha, correct 20 tuples as labels for Baran, craft DCs for HoloClean, and develop PPL programs for PClean. Note that these baselines do not support expressions in UCs, and BClean does not use inputs from them, such as DCs or PPL programs.

Metrics. We use the following metrics to measure accuracy:

- **Precision:** The fraction of the correctly repaired errors over the total number of modified cells.
- **Recall:** The fraction of the correctly repaired errors over the total number of errors labeled with ground truth.

Table 2: Statistics of datasets. Error types include typos (T), missing values (M), inconsistency (I), and swapping value errors (S).

Dataset	Size (#rows, #columns, #cells)	Noise rate	Error types	#UCs (BClean)	#DCs (HoloClean)	#lines of PPL (PClean)	# labels of tuples (Raha + Baran)
Hospital	(1000, 15, 15k)	~5%	T, M, I	15	13	51	20+20
Flights	(2376, 6, 14k)	~30%	T, M	6	4	36	20+20
Soccer	(200000, 10, 2M)	~1%	T, M, I	10	4	37	20+20
Beers	(2410, 11, 27k)	~13%	T, M, I	11	6	26	20+20
Inpatient	(4017, 11, 44k)	~10%	T, M, I, S	11	3	54	20+20
Facilities	(7992, 11, 88k)	~5%	T, M, I, S	11	8	48	20+20

Table 3: List of user constraints (UCs).

DataSet	UCs (max. and min. length constraints for all textual attributes, not-null value constraints for all attributes.)
Hospital	$\wedge([1-9][0-9]\{4, 4\})[\text{ProvideNumber, ZipCode }]; \wedge([1-9][0-9]\{9, 9\})[\text{PhoneNumber }];$
Flights	$([1-9]:[0-5][0-9][s][ap].[m]. 1[0-2]:[0-5][0-9][s][ap].[m]. 0[1-9]:[0-5][0-9][s][ap].[m].)[\text{sched_dep_time, act_dep_time, sched_arr_time, act_arr_time}]$
Soccer	$([1][9][6-9][0-9])[\text{birthyear }]; ([2][0][0-9][0-9])[\text{season }]$
Beers	$\backslash d + \backslash . \backslash d + (\backslash d+)[\text{ounces, abv }]$
Inpatient	N/A
Facilities	N/A

Table 4: Precision (P), recall (R), and F1-score (F1) of data cleaning methods.

Method	Hospital			Flights			Soccer			Beers			Inpatient			Facilities		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
BClean _{UC}	1.000	0.935	0.966	0.807	0.729	0.766	0.927	0.982	0.954	0.880	0.065	0.121	0.934	0.883	0.908	0.810	0.805	0.807
BClean	0.998	0.956	0.976	0.852	0.816	0.834	0.928	0.979	0.952	0.916	0.887	0.901	0.909	0.845	0.876	-	-	-
BClean _{PI}	1.000	0.960	0.980	0.831	0.780	0.805	0.919	0.986	0.951	0.948	0.949	0.949	0.934	0.883	0.908	0.810	0.805	0.807
BClean _{PIP}	0.997	0.903	0.948	0.830	0.784	0.807	0.845	0.931	0.885	0.948	0.882	0.914	0.929	0.791	0.855	0.753	0.730	0.741
PClean	1.000	0.927	0.962	0.907	0.884	0.895	0.184 [†]	0.672 [†]	0.289 [†]	0.028 [†]	0.028 [†]	0.028 [†]	0.576 [†]	0.460 [†]	0.512 [†]	-	-	-
HoloClean	1.000	0.456	0.626	0.742	0.352	0.477	-	-	-	1.000	0.024	0.047	0.966	0.219	0.357	1.000	0.612	0.759
Raha+Baran	0.971	0.585	0.730	0.829	0.650	0.729	0.768	0.103	0.182	0.873	0.872	0.873	0.643	0.442	0.524	0.499	0.309	0.382
Garf	1.000	0.556	0.715	0.968	0.012	0.024	0.667	0.534	0.583	0.973	0.011	0.021	0.971	0.091	0.166	0.963	0.281	0.435

[†] We invite people familiar with PClean to author the data models in PPL.
 – Out-of-memory, out-of-runtime (24h), or no repairs.

- **F1-score:** The harmonic mean of precision and recall.

To measure time and user cost, we collect the following metrics:

- **User time:** The total time spent by users to inject prior knowledge into each baseline method. Initially, two experts are trained to understand the inputs and outputs of each baseline, including DCs of HoloClean, UCs of BClean, PPLs of PClean, and labeling of Raha+Baran. User Time represents the average time these experts spend to inject prior knowledge into a novel, unseen dataset. This metric assesses the usability of the system for trained users. We exclude training time as it varies across individuals.
- **Execution time:** The total runtime of the system modules to complete the data cleaning task, given a dataset and the users’ injected prior knowledge.

Parameters. For BClean, we set $\lambda = 1$, $\beta = 2$, and $\tau = 0.5$. These values reflect the principle that the lower the confidence in a tuple’s relationship due to UC violations, the less trustworthy it is considered. A tuple with a confidence score lower than τ is deemed untrustworthy, with a penalty score of $\beta = 2$ applied to the relationship. Conversely, a credible relationship score is set to $\lambda = 1$. For other baselines, we apply the default settings in PClean and HoloClean. We set labels of tuples to 20 + 20 following the suggestions from the Raha and Baran papers.

Environment. We conduct the experiments on a single machine equipped with 128GB RAM and 32 Intel(R) Core(TM) Xeon(R) 6326 CPU processors, each running at 2.90GHz. We run our experiments

Table 5: Precision / Recall / F1-score on sampled Soccer

BClean	HoloClean	PClean	Raha+Baran
0.345/0.931/0.504	0.919/0.551/0.689	0.15/0.665/0.244	0.523/0.133/0.212

5 times and reported the average performance for both effectiveness and efficiency.

7.2 Effectiveness and Efficiency

7.2.1 Data Cleaning Quality. We assess BClean’s precision, recall, and F1-score against other competitors. As illustrated in Table 4, BClean surpasses other techniques when the dataset has adequate relational information and a low error rate. BClean achieves superior precision, recall, and F1-scores in both Hospital and Soccer datasets, which can be attributed to the high-quality BN generated by BClean and its compensatory scoring model. For the Soccer dataset, BClean significantly outperforms other techniques. The poor performance of Raha+Baran can be attributed to the propagation of detection errors to the correction stage. BClean also outperforms Garf, as the latter primarily focuses on generating explainable rules for error correction.

Due to the out-of-memory issue encountered with HoloClean, we randomly sample 50,000 tuples from the Soccer dataset and evaluate the cleaning quality across all baselines. As shown in Table 5, although BClean’s performance is inferior to HoloClean due to a limited relational context, its recall significantly exceeds that of HoloClean.

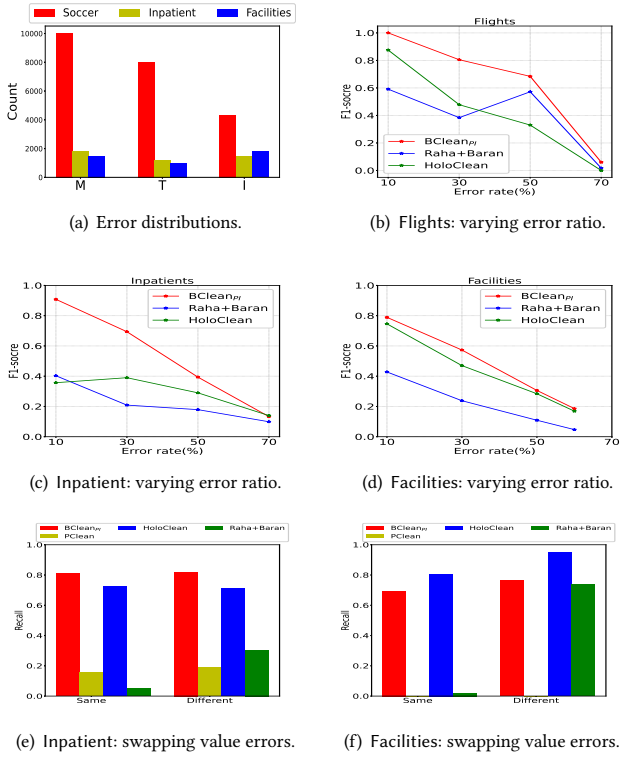


Figure 4: Error analysis results.

For the Flights dataset, despite its high error rate in comparison to the other datasets, BClean maintains commendable performance, achieving results comparable to the best-performing method.

Among the four versions of BClean, the efficiency-optimized ones, BClean_{PI} and BClean_{PIP}, display similar performance to the non-optimized BClean. These optimizations do not result in any significant impact on cleaning quality. Although BClean_{UC} does not incorporate UCs, it remains competitive owing to the robustness of its BN construction and inference in the presence of noise, e.g., achieving an F1-score of 0.966 on Hospital, compared to 0.626, 0.730, and 0.715 obtained by HoloClean, Raha+Baran, and Garf, respectively.

Regarding other methods, PClean achieves the best precision, recall, and F1-scores on the Flights dataset, however, this is largely reliant on precise domain knowledge modeling by experts. For the Soccer dataset, users find it challenging to articulate data distributions. Users can only estimate the distributions based on their observations and summaries, which results in less than satisfactory performance as evidenced by the results on Soccer.

7.2.2 Error Analysis. Next, we study how error types (T, M, and I, as specified in the error injection of the experiment setup) and ratio affect the performance. The distribution of errors across six datasets is displayed in Figure 4(a).

Varying error types. To assess the robustness of the competing methods, we evaluate BClean’s recall for each type of error within the Soccer, Inpatient, and Facilities datasets (Table 6). We do not report

Table 6: Recall for different types of errors (T, M, and I).

Method	Soccer			Inpatient			Facilities		
	T	M	I	T	M	I	T	M	I
BClean _{PI}	0.997	1.000	0.990	0.840	1.000	0.843	0.683	0.900	0.837
PClean	1.000	0.568	0.953	0.323	0.760	0.477	0.0	0.0	0.0
HoloClean	0.749	1.000	0.923	0.954	0.612	0.949	0.804	1.000	0.851
Raha+Baran	0.047	0.244	0.018	0.491	0.890	0.109	0.295	0.501	0.213

precision here, as it is challenging to determine which type of error a corrected value from the baselines originally belongs to. As shown in Table 6, BClean consistently outperforms the other methods for each type of error, with higher recalls of 0.8, 0.9, and 0.9 on average for T, M, and I, respectively, and up to 0.4, 0.5, and 0.5 higher than the best of the other methods. This showcases BClean’s robustness against various types of errors.

Varying error ratio. We assess the competing methods by varying the error ratio from 10% to 70%, as depicted in Figure 4(b) – 4(d), following the evaluation in [42]. The general trend reveals that all methods experience a decline in F1-score with increasing error ratios. However, BClean demonstrates greater robustness and consistently outperforms the others. Even at a high error ratio of 70%, BClean maintains a relative performance advantage, with an F1-score across the three datasets on average 0.1 higher than the best-performing alternative.

Swapping value errors. We evaluate the performance of BClean and other methods under conditions where 10% and 5% swapping value errors are randomly injected into the Inpatient and Facilities datasets. As illustrated in Figure 4(e), BClean outperforms the others with a recall on average 0.1 higher than the best among the other methods. This validates BClean’s capability to handle swapping value errors, primarily due to the use of Bayesian inference featuring a BN and a compensatory scoring model.

7.2.3 Runtime Analysis. We report runtime statistics for all the competitors across the datasets. As demonstrated in Table 7, the user time required for BClean is considerably less than that for all the competing methods, with the exception of Raha+Baran. Due to its easy-to-deploy design, BClean offers a more user-friendly learning curve and reduced complexity compared to the other methods.

Table 2 shows that BClean users need to compose a few format-based UCs, focusing primarily on single attributes. Garf requires no user input because it is based on a generative model, while Raha+Baran has the least user involvement among the competing methods.

Concerning execution time, BClean takes slightly longer than PClean but is faster than both HoloClean and Raha+Baran on the Hospital and Flights datasets. However, when processing the larger Soccer dataset with 2M cells, the basic BClean variant encounters efficiency issues, taking over 10 hours to complete. The optimized versions, BClean_{PI} and BClean_{PIP}, mitigate this problem by implementing partition inference and pruning. Their execution time is roughly on par with that of PClean. When considering the total time cost, which includes both user and execution time, BClean proves to be significantly more time-efficient than PClean and HoloClean. Although Raha+Baran boasts a shorter total time than BClean, its lower quality reveals its limitations. It’s also worth noting that Raha+Baran requires significant time to infer correct

Table 7: Runtime of data cleaning methods, including user time (user) and execution time (exec).

Method	Time	Hospital	Flights	Soccer	Beers	Inpatient	Facilities
PClean	user	≥ 72h	≥ 72h	≥ 72h	≥ 72h	≥ 72h	≥ 72h
	exec	16s	7s	30m 44s	2m55s	3m17s	1m32s
HoloClean	user	15h	12h	14h	15h	15h	15h
	exec	1m 40s	36s	-	1m37s	4m14s	6m2s
Raha+Baran	user	30m	30m	30m	30m	30m	30m
	exec	1m 46s	41s	8m 59	3m2s	10m36s	10m55s
Garf	user	0	0	0	0	0	0
	exec	5m24s	1m57s	18h30m	2m8s	26m48s	30m10s
BClean	user	5h	2h	3h	2h	3h	3h
	exec	25s	17s	10h 48m	1m40s	7h41m	≥ 72h
BClean _{PI}	user	5h	2h	3h	2h	3h	3h
	exec	22	12s	30m 42s	31s	7m57s	17m16s
BClean _{PIP}	user	5h	2h	3h	2h	3h	3h
	exec	22s	12s	27m 46s	30s	7m2s	14m35s

values, as labeling tuples requires looking up values from other attributes and tuples, and necessitates user checks for contextual information. Garf, on the other hand, has a high execution time due to the generation of rules for interpretability.

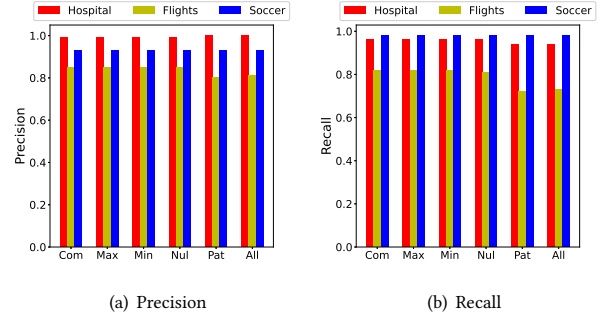
7.3 Analysis of User Interactions

BClean involves user-interactions in two operations: defining UCs and manipulating the BN. We evaluate their impacts on the performance.

7.3.1 Impact of User Constraints. We assess the significance of each UC by measuring the system’s performance after the removal of some or all types of constraints. UCs consist of max length (Max), min length (Min), allowing null values (Nul), and patterns (Pat). We evaluate the system’s performance with each constraint type removed, as well as after removing all constraints (All), and compare them to the complete version (Com).

As demonstrated in Figure 5, we find that the Pat constraint is the most influential. Changes in precision and recall are minor or even non-existent when Max, Min, and Nul are incomplete. When the user cannot provide the correct pattern (regular expression), or when no pattern exists to describe the dataset, a decline in precision and recall becomes apparent, as shown in Figure 5 (a) and (b) respectively. For optimal performance, BClean still needs expert input. Constraints provided by experts can directly exclude candidate values that do not meet the requirements before the inference process begins. This technique is extremely useful for data with many instances where the correct value is less than the error value and the error value does not meet the requirements, such as the Flights dataset.

In more detail, without pattern constraints, we receive two candidates with probabilities $Pr[g_0 = "7 : 10 a.m."] = 0.37$ and $Pr[g_1 = "7 : 21am"] = 0.57$, and ground g_1 to be the final result in Flights. As shown in Figure 5, the system achieves sub-optimal inference and less effective error detection with lower precision (0.77) and recall (0.72) in the absence of a pattern. When a pattern is present, $Pr[g_1]$ is set to 0 prior to inference to prevent it from appearing in the final result. Thus, we only infer the ground truth repair within a restricted domain, yielding higher precision and recall. For Hospital, the pattern filters out the candidate value $g_3[\text{ZipCode}] = "1xx18"$, as the ZipCode in Hospital must be a five-digit numeric string. This situation also occurs in Soccer.

**Figure 5: Effect of incomplete UCs on precision and recall.****Table 8: Varying λ on Hospital, $\beta = 2$, $\tau = 0.5$.**

λ	0	1	2	5	10	15
F1	0.98096	0.98096	0.98096	0.98096	0.98096	0.98096

Table 9: Varying β on Hospital, $\lambda = 1$, $\tau = 0.5$.

β	0	1	2	10	50
F1	0.97996	0.98096	0.98096	0.98096	0.98096

Table 10: Varying τ on Hospital, $\lambda = 1$, $\beta = 2$.

τ	0.1	0.3	0.5	0.7	0.9
F1	0.98096	0.98096	0.98096	0.97996	0.97996

In conclusion, while the absence of user constraints results in a slight decline in cleaning performance, the overall reduction remains within an acceptable range. This demonstrates that BClean is user-friendly and suitable for users with limited expertise.

7.3.2 Impact of Network Manipulation. We conduct an experiment involving the modification of the BN through user interaction. The results indicate no difference between the states before and after the modification for the Hospital and Soccer datasets. After users’ adjustment on BN, the Hospital dataset adds an edge state \rightarrow state_avg, resulting in one additional cell being cleaned, whereas the Soccer dataset displays no changes. However, for the Flight dataset, the BN structure obtained by the greedy search or FD-based method is incorrect, yielding cleaning results with a precision of only 0.217 and recall of 0.374. Following user adjustments, the results are improved to 0.852 precision and 0.816 recall, as shown in Table 4. In addition, such adjustments are with less than 5 minutes, suggesting that the user’s effort is marginal.

7.4 Parameter Tuning

In BClean, there are three parameters: λ , β , and τ . To assess their impact, we fix two parameters at a time and vary the other. Tables 8 – 10 report the F1-score on the Hospital dataset. All three parameters have a minimal effect on the F1-score. This observation highlights that the performance of BClean is stable and largely unaffected by parameter changes, aligning with the easy-to-deploy design of BClean, which requires users to spend little effort on parameter tuning.

8 RELATED WORK

Data cleaning. Data cleaning is an extensively-studied problem. We can broadly categorize the methods into five types.

- **Generative methods:** Generative data cleaning methods integrate error detection and correction components, often leveraging probabilistic inference to iteratively clean datasets towards maximum likelihood [16, 18, 29, 52]. Other notable methods in this category use PPLs [35, 44] or user queries [7, 41]. Generative models are data-driven and do not require any user labels, but many of them necessitate some prior knowledge to achieve good performance.
- **Rule-based methods:** These employ data quality rules that check data inconsistency, such as FDs [2], conditional FDs [9, 20], DCs [13], user-defined functions [15, 27], and regular expression entities [24]. Rule discovery algorithms [21–23] have been proposed to mine these rules from relatively clean training data. Correspondingly, error detection [26] and data cleaning algorithms [19, 62] have been designed to exploit these rules to efficiently locate and correct errors in big data.
- **ML pipeline:** ML models have been widely used for data clean. For instance, SCAREd [61] adopts ML and likelihood methods. In CleanML[36], five error types were evaluated, and the impact of data cleaning on the ML pipeline was discussed. Picket [39] employs a self-supervised strategy to rectify data errors in ML pipelines. [17] proposed a novel framework based on Shapley values to interpret the results of any data cleaning module. ActiveClean [33] is a semi-supervised model with convex loss functions for data cleaning.
- **ML imputation:** Additionally, many ML methods focus on imputing missing values. A common approach is the denoising autoencoder [57]. More recent approaches use generative adversarial networks (GANs) [63], the attention mechanism [54, 59], or both [30]. However, many of them require considerable effort in tuning hyperparameters and training the model, especially for GAN-based methods.
- **Hybrid models:** These methods blend logical rules and ML models so that rules are generated and used as features, and ML models are applied for prediction using these features as inputs, e.g., [28, 41, 42, 50]. For instance, Raha+Baran [41, 42] applies several feature engineering strategies (FDs, outliers, and patterns) to generate features and then employs ML models for prediction. HoloDetect [28] uses data augmentation and neural networks to detect errors while minimizing human involvement.

In contrast to the aforementioned methods, BClean employs a novel direction, executing approximate probabilistic inference with an easy-to-deploy configuration that includes UCs (e.g., regular expressions and statistical indicators), and a BN construction algorithm. User interactions are leveraged to fine-tune and retrieve the cleaning solution.

Bayesian network construction. Existing BN structure learning methods include automatic greedy search [46, 55, 56] and user-constructed networks [35, 44]. Automatic methods tend to be weakly robust to dirty data due to erroneous propagation, while user-constructed methods often incur high user costs when dealing with large datasets. In contrast, BClean facilitates fine-tuning on structure learning via user interaction, which allows for minimal

domain knowledge requirements while ensuring high data cleaning accuracy using the BN.

Bayesian inference. Bayesian inference can be categorized into exact inference and approximate inference. Exact inference yields posterior probability distributions, but methods like variable elimination and belief propagation [46] can be computationally intensive and susceptible to erroneous propagation. Approximate inference, based on sampling techniques such as Gibbs sampling [64], typically trades runtime improvement for accuracy. Exact inference is challenging to execute outside clean data, and approximate inference can propagate errors when sampling dirty data. BClean generates approximate inference by partitioning the BN, thereby balancing accuracy with runtime efficiency.

9 CONCLUSION

We introduced BClean, an easy-to-deploy Bayesian data cleaning system that requires minor user effort in specifying domain knowledge. BClean automatically constructs a Bayesian network from an observed dataset and optionally adjusts the network through user interaction. User constraints, which cover a wide range of constraint forms such as dependency rules (e.g., functional dependencies and denial constraints) and arithmetic expressions, can be specified to reflect domain knowledge. BClean detects erroneous data and infers correct values using Bayesian inference, which leverages the Bayesian network and a compensatory scoring model. To reduce the processing cost of data cleaning, we proposed a set of optimization techniques. Our experiments demonstrated the effectiveness and efficiency of BClean and its superiority over alternative solutions.

REFERENCES

- [1] Z. Abedjan, X. Chu, D. Deng, R. C. Fernandez, I. F. Ilyas, M. Ouzzani, P. Papotti, M. Stonebraker, and N. Tang. Detecting data errors: Where are we and what needs to be done? *PVLDB*, 9(12):993–1004, 2016.
- [2] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [3] M. Ale Ebrahim Dehkordi, J. Lechner, A. Ghorbani, I. Nikolic, E. Chappin, and P. Herder. Using machine learning for agent specifications in agent-based models and simulations: A critical review and guidelines. *Journal of Artificial Societies and Social Simulation*, 26(1), 2023.
- [4] A. Ankan and A. Panda. pgmpy: Probabilistic graphical models using python. In *SCIPY*, 2015.
- [5] A. Bartoli, A. D. Lorenzo, E. Medvet, and F. Tarlao. Regex Generator++. <http://regex.inginf.units.it/>.
- [6] A. Bartoli, A. D. Lorenzo, E. Medvet, and F. Tarlao. Inference of regular expressions for text extraction from examples. *IEEE Trans. Knowl. Data Eng.*, 28(5):1217–1230, 2016.
- [7] L. E. Bertossi, S. Kolahi, and L. V. S. Lakshmanan. Data cleaning and query answering with matching dependencies and matching functions. In *ICDT*, pages 268–279, 2011.
- [8] E. Bingham, J. P. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. Szerlip, P. Horsfall, and N. D. Goodman. Pyro: Deep universal probabilistic programming. *Journal of Machine Learning Research*, 20(1):973–978, 2019.
- [9] P. Bohannon, W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for data cleaning. In *ICDE*, pages 746–755, 2007.
- [10] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1):1–32, 2017.
- [11] Centers for Medicare & Medicaid Services. Provider data catalog. <https://data.cms.gov/provider-data/>.
- [12] M. Chickering, D. Geiger, and D. Heckerman. Learning bayesian networks: Search methods and experimental results. In *AISTATS*, 1995.
- [13] X. Chu, I. F. Ilyas, and P. Papotti. Holistic data cleaning: Putting violations into context. In *ICDE*, pages 458–469, 2013.
- [14] X. Chu, J. Morcos, I. F. Ilyas, M. Ouzzani, P. Papotti, N. Tang, and Y. Ye. KATARA: A data cleaning system powered by knowledge bases and crowdsourcing. In

- SIGMOD*, pages 1247–1261, 2015.
- [15] M. Dallachiesa, A. Ebaïd, A. Eldawy, A. K. Elmagarmid, I. F. Ilyas, M. Ouzzani, and N. Tang. NADEEF: a commodity data cleaning system. In *SIGMOD*, pages 541–552, 2013.
- [16] S. De, Y. Hu, V. V. Meduri, Y. Chen, and S. Kambhampati. Bayeswipe: A scalable probabilistic framework for improving data quality. *Journal of Data and Information Quality*, 8(1):5:1–5:30, 2016.
- [17] D. Deutch, N. Frost, A. Gilad, and O. Sheffer. Explanations for data repair through shapley values. In *CIKM*, pages 362–371, 2021.
- [18] P. Doshi, L. G. Greenwald, and J. R. Clarke. Using bayesian networks for cleansing trauma data. In *FLAIRS*, pages 72–76, 2003.
- [19] W. Fan, F. Geerts, and X. Jia. Semandaq: a data quality system based on conditional functional dependencies. *PVLDB*, 1(2):1460–1463, 2008.
- [20] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *ACM Transactions on Database Systems*, 33(2):6:1–6:48, 2008.
- [21] W. Fan, F. Geerts, L. V. S. Lakshmanan, and M. Xiong. Discovering conditional functional dependencies. In *ICDE*, pages 1231–1234, 2009.
- [22] W. Fan, F. Geerts, J. Li, and M. Xiong. Discovering conditional functional dependencies. *IEEE Transactions on Knowledge and Data Engineering*, 23(5):683–698, 2011.
- [23] W. Fan, Z. Han, Y. Wang, and M. Xie. Parallel rule discovery from large datasets by sampling. In Z. Ives, A. Bonifati, and A. E. Abbadi, editors, *SIGMOD*, pages 384–398, 2022.
- [24] W. Fan, P. Lu, and C. Tian. Unifying logic rules and machine learning for entity enhancing. *Science China Information Sciences*, 63(7), 2020.
- [25] W. Fan, S. Ma, N. Tang, and W. Yu. Interaction between record matching and data repairing. *Journal of Data and Information Quality*, 4(4):1–38, 2014.
- [26] W. Fan, C. Tian, Y. Wang, and Q. Yin. Parallel discrepancy detection and incremental detection. *PVLDB*, 14(8):1351–1364, 2021.
- [27] F. Geerts, G. Mecca, P. Papotti, and D. Santoro. The LLUNATIC data-cleaning framework. *PVLDB*, 6(9):625–636, 2013.
- [28] A. Heidari, J. McGrath, I. F. Ilyas, and T. Rekatsinas. Holodetect: Few-shot learning for error detection. In *SIGMOD*, pages 829–846, 2019.
- [29] Y. Hu, S. De, Y. Chen, and S. Kambhampati. Bayesian data cleaning for web data. *CoRR*, abs/1204.3677, 2012.
- [30] J. Kawagoshi, Y. Dong, T. Nozawa, and C. Xiao. CAGAIN: Column attention generative adversarial imputation networks. In *DEXA*, volume 14147 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 2023.
- [31] S. Kolahi and L. V. S. Lakshmanan. On approximating optimum repairs for functional dependency violations. In *ICDT*, pages 53–62, 2009.
- [32] N. Koudas, A. Saha, D. Srivastava, and S. Venkatasubramanian. Metric functional dependencies. In Y. E. Ioannidis, D. L. Lee, and R. T. Ng, editors, *ICDE*, pages 1275–1278. IEEE Computer Society, 2009.
- [33] S. Krishnan, J. Wang, E. Wu, M. J. Franklin, and K. Goldberg. Activeclean: Interactive data cleaning for statistical modeling. *PVLDB*, 9(12):948–959, 2016.
- [34] J. Kubica and A. W. Moore. Probabilistic noise identification and data cleaning. In *ICDM*, pages 131–138, 2003.
- [35] A. K. Lew, M. Agrawal, D. A. Sontag, and V. Mansinghka. Pclean: Bayesian data cleaning at scale with domain-specific probabilistic programming. In *AISTATS*, pages 1927–1935, 2021.
- [36] P. Li, X. Rao, J. Blase, Y. Zhang, X. Chu, and C. Zhang. Cleanml: A study for evaluating the impact of data cleaning on ML classification tasks. In *ICDE*, pages 13–24, 2021.
- [37] W. Li, L. Li, Z. Li, and M. Cui. Statistical relational learning based automatic data cleaning. *Frontiers of Computer Science*, 13(1):215–217, 2019.
- [38] X. Li, X. L. Dong, K. Lyons, W. Meng, and D. Srivastava. Truth finding on the deep web: Is the problem solved? *arXiv preprint arXiv:1503.00303*, 2015.
- [39] Z. Liu, Z. Zhou, and T. Rekatsinas. Picket: guarding against corrupted data in tabular data during learning and inference. *Vldb Journal*, 31(5):927–955, 2022.
- [40] P.-L. Loh and P. Bühlmann. High-dimensional learning of linear causal networks via inverse covariance estimation. *Journal of Machine Learning Research*, 15(1):3065–3105, 2014.
- [41] M. Mahdavi and Z. Abedjan. Baran: Effective error correction via a unified context representation and transfer learning. *PVLDB*, 13(11):1948–1961, 2020.
- [42] M. Mahdavi, Z. Abedjan, R. Castro Fernandez, S. Madden, M. Ouzzani, M. Stonebraker, and N. Tang. Raha: A configuration-free error detection system. In *SIGMOD*, pages 865–882, 2019.
- [43] C. Mayfield, J. Neville, and S. Prabhakar. A statistical method for integrated data cleaning and imputation. 2009.
- [44] B. Milch, B. Marthi, S. Russell, D. Sontag, D. L. Ong, and A. Kolobov. 1 bloc: Probabilistic models with unknown objects. *Statistical Relational Learning*, page 373, 2007.
- [45] P. H. Oliveira, D. S. Kaster, C. T. Jr., and I. F. Ilyas. Batchwise probabilistic incremental data cleaning. *CoRR*, abs/2011.04730, 2020.
- [46] S. Parsons. *Probabilistic Graphical Models: Principles and Techniques* by daphne koller and nir friedman, MIT press, 1231 pp., \$95.00, ISBN 0-262-01319-3. *Knowledge Engineering Review*, 26(2):237–238, 2011.
- [47] J. Peng, D. Shen, N. Tang, T. Liu, Y. Kou, T. Nie, H. Cui, and G. Yu. Self-supervised and interpretable data cleaning with sequence generative adversarial networks. *PVLDB*, 16(3):433–446, 2022.
- [48] J. Rammelaere and F. Geerts. Explaining repaired data with cfd. *PVLDB*, 11(11):1387–1399, 2018.
- [49] G. Raskutti and C. Uhler. Learning directed acyclic graph models based on sparsest permutations. *Stat*, 7(1):e183, 2018.
- [50] T. Rekatsinas, X. Chu, I. F. Ilyas, and C. Ré. Holoclean: Holistic data repairs with probabilistic inference. *PVLDB*, 10(11):1190–1201, 2017.
- [51] E. K. Rezig, M. Ouzzani, A. K. Elmagarmid, W. G. Aref, and M. Stonebraker. Towards an end-to-end human-centric data cleaning framework. In *HILDA*, pages 1–7, 2019.
- [52] C. D. Sa, I. F. Ilyas, B. Kimelfeld, C. Ré, and T. Rekatsinas. A formal framework for probabilistic unclean databases. In *ICDT*, pages 6:1–6:18, 2019.
- [53] P. L. Spirtes and C. Glymour. An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9:62 – 72, 1991.
- [54] S. Tihon, M. U. Javaid, D. Fourure, N. Posocco, and T. Peel. DAEMA: denoising autoencoder with mask attention. In *ICANN*, pages 229–240, 2021.
- [55] I. Tsamardinou, L. E. Brown, and C. F. Aliferis. The max-min hill-climbing bayesian network structure learning algorithm. *Machine Learning*, 65(1):31–78, 2006.
- [56] K. Tzoumas, A. Deshpande, and C. S. Jensen. Lightweight graphical models for selectivity estimation without independence assumptions. *PVLDB*, 4(11):852–863, 2011.
- [57] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P. Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11:3371–3408, 2010.
- [58] J. Wang and N. Tang. Towards dependable data repairing with fixing rules. In *SIGMOD*, pages 457–468, 2014.
- [59] R. Wu, A. Zhang, I. F. Ilyas, and T. Rekatsinas. Attention-based learning for missing data imputation in holoclean. In *MLSys*, 2020.
- [60] S. Wu, S. Sanghavi, and A. G. Dimakis. Sparse logistic regression learns all discrete pairwise graphical models. In *NeurIPS*, pages 8069–8079, 2019.
- [61] M. Yakout, L. Berti-Équille, and A. K. Elmagarmid. Don't be scared: use scalable automatic repairing with maximal likelihood and bounded changes. In *SIGMOD*, pages 553–564, 2013.
- [62] M. Yakout, A. K. Elmagarmid, J. Neville, M. Ouzzani, and I. F. Ilyas. Guided data repair. *PVLDB*, 4(5):279–289, 2011.
- [63] J. Yoon, J. Jordon, and M. van der Schaar. GAIN: missing data imputation using generative adversarial nets. In *ICML*, pages 5675–5684, 2018.
- [64] C. Zhang and C. Ré. Dimmwwitted: A study of main-memory statistical analytics. *PVLDB*, 7(12):1283–1294, 2014.
- [65] Y. Zhang, Z. Guo, and T. Rekatsinas. A statistical perspective on discovering functional dependencies in noisy data. In *SIGMOD*, pages 861–876, 2020.
- [66] B. Zhao, B. I. P. Rubinstein, J. Gemmell, and J. Han. A bayesian approach to discovering truth from conflicting sources for data integration. *PVLDB*, 5(6):550–561, 2012.
- [67] Y. Zheng, J. Wang, G. Li, R. Cheng, and J. Feng. QASCA: A quality-aware task assignment system for crowdsourcing applications. In *SIGMOD*, pages 1031–1046, 2015.