

# A Survey on Sampling and Profiling over Big Data (Technical Report)

Zhicheng Liu\*, Aoqian Zhang<sup>§</sup>

\*Tsinghua University, lzc18@mails.tsinghua.edu.cn

<sup>§</sup>University of Waterloo, aoqian.zhang@uwaterloo.ca

arXiv:2005.05079v1 [cs.DB] 8 May 2020

**Abstract**—Due to the development of internet technology and computer science, data is exploding at an exponential rate. Big data brings us new opportunities and challenges. On the one hand, we can analyze and mine big data to discover hidden information and get more potential value. On the other hand, the 5V characteristic of big data, especially Volume which means large amount of data, brings challenges to storage and processing. For some traditional data mining algorithms, machine learning algorithms and data profiling tasks, it is very difficult to handle such a large amount of data. The large amount of data is highly demanding hardware resources and time consuming. Sampling methods can effectively reduce the amount of data and help speed up data processing. Hence, sampling technology has been widely studied and used in big data context, e.g., methods for determining sample size, combining sampling with big data processing frameworks. Data profiling is the activity that finds metadata of data set and has many use cases, e.g., performing data profiling tasks on relational data, graph data, and time series data for anomaly detection and data repair. However, data profiling is computationally expensive, especially for large data sets. Therefore, this paper focuses on researching sampling and profiling in big data context and investigates the application of sampling in different categories of data profiling tasks. From the experimental results of these studies, the results got from the sampled data are close to or even exceed the results of the full amount of data. Therefore, sampling technology plays an important role in the era of big data, and we also have reason to believe that sampling technology will become an indispensable step in big data processing in the future.

**Index Terms**—Big Data, Large Amount, Sampling, Data Profiling

## I. INTRODUCTION

**W**E are in the era of big data. With the development of computer science and internet technology, data is exploding at an exponential rate. According to statistics, Google processes more than hundreds of PB data per day, Facebook generates more than 10 PB of log data per month, Baidu processes nearly 100 PB of data per day, and Taobao generates dozens of terabytes online transaction data every day [1]. In May 2011, the McKinsey Global Institution (MGI) released the report<sup>1</sup> which said that big data has great potential in the European Public Sector, US Health Care, Manufacturing, US Retail Industry and Location-based Services. MGI estimates in the report that the mining and analysis of big data will generate 300 billion in potential value per year in the US medical sector and more than 149 billion in the European public sector [2]. It

can be seen that there is great value behind big data. Therefore, mining the hidden value under big data makes a lot of sense.

Big data is something so huge and complex that it is difficult or impossible for traditional systems and tools to process and work on it [3]. In the latest development, IBM uses "5Vs" model to depict big data. In the "5Vs" model, Volume means the amount of data and it is the most direct difficulty faced by traditional systems; Velocity means that data is generated quickly; Variety means that data sources and data types are diverse including structural, semi-structured, and unstructured data; Value is the most important feature of big data, although the value density of data is low; Veracity refers to that data quality of big data where there is dirty data. Because big data is so large that data analysis and data mining based on big data require high computing power and storage capacity. In addition, some classical mining algorithms that require several passes over the whole dataset may take hours or even days to get result [4].

### A. Data Sampling

At present, there are two major strategies for data mining and data analysis: sampling and using distributed systems [5]. The existing big data processing framework includes batch processing framework like Apache Hadoop, streaming data processing framework like Apache Storm, hybrid processing framework like Apache Spark and Apache Flink. Sampling is a scientific method of selecting representative sample data from target data. Designing a big data sampling mechanism is to reduce the amount of data to a manageable size for processing [6]. Even if computer clusters are available, we can use sampling such as block-level sampling to speed up big data analysis [7].

Different from distributed systems, sampling is a kind of data reduction method like filtering. Distributed systems increase computing power by adding hardware resources. However, a huge computing cost is not always affordable in practice. It is highly demanded to perform the computing under limited resources. In this sense, sampling is very useful. Since the full amount of data is not used, the approximate result is obtained from the sample data. Such approximate result is quite useful in the context of big data. The computational challenge of big data means that sampling is essential and the sampling methods chosen by researchers is also important [8]. Besides, the biases caused by sampling are also something need to be considered.

<sup>1</sup>The Next Frontier of Big Data: Innovation, Competition, and Productivity

Sampling or re-sampling is to use less data to get the overall characteristics of the whole dataset. Albattah [9] studies the role of sampling in big data analysis. He believes that even if we can handle the full amount of data, we don't have to do this. They focus on how sampling will play its role in specific fields of Artificial Intelligence and verify it by doing experiments. The experimental results show that sampling not only reduces the data processing time, but also get better results in some cases. Even though some examples of sampling are not as effective as the original dataset, they are obviously negligible compared to the greatly reduced processing time. As stated in [9], we believe that sampling can improve big data analysis and will become a preprocessing step in big data processing in the future.

When it comes to sampling, how to determine the sample size is a very important factor, and different scholars have proposed many methods to determine the sample size [10]–[13]. And we also have to consider sampling bias when using sampling techniques. In addition, some scholars have also studied the application of sampling techniques in the context of big data, e.g., combining sampling with distributed storage, big data computing frameworks. And these will be introduced in detail in Section III.

### B. Data Profiling

Data mining is an emerging research area, whose goal is to extract significant patterns or interesting rules from large data sets [14]. Data profiling gathers metadata of data that can be used to find data to be mined and import data into various tools for analysis, which is an important preparatory task [15]. There is currently no formal, universal or widely accepted definition of distinction between data profiling and data mining. Abedjan et al. [16] think data profiling is used to generate metadata for data sets that are used to help understand data sets and manage data sets. However, data mining is used to mine the hidden knowledge behind the data, which is not so obvious. Of course, data profiling and data mining also have some overlapping tasks, such as association rule mining and clustering. In summary, the goal of data profiling is to generate summary information about the data to help understand the data, and the goal of data mining is to mine the new insights of the data.

There are many use cases of data profiling, such as data profiling for missing data imputation [17], [18] or erroneous data repairing in relational database [19]. However, data profiling itself has to face computational challenges, especially when it comes to large data sets. Hence, how to alleviate the computational challenges of data profiling is very significant in era of big data. As mentioned above, sampling for big data profiling is very valuable and meaningful. We will give a brief introduction for data profiling in Section II.

### C. Sampling for Data Profiling

In this paper, we focus on the sampling techniques used for big data profiling. Certainly, we will first introduce data profiling and sampling technology separately. Among them, data profiling has been associated with outstanding survey

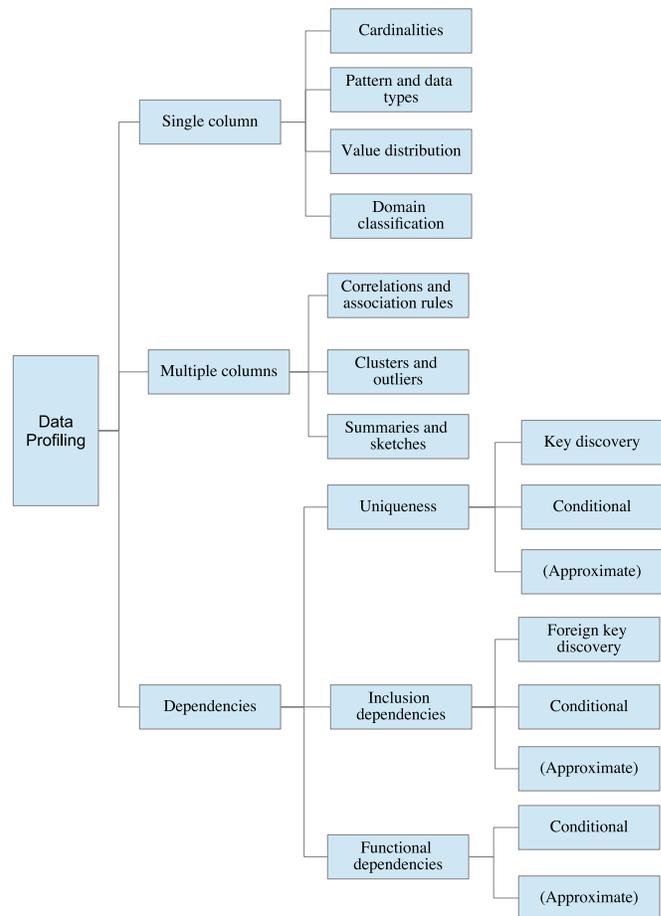


Fig. 1: A classification of typical data profiling tasks [16].

papers such as [16]. Finally, our core content is to introduce the application of sampling in data profiling tasks when facing large data sets.

In [16], the research on data profiling around the relational database is fully investigated and introduced. The classification of data profiling (see Figure 1) is given in [16]. We will investigate the sampling techniques for important data profiling tasks in single column, multiple columns and dependency according to the classification of data profiling in [16]. Some traditional sampling methods are introduced in [10], and methods of determining the sample size are mainly introduced, but less attention is paid to sampling in big data context. Therefore, when discussing the sampling technology below, we will supplement some applications and information of sampling in the big data scenario, e.g., block-based sampling.

Specifically, in order to ensure the comprehensiveness of the survey, we follow the systematic search method provided in [16], a comprehensive summary of data profiling techniques. As also illustrated in Figure 1 of our manuscript, Abedjan et al. [16] categorize the data profiling approaches into three aspects, from the elementary columns to the complex ones, i.e., (1) data profiling for single columns, (2) data profiling for multiple columns, and (3) data profiling for dependencies. While the sampling techniques for data profiling are not emphasized in [16], in our paper, we extensively select the studies on sam-

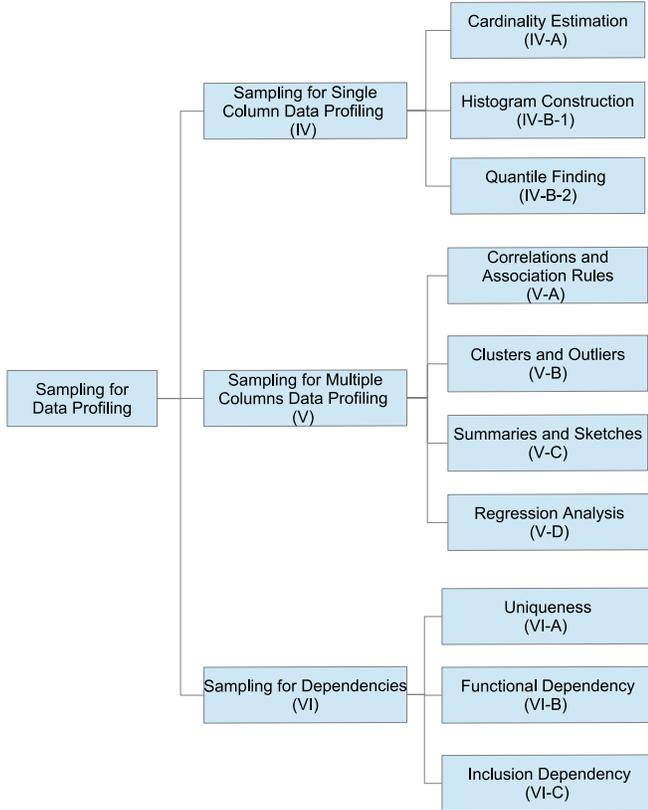


Fig. 2: A systematic search method for selecting studies, following the categorization in Figure 1 by [16].

pling for data profiling in the aforesaid categories, respectively. Figure 2 presents the systematic search method for selecting studies, following the categorization in [13]. Following this method, we summarize the typical methods selected in each category in Table III.

The remaining of this paper is organized as follows. In Section II, we introduce the relevant knowledge of data profiling. In Section III, we introduce sampling techniques and some important factors in sampling techniques. Next we introduce the application of sampling for single-column data profiling tasks in Section IV, multi-column data profiling tasks in Section V and dependencies in Section VI based on the classification of data profiling tasks in [16]. Finally, in Section VII, we summarize the content of the article and propose some future works. The organizational structure of this article is shown in Figure 2.

## II. DATA PROFILING

Before using or processing data, it is very important to have a general understanding of the data. Data profiling is the activity that finds metadata of data set [16], [20], [21], therefore it can provide basic information about data to help people understand the data. Data profiling is an important area of research for many IT experts and scholars. Data profiling has many classic use cases, such as data integration, data quality, data cleansing, big data analysis, database management, query optimization [16], [20]. Abedjan et al. [16] mainly investigates data profiling for relational data. However, in addition to

relational databases, many non-relational databases need data profiling [20], such as time series data [22]–[24], graph data [25]–[27], or heterogeneous data in dataspace [28]–[30].

Data profiling tasks are classified in [16] and [20]. Abedjan et al. [16] classify the data profiling tasks of single data source, and divides the tasks of data profiling into single column data profiling, multiple columns data profiling and dependency (see Figure 1). In fact, dependencies belong to multiple columns data profiling tasks. Abedjan et al. [16] put dependencies separately into a large category and discuss it in detail. Naumann [20] classifies data profiling from single data source to multiple data sources.

There are three challenges for data profiling: managing the input, performing the computation and managing the output [16], [20], [31]. In this article we focus on the second challenge, performing the computation, i.e., the computational complexity of data profiling. The computational complexity of data profiling depends on the number of rows and columns of data. When the data set is very large, the calculation of data profiling can be very expensive. This is why we care about sampling for big data profiling, in order to reduce the computational pressure and speed up the process of data profiling.

## III. SAMPLING TECHNIQUES

In this section, we introduce common sampling techniques in III-A, application of sampling technology in big data context in III-B, methods of determining sample size in III-C and resolutions of reducing sampling bias in III-D.

### A. Common Sampling Techniques

Sampling refers to estimating the characteristics of the entire population through the representative subsets within the population [10]. From a big perspective, sampling involves probability and non-probability sampling. Probability sampling means that every unit in a finite population has a certain probability to be selected, and it does not necessarily require equality. Non-probability sampling is generally based on subjective ideas and inferences, e.g., common web questionnaires [32], [33]. The sampling methods mentioned below are all probability sampling methods. Sampling is often used in data profiling [16], data analysis [34], data mining [6], data visualization [35], machine learning [36] etc. The advantage of sampling is that algorithms or models can be conducted using subset instead of the whole data set. There are some commonly used sampling techniques including simple random sampling [37], stratified sampling [38], systematic sampling [39], cluster sampling [40], oversampling and undersampling [41], [42], reservoir sampling [43], etc. Table I gives an overview of these common sampling methods.

### B. Sampling in Big Data Context

In the era of big data, the application of sampling is particularly important due to the large amount of big data. And sampling can be performed with the help of big data computing framework. For example, He et al. [44] use MapReduce

TABLE I: Common sampling methods

Sampling method	Description
Simple random sampling	Extracting a certain number of samples and each tuple is selected with equal probability.
Stratified random sampling	Tuples are divided into homogeneous groups and sample from each group.
Systematic sampling	Sampling at regular intervals until the sample size is satisfied.
Cluster sampling	Tuples are divided into non-overlapping groups and randomly select some groups as samples.
Oversampling and Undersampling	Oversampling randomly repeat the minority class samples, while Undersampling randomly discard the majority class samples to balance the data.
Reservoir sampling	Adding tuples into the reservoir of a fixed size with unknown size of the entire data set.

to sample from the data which contains uncertainty. He et al. [45] propose a block-based sampling (I-sampling) method for large-scale dataset when the whole dataset is already assigned on a distributed system. The processing flow of I-sampling is shown in Figure 3.

The traditional sampling methods like simple random sampling, stratified sampling and systematic sampling are records-based. These records-based sampling methods require the complete pass over the whole dataset. Hence, they are commonly used for small or medium scale datasets on single computer. Even though the whole dataset is already assigned on a distributed system, it is very difficult to get a high-quality partition of the original dataset [46]. In the era of big data, data profiling tasks can be carried out on distributed systems, e.g., data profiling tasks on HDFS data. Therefore, block-based sampling proposed in [45] can be used as a promising sampling method for data stored in distributed machines.

He et al. [45] propose a block-based sampling method for large-scale dataset. They take block-based sampling as one of components of their big data learning framework which is called the asymptotic ensemble learning framework [47]. However, the block-based sampling method is suitable for data that is randomly ordered and not for those records that are stored in an orderly manner. In order to solve this problem, they propose a general block-based sampling method named I-sampling.

I-sampling contains four steps to get sample. Firstly, I-sampling divides large-scale dataset into non-overlapping primary data blocks  $A_i$ . Secondly, I-sampling shuffles primary data blocks  $A_i$  to get shuffling data blocks  $B_i$ . The purpose of shuffling is to disrupt the order of original data. Thirdly, I-sampling randomly selects data from  $B_i$  and put it into the basic blocks to get a block pool  $C$ . Finally, a certain number of basic blocks are randomly selected from the block pool, and the data in these blocks is taken as sample. In experiments, He et al. [45] demonstrate that the block-based sampling has the basically equal means and variances with simple random sampling. Besides, the distribution of I-sampling data is approximately the same with original dataset. And RMSEs of extreme learning machine based on records-based random sampling and I-sampling are nearly the same. The processing flow of I-sampling is shown in Figure 3.

As a matter of fact, data contains uncertainty in many applications. For example, when we do experiments, such as sampling, uncertainty occurs because there are many potential results for sampling. Uncertainty means the diversity of potential outcomes, which is unknown to us. And dealing with big data with uncertainty distribution is one of the most important

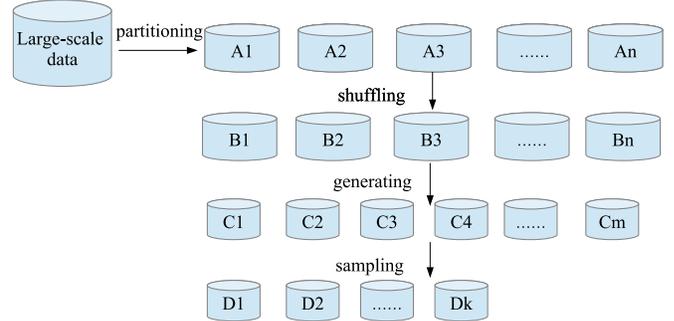


Fig. 3: I-sampling workflow [45].

issues of big data research [44]. The sample quality affects the accuracy of data profiling. The following example shows how to use MapReduce to accelerate sampling from a big data set with uncertainty distribution, and select Minimal Consistent Subset with better sample quality. Minimal Consistent Subset (MCS) is a consistent subset with a minimum number of elements.

He et al. [44] use MapReduce to sample from the data which contains uncertainty. They propose a Parallel Sampling method based on Hyper Surface (PSHS) for big data with uncertainty distribution to get the MCS of the original sample set. Hyper Surface Classification (HSC) is a classification method based on Jordan Curve Theorem and put forward by He et al. [48]. The MCS of HSC is a sample subset by selecting one and only one representative sample from each unit included in the hyper surface. Some samples in the MCS are replaceable, while others are not, leading to the uncertainty of elements in MCS [44]. Because of large-scale of data, they use MapReduce for parallel sampling. MapReduce is a well-known distributed computing framework for big data today.

PSHS algorithm needs to execute three kinds of MapReduce jobs. In the first task, based on the value of each dimension of the data, the map function places each sample in the region to which it belongs. The reduce function determines whether this area is pure, and labels each area with a corresponding label: pure or impure. In the second task, the corresponding decision tree is generated and the samples that have no effect on the generated decision tree must be removed. The third task is the sampling task, where the map function is to place each sample in the pure region it belongs to according to the rules. In pure regions, these samples have the same effect on building the classifier, hence the reduce task is to randomly select one and only one from each region for building the MCS. The Minimal Consistent Subset selected by this parallel sampling is a good

representation of the original data.

### C. Determination of Sample Size

It is very important to select effective samples [49]. If the sample size is too small, it may get an incorrect conclusion. If the sample size is too large, the calculation time is too long. Therefore, when performing machine learning algorithms, data mining algorithms or data profiling tasks on large-scale dataset, how to choose the appropriate sampling method and determine the sample size are important factors in determining whether the correct result (within the allowable error range) can be obtained.

There are some classic methods to determine the sample size. Singh and Masuku [10] have detailed and summarized these traditional methods. For example, you can take the sample size in other similar studies as the size of the sample in your study. Furthermore, you can determine the size of the sample according to the published tables. These tables determine the size of sample according to some given evaluation indicators and the size of the original dataset. Some commonly used evaluation indicators include the level of precision, the level of confidence or risk, the degree of variability, etc. Another method to assure size of sample is to calculate the size of sample according to some simple calculation formulas, which calculate the size of sample based on sampling error, confidence, and P-value. A simple formula (1) for estimating the sample size [10] is as follows, where  $n$  is the sample size,  $N$  is the amount of raw data,  $e$  is the level of precision, a 95% confidence level and  $P = .5$  are assumed.:

$$n = \frac{N}{1 + N * e^2} \quad (1)$$

When data mining algorithms are performed based on massive amounts of data, much of current research prefers to scale up data mining algorithms to deal with computational (time and memory) constraints, but some scholars focus on selecting how many samples to conduct data mining algorithms. In data mining algorithms, a common formula used to estimate the number of samples is Probably Close Enough (PCE). The convergence conditions are determined using PCE standard to obtain the best sample size for sampling. PCE is calculated as formula (2).

$$Pr[|acc(D) - acc(D_i)| \geq \epsilon] \leq \delta \quad (2)$$

Where  $D_i$  represents sample data,  $D$  represents all data,  $\epsilon$  represents the error range threshold of accuracy, and  $\delta$  represents probability.

Furthermore, Satyanarayana [11] proposes a dynamic adaptive sampling method for estimating the amount of data required for the learning curve to converge at each iteration. The author applies Chebyshev inequality to derive an expression that will estimate the number of instances at each iteration, which takes advantage of classification accuracy in order to get more precise estimates of the next sample. The expression is formula (3), where  $D_i$  is sample under consideration,  $acc(x_i)$  is classification accuracy of each instance,  $\epsilon$  is approximation parameter and  $\delta$  is probability of failure. And Satyanarayana [11] uses the formula (4) to check convergence at each

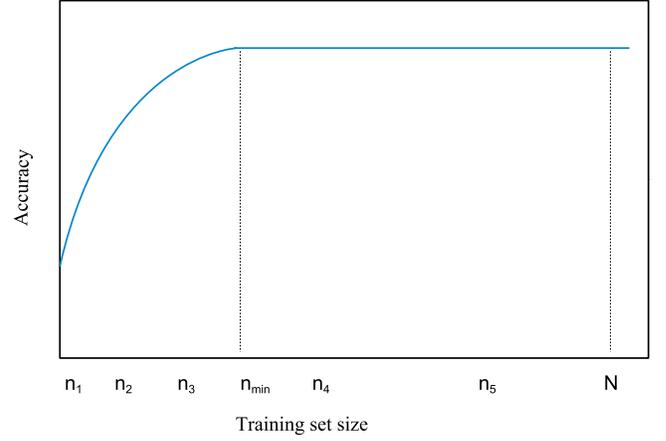


Fig. 4: Learning curves [50].

iteration, where  $D_i$  is the sample of current iteration and  $D_{i-1}$  is the sample of last iteration.

$$m \geq \frac{2}{\frac{1}{|D_i|} \sum_{i=1}^{D_i} acc(x_i)} \left[ \frac{1}{\epsilon^2} \log \frac{1}{\delta} \right] \quad (3)$$

$$\left| \frac{1}{|D_i|} \sum_{i=1}^{D_i} acc(x_i) - \frac{1}{|D_{i-1}|} \sum_{i=1}^{D_{i-1}} acc(x_i) \right| < \epsilon \quad (4)$$

When sampling is used in machine learning, the most appropriate number of samples is to make the accuracy rate reach the maximum value and increasing the number of samples can no longer improve the accuracy of the learning algorithm. The corresponding figure is Figure 4, where  $n_{min}$  is the minimum sample size. In this case, there is a method for determining the minimum number of samples called sequential sampling. Sequential sampling refers to sample sequentially and stop sampling until a certain criterion is met. John and Langley [12] propose a method called Arithmetic Sampling. This method uses a schedule  $Sa = (n_0, n_0 + n_\delta, n_0 + 2n_\delta, n_0 + 3n_\delta, \dots, N)$  to find the minimum sample size, where  $n_0$  is the starting sample size and  $n_\delta$  is fixed interval. Provost et al. [50] think that the main drawback is that if the minimum size of sample is a large multiple of  $n_\delta$ , it will take many runs to reach convergence. Obviously, if  $n_\delta$  is too small, it will take many iterations to get convergence, and if  $n_\delta$  is too large, it may skip the optimal size of sample.

Singh et al. [13] propose another sequential sampling strategy for classification problem. They mention that data for training machine learning models typically originates from computer experiments such as simulations. And computer simulations are often computationally expensive. In order to ease the computation pressure, they use sampling to get as little data as possible. The sequential sampling starts with an initial small data set  $X^\delta$ , and it will iteratively increase the sample by taking training points at well-chosen locations  $\delta$  in the input space until stopping criteria is reached. The sequential sampling strategy chooses a representative set of data samples that focuses the sampling on those locations in the input space where the class labels are changing more

rapidly, while making sure that no class regions are missed [13]. The sample update formula is formula (6) where class labels  $Y^\delta$  obtained by formula (5) are result of simulator evaluates  $X^\delta$ . With sequential sampling strategy, small and high quality samples can be obtained.

$$Y^\delta := f(X^\delta) \quad (5)$$

$$S := S \cup (X^\delta, Y^\delta) \quad (6)$$

#### D. Sampling Error and Sampling Bias

Sampling error is when a randomly chosen sample does not reflect the underlying population purely by chance and sampling bias is when the sample is not randomly chosen at all [51]. Sampling bias is one of the causes of sampling error. These two are often confused by some scholars. Sampling bias is caused by the failure of the sampling design, which cannot truly extract the sample randomly from the population [52].

There is a typical case of sampling error. The large Nurses Health Study tracked 48,470 postmenopausal women for 10 consecutive years, aged between 30 and 63 years old. The study concluded that hormone replacement therapy can reduce the incidence of severe coronary heart disease by nearly half [53]. Despite the large sample size, the study failed to recognize the atypical nature of the sample and the confusion of estrogen therapy with other active health habits [54]. This also illustrates the importance of proper sampling methods and the collected samples to get the right conclusions.

To be able to correctly select the samples that represent the original data, Kim and Wang [55] focus on and solve the problem of selection bias in the process of sampling. Since big data is susceptible to selection bias, they propose two ways to reduce the selection bias. One is based on the inverse sampling method. This method is divided into two stages. The first stage is to sample directly from the big data. The sample is easily affected by the selection bias, thereby it is necessary to calculate the importance of each element in the sample to determine selection bias. In the second stage, the data sampled from the first stage is re-sampled according to the importance weight of each element. In this way, they have achieved the goal of realizing the correction of the selection deviation. The other is the idea of using data integration. They propose to use the survey data and big data to correct the selection bias by means of the auxiliary information of survey data.

From the perspective of official statisticians, Tam et al. [56] believe that big data is challenged by self-selection bias. Self-selection bias causes biased sample with non-probability sampling. Inferences from big data with this bias will be affected. Thus, they outline methods for adjusting self-selection bias to estimate proportions, e.g., using pseudo weights and super population models [57].

As a matter of fact, the case of 2016 US presidential election studied in [58] is precisely because of the existence of self-selection bias, which ultimately leads to data deceiving us. Therefore, to get the correct conclusion from the data, you need to ensure the quality of the data. Probability sampling can guarantee the quality of the data. When probability sampling cannot be satisfied, the data will be affected by Law of Large

Populations (LLP). The large amount of data  $N$  will affect our estimation error. In summary, when doing data sampling, data quality must be taken into account, and those high quality data sets should be given higher weight. This will prevent our statistical inferences from being affected.

#### IV. SAMPLING FOR SINGLE COLUMN DATA PROFILING

Single column data profiling tasks are divided into cardinalities, value distributions, patterns, data types, and domains [79]. Table II [16] lists typical metadata that may result from single-column data profiling. For some single-column data profiling tasks, such as decimals which calculates maximum number of decimals in numeric values, simple sampling methods cannot guarantee reliable results. And for identifying a domain of one column, it is often more difficult and not fully automated [80]. Among them, cardinality, histograms and quantiles are often used for query optimizers, therefore sampling techniques are more commonly used in these tasks. Specifically, in Section IV-A, we introduce sampling for cardinality estimation. Section IV-B presents sampling for value distribution. More advanced statistics include the probabilistic correlations on text attributes [81].

##### A. Sampling for Cardinality Estimation

Cardinalities or counts of values in a column are the most basic form of metadata [16]. Cardinalities usually include number of rows, number of null values and number of distinct values, which is the most important type of metadata [82]. For some tasks, such as number of rows and number of null values, a single pass over a column can get the exact result. However, finding the number of distinct values may require to sort or hash the value of column [80]. Similarly, when facing large data sets, statistics of the number of distinct values of an attribute have to face the pressure of memory and calculation. Therefore, the estimation of the number of distinct values based on sampling has been studied [59]–[61].

Haas et al. [59] propose several sampling-based estimators to estimate the number of different values of an attribute in a relational database. They use a large number of attribute value distributions from various actual databases to compare these new estimators with those in databases and statistical literature. Their experimental results prove that no estimator is optimal for all attribute value distributions. And from their experimental results, it can be seen that the larger the sampling fraction, the smaller the estimated mean absolute deviation will be. They therefore propose a sampling-based hybrid estimator  $\hat{D}_{hybrid}$  and get the highest precision on average at a given sampling fraction.

Similar to Haas et al., Charikar et al. [60] also obtain a negative result in the experiment that no estimator based on sampling can guarantee small errors on the input data of different distributions, unless a larger sampling fraction is performed on the input data. They therefore propose a new estimator Guaranteed-Error Estimator (GEE), which is provably optimal. Although its error on the input of different distributions is small, it does not make use of the knowledge of different distributions. For example, in the case of low-skew

TABLE II: Overview of single-column profiling tasks [16]

Category	Task	Description
Cardinalities	num-rows	Number of rows
	value length	Measurements of value lengths (minimum, maximum, median, and average)
	null values	Number or percentage of null values
	distinct	Number of distinct values; sometimes called "cardinality"
Value distributions	uniqueness	Number of distinct values divided by the number of rows
	histogram	Frequency histograms (equi-width, equi-depth, etc.)
	constancy	Frequency of most frequent value divided by number of rows
	quartiles	Three points that divide the (numeric) values into four equal groups
Patterns, data types, and domains	first digit	Distribution of first digit in numeric values
	basic type	Generic data type, such as numeric, alphabetic, alphanumeric, date, time
	data type	Concrete DBMS-specific data type, such as varchar, timestamp.
	size	Maximum number of digits in numeric values
	decimals	Maximum number of decimals in numeric values
	patterns	Histogram of value patterns (Aa9...)
	data class	Semantic, generic data type, such as code, indicator, text, date/time, quantity, identifier
	domain	Classification of semantic domain, such as credit card, first name, city, phenotype

TABLE III: Summary of sampling for big data profiling tasks

Data Profiling	Sampling-based method	
Single column	Cardinality Estimation	$\hat{D}_{hybrid}$ [59], GEE [60], AE [60], Distinct sampling [61]
	Histograms	Random sampling [62], Backing sample [63]
	Quantiles	Non-uniform random sampling [64], Improved random sampling [65]
Multiple columns	Correlations and association rules	Sequential random sampling without replacement [14], Two-phased sampling [4], ISbFIM [66]
	Clusters and outliers	Biased sampling [67]
	Summaries and sketches	Error-bounded stratified sampling [68], [69]
	Regression analysis	IBOSS [70], Random sampling without replacement [71]
Dependency	Uniqueness	GORDIAN [72], HCA-Gordian [73]
	Functional dependencies	AID-FD [74], HYFD [75], CORDS [76], BRRSC [77]
	Inclusion dependencies	FAIDA [78]

data with a large number of distinct values, GEE performs not very well in practice. They further propose a new heuristic version of GEE called Adaptive Estimator (AE), which avoids the problems encountered by GEE.

Different from the previous research using random sampling, Gibbons [61] proposes distinct sampling to accurately estimate the number of distinct values. Distinct sampling can collect distinct samples in a single scan of the data, and the samples can be kept up to date in the state of data deletions and insertions. On a truly confidential data set Call-center, distinct sampling uses only 1% of the data, and can achieve a relative error of 1% -10%, while increasing the speed of report generation by 2-4 orders of magnitude. They compare distinct sampling with GEE, AE in the experiment and prove that in real-world data sets, distinct sampling performs much better than GEE and AE.

It is worth noting that Harmouch and Naumann [82] conduct an experimental survey on cardinality estimation. In the experiment, they use the GEE [60] as an example of evaluation. They perform experiments on synthetic and real-world data sets. It can be seen from the experimental results that the larger the sampling fraction, the smaller the average estimation relative error. And when GEE wants to reach 1% relative error, it needs to collect more than 90% of the data. In conclusion, when faced with large data sets, cardinality estimation requires high memory, and sampling can reduce memory consumption, but cannot guarantee reasonable accuracy all input distributions.

### B. Sampling for Value Distribution

Value distribution is a very important part of single-column data profiling. Histogram and quantile are two typical forms used to represent value distribution. The histogram is used to describe the distribution of data, while quantile refers to dividing the data into several equal parts.

1) *Sampling for Histogram Construction:* Many commercial database systems maintain histograms to summarize the contents of large relations and permit efficient estimation of query result sizes for use in query optimizers [63]. Histogram can be used to describe the frequency distribution of attributes of interest, which groups attributes values into buckets and approximates true attribute values and their frequencies based on summary statistics maintained in each bucket [83]. However, the database is updated frequently, hence the histogram also needs to be updated accordingly. Recalculating histograms is expensive and unwise for large relations.

Gibbons et al. [63] propose sampling-based approaches for incremental maintenance of approximate histograms. They use a "backing sample" to update histograms. Backing sample is a random sample of the relation which is kept up to date in the presence of databases updates, which is generated by uniform random sampling. Therefore, random sampling can help to speed up histogram re-computation. For example, SQL Server recomputes histograms based on a random sample from relations [62].

Chaudhuri et al. [62] focus on how much sample is enough to construct a histogram. They propose a new error metric called the max error metric for approximate equip-depth

histogram. The max error metric is formula (7) shown below, where  $b_j$  is number of values in bucket  $j$ ,  $k$  is the number of buckets and  $n$  is the number of records. A  $k$ -histogram is said to be a  $\delta$ -deviant histogram when  $\Delta_{max} \leq \delta$ . And size of sample  $r$  is calculated as the following formula (8), where  $\delta \leq \frac{n}{k}$  and  $\gamma$  is predefined probability.

$$\Delta_{max} = \max_{1 \leq j \leq k} |b_j - \frac{n}{k}| \quad (7)$$

$$r \geq \frac{4n^2 \ln(\frac{2n}{\gamma})}{k\delta^2} \quad (8)$$

As mentioned above, the histogram can be used to represent the distribution of data. In exploratory data analysis, analysts want to find a specific distribution from a large number of candidate histograms. The traditional approach is "generate and test", i.e., generating all possible histograms, and then testing whether these histograms meet the requirements. This approach is undesirable when the data set is large. Therefore, Macke et al. [84] propose a sampling-based approach to identify the top  $k$  closest histograms called HistSim. The idea of HistSim is using random sampling method without replacement to collect samples for histogram constructing. Then they normalize the representation vector of the histogram, and use  $l1$  distance to calculate the similarity. Furthermore, they propose FastMatch, which combines HistSim and block-based sampling method, and obtain near-perfect accuracy with up to 35 speedup over approaches that do not use sampling on several real-world datasets in the experiment.

2) *Sampling for Quantile Finding*: Quantiles can be used to represent the distribution of single column value. Quantiles are used by query optimizers to provide selectivity estimates for simple predicates on table values [85]. Calculating exact quantiles on large data sets is time consuming and requires a lot of memory. For example, quantile finding algorithm in [86] requires to store at least  $N/2$  data elements to find the median, which is memory unacceptable for large-scale data.

Therefore, Manku et al. [64] present a novel non-uniform random sampling to find approximate quantile. They apply non-uniform random sampling to reduce memory requirements. Non-uniform means that the probability of selecting each element in the input is different. They set the earlier elements in the input sequence with larger probability than those arrive later. And the process of quantile finding is shown in Figure 5. When the data arrives, they randomly select an element in each data block and put it into buffers. Then based on sample, deterministic algorithms are performed to find quantiles.

However, simply using random sampling method and calculating the quantiles on the sample may not be accurate enough on sensor networks. Hence, Huang et al. [65] propose a new sampling-based quantile computation algorithm for sensor networks to reduce the communication cost. To improve accuracy, they augment the random sample with additional information about the data. They analyze how to add additional information to the random sample under the flat model and the tree model. For example, in the flat model, each node first samples each data value independently with a certain probability  $p$  and computes its local rank. Then the samples and their local ranks

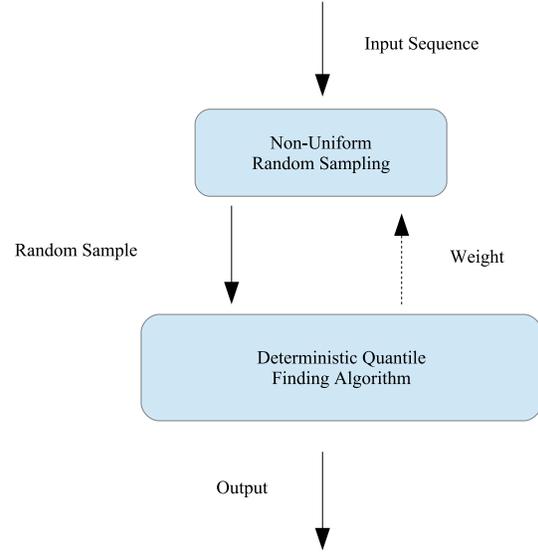


Fig. 5: Sampling for quantile finding [64].

are sent to base station. The base station estimates rank for any value it receives and then quantile queries can be solved. In the end, they prove through experiments that the quantile computation in Sensor Networks based on this new sampling method reduces one to two orders of magnitude in terms of the total communication cost compared with the previous method.

## V. SAMPLING FOR MULTIPLE COLUMNS DATA PROFILING

As shown in Figure 1, the content of the multiple columns data profiling tasks includes association rule mining [87], clusters and outliers [88], summaries and sketches [16]. Besides, statistical methods such as regression analysis [89] can be used to perform multiple columns analysis, analyzing the relationship between these columns. Specifically, in Section V-A, we investigate sampling for discovering association rules. Section V-B presents the content of sampling for clusters and outliers. And sampling for summaries and sketches is introduced in Section V-C. Then, in Section V-D, we introduce sampling for helping perform regression analysis.

### A. Sampling for Discovering Association Rules

The discovery of association rules is a typical problem in data profiling for multiple columns. The algorithm currently used to find association rules needs to scan the database several times. For large data sets, the time overhead of scanning several times is hard to accept. Large amount of data leads to input data, intermediate results and output patterns can be too large to fit into memory and prevents many algorithms from executing [66]. Some scholars have proposed using parallel or distributed methods to solve the problem of data volume [90], [91]. But it is difficult to design parallel or distributed algorithms.

Therefore, Zaki et al. [14] use sampling to get samples of transaction and find the association rules based on the obtained samples. They take sequential random sampling without replacement as their sampling method and use Chernoff bounds

to obtain sample size. Finally, they experimentally prove that sampling can speed up the discovery of association rules by more than an order of magnitude and provide high accuracy for association rules.

Chen et al. [4] propose a two-phased sampling-based algorithm to discover association rules in large databases. At the first stage, a large initial sample of transactions is randomly selected from databases, which is applied to calculate support of each individual item. And these estimated supports are used to trim the initial sample to a smaller final sample  $S_0$ . At the second stage, association-rule algorithm is performed against the final sample  $S_0$  to get association rules according to provided minimum support and confidence. In the experiment, the authors prove 90-95% accuracy obtained using the final sample  $S_0$  and the size of sample is only 15-33% of the whole databases. This again proves that sampling can be used to speed up data analysis and big data profiling.

Wu et al. [66] propose an Iterative Sampling based Frequent Itemset Mining method called ISbFIM. The same as [14], Wu et al. [66] use random sampling as the sampling method. But the difference is that they use iterative sampling to get multiple subsets and find frequent items from these subsets. They can guarantee that the most frequent patterns for the entire data set have been enumerated and implement a Map-Reduce version of ISbFIM to demonstrate its scalability on big data. Because the volume of input data is reduced, the problem that input data, intermediate results, or the final frequent items cannot be loaded into memory is solved. And the traditional exhaustive search-based algorithms like Apriori can be fitted for big data context.

### B. Sampling for Clustering and Anomaly Detection

Clustering is to segment similar records into the same group according to certain characteristics, and those records that cannot be classified into any group may be abnormal points. The challenge that clustering technology encounters in the era of big data is also the problem of data volume, and the clustering operation itself consumes a lot of calculations. Shirkhoshidi et al. [92] divide big data clustering into two categories: single-machine clustering and multiple-machine clustering. Single column reduces the amount of data by using data reduction methods, e.g., sampling and dimensionality reduction. Multi-machine clustering refers to the use of parallel distributed computing frameworks, e.g., MapReduce and cluster resources to increase computing power.

Kollios et al. [67] propose biased sampling to speed up clustering and anomaly detection on big data. Unlike the previous work, they consider the data characteristics and analysis goals during the sampling process. Based on the tasks of clustering and anomaly detection, Kollios et al. [67] consider the data density problem in the dataset. They propose a biased sampling method to improve the accuracy of clustering and anomaly detection. The biased sampling is to make the data points in each cluster and the abnormal points have a higher probability of being selected. In order to achieve this goal, they use the density estimation method to estimate density around the data points. In the experiment, they prove that density-

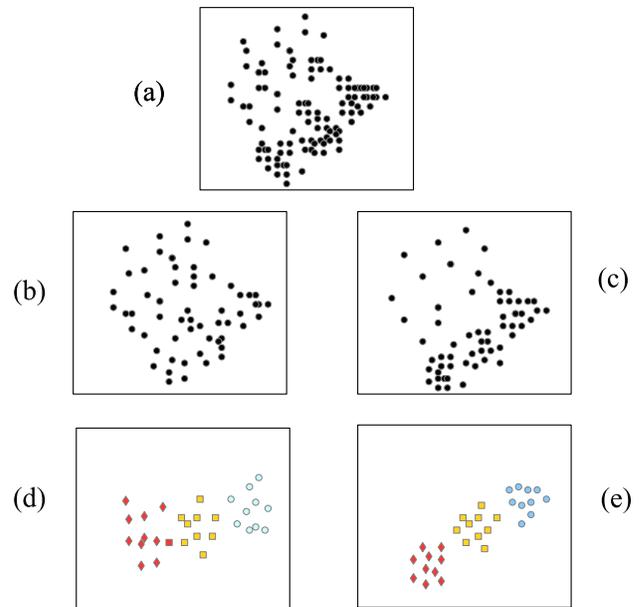


Fig. 6: Application of biased sampling in clustering tasks [67].

based sampling has a better effect on clustering than uniform sampling.

Figure 6 shows the use of biased samples in clustering. Figure 6(a) is the distribution of the original data and there are three classes with higher density. Figure 6(b) is the result of random sampling on the original data set. Figure 6(c) is the result of applying the biased sampling to the original data. Figure 6(d) shows 10 data points selected from each of the three categories clustered based on the random sampling, and Figure 6(e) shows 10 data points selected from each of the three categories clustered based on the biased sampling method. After comparison with the categories in the original data, it is found that the clustering results of the biased samples are more accurate.

### C. Sampling for Summaries and Sketches

Summaries or sketches can be performed by sampling or hashing data values to a smaller domain [16]. Although different scholars have applied different sampling algorithms, the most commonly used sampling algorithm among data scientists is random sampling [69]. The main reason is that random sampling is the best and easiest to use, which is the only technique commonly used by data scientists to quickly gain insights from big data sets.

Rojas et al. [69] first interview 22 data scientists working on large data sets and find that they basically use random sampling or pseudo-random sampling. Certainly, these data scientists believe that other sampling techniques may achieve better results than random sampling. These scientists perform a data exploration task that used different sampling methods to support classification of more than 2 million generated samples from data records of Wikipedia article edit. Research has shown that sampling techniques other than random sampling can generate insights into the data, which can help focus on the different characteristics of the data without affecting

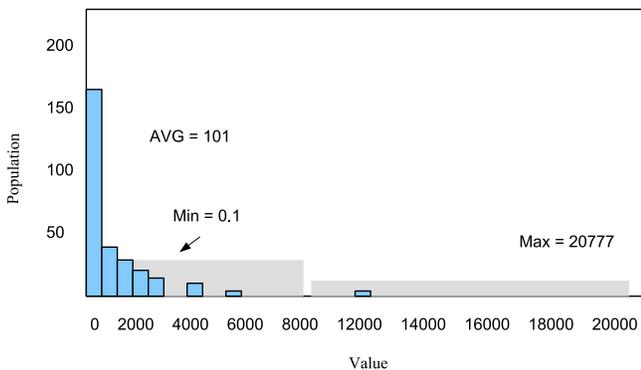


Fig. 7: Sparseness of one representative production data [93].

the quality of data exploration and helping people understand the data. This shows that with the application of sampling, Summaries or sketches of data can be created to help scientist observe and understand the data.

Aggregated queries are also a way to generate summaries of data. Aggregate queries are computationally expensive which need to traverse the data. In the era of big data, a single machine often cannot make such a large amount of data. Therefore, aggregate queries for big data are often performed on distributed systems that scales to thousands of machines. The commonly used distributed computing frameworks are Hadoop, spark, etc. Although distributed systems provide tremendous parallelism to improve performance, the processing cost of aggregated queries remains high [68]. Investigation in one cluster of [68] reveals that 90% of 2,000 data mining jobs are aggregation queries. These queries consume two-thousand machine hours on average, and some of them take up to 10 hours.

Therefore, Yan et al. [68] use sampling technique to reduce the amount of data. When error bounds cannot be compromised and data is sparse, they think that conventional uniform sampling often yields high sampling rates and thus deliver limited or no performance gains. For example, uniform sampling with 20% error bound and 95% confidence needs to consume 99.91% of the data whose distribution is shown in Figure 7. Hence, they propose error-bounded stratified sampling, which is a variant of stratified sampling [93] and relies on the insight, i.e., prior knowledge of data distribution, to reduce sample size. Error bound means that the real value has a large probability of falling within an interval. Sparse data means that the data is generally limited but wide-ranging.

Taking the data distribution in Figure 7 as an example, error-bounded stratified sampling can divide the data into two groups. One group covers the header data and the other covers the tail data. Because the data range of the first group is small, the sampling rate is also small. Although the data range of the second group is large, the data basically falls in the first group. Even if the data of the second group is all taken as a sample, the overall sampling rate is still low. It is worth mentioning that the technique has been implemented into Microsoft internal search query platform.

#### D. Sampling for Regression Analysis

Statistical analysis such as regression analysis can be used to analyze the relationship between multiple columns in a relation. Sauter [94] think that statistics are learned from data. Statistics methods are often used for data profiling, which have encountered the problem of excessive data volume in the era of big data. Statistical analysis of the entire big data set requires a certain amount of calculation and time.

Under the computational pressure of large data sets, many traditional statistical methods are no longer applicable. Although sampling can help with data reduction, how to avoid sampling errors caused by sampling needs to be considered. For example, [70] mention that in the context of linear regression, traditional sub-sampling methods are prone to introduce sampling errors and affect the covariance matrix of the estimator. Hence, they propose information-based optimal subdata selection method called IBOSS. The goal of IBOSS is to select data points that are informative so that small-sized subdata retains most of the information contained in the complete data. Simulation experiments prove that IBOSS is faster and suitable for distributed parallel computing.

Jun et al. [71] propose to use sampling to divide big data into some sub data sets in regression problem for reducing the computing burden. The traditional statistical analysis of big data is to sample from big data, and then perform statistical analysis on the sample to infer the population. Jun et al. [71] divide the big data closed to population into some sub data sets with small size closed to sample which is proper for big data analysis. They treat the entire data set as a population and the sub set as a sample to reduce computing burden. And they select regression analysis to perform experiments. The traditional processing is shown in Figure 8, and their design is shown in Figure 9.

Their design consists of three steps: the first step is to first generate  $M$  sub-data sets using random samples without replacement; the second step is to calculate the regression parameters of each sub-data set and calculate the average of regression parameters of the  $M$  sub-data sets; the third step is to use the averaged parameters obtained in the second step to estimate regression parameters on the entire data set. This design that combines sampling and parallel processing helps them speed up regression analysis on big data. By experimenting with the data set from the simulation and UCI machine learning repository, the author proves that the regression parameters obtained by distributed calculation on random samples are close to the regression parameters calculated on entire data set. This provides a reference for statistical analysis on the entire large data set.

## VI. SAMPLING FOR DEPENDENCIES

A dependency is a metadata that describes the relationship between columns in relation, based on either value equality or similarity [95]. There are many use cases for dependencies. For example, unique column combinations are used for finding key attributes in relation [72], and functional dependencies can be used for schema normalization [96] or consistent query answering [97], while inclusion dependencies can suggest how

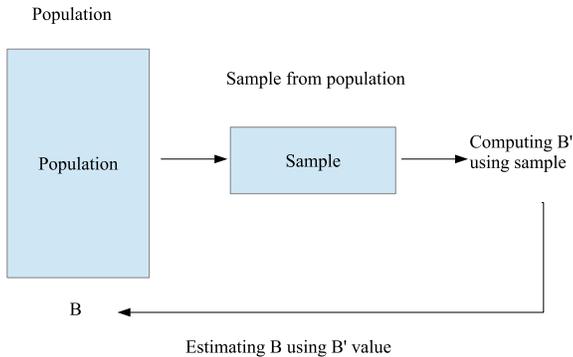


Fig. 8: Traditional big data regression analysis [71].

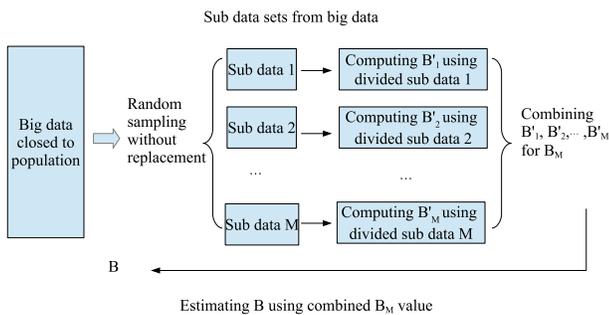


Fig. 9: Sampling-based partitioning big data regression analysis [71].

to join two relations [16]. Inclusion dependencies together with functional dependencies form the most important data dependencies used in practice [98]. But discovery of dependencies is time consuming and memory consuming. Many functional dependencies discovery algorithms are not suitable for large data sets. Sampling could be employed to estimate the support and confidence measures of data dependencies [99], [100]. By sampling, you can select a small enough representative data set from the big data set. Hence, the choice of sampling method is very important, which help to ensure that the estimated inaccuracy rate is below a predefined bound with high confidence. Specifically, based on the classification for dependency in [16], we investigate sampling for unique column combinations in Section VI-A, functional dependency in Section VI-B, inclusion dependency in Section VI-C.

#### A. Sampling for Discovery of Unique Column Combinations

An important goal in data profiling is to find the right key for the relational table, e.g., primary key. The step before key discovery is to discover unique column combinations. Unique column combinations are sets of columns whose values uniquely identify rows, which is an important data profiling task [101]. But discovery of unique column combinations is computationally expensive, which is suitable for small dataset or samples of large dataset. For large data set, sampling is a promising method for knowledge discovery [102]. Based on sampling-based knowledge discovery, it is necessary to first select samples from the entire data set and obtain knowledge

from the samples, and then use the entire data set to verify that the acquired knowledge is correct.

A typical algorithm for identifying key attributes is GORDIAN proposed by Sismanis et al. [72]. The main idea of GORDIAN is to turn the problem of keys identification into cube computation problem, and then find non-keys through cube computation. Finally, GORDIAN calculates the complement of the non-keys set to obtain the desired set of keys. Therefore, the GORDIAN algorithm can be divided into three steps: (i) create the prefix tree through a single pass over the data; (ii) find maximal non-uniques by traversing the prefix tree with pruning; (iii) get minimal keys from set of maximal non-uniques. In order to make GORDIAN scalable to large datasets, Sismanis et al. combine GORDIAN with sampling. Experiments have shown that sampling-based GORDIAN can find all true keys and approximate keys using only a relatively small number of samples.

GORDIAN algorithm is further developed by Abedjan and Naumann [73] to discover unique column combinations. Since the existing algorithms are either too violent or have high memory requirements and cannot be applied to big data sets. A hybrid solution HCA-Gordian, which combines Gordian algorithm [72] and their new algorithm the Histogram-Count-based Apriori Algorithm (HCA), is proposed by Abedjan and Naumann [73] to discover unique column combinations. GORDIAN algorithm is used to find composite keys and the HCA is an optimized bottom-up algorithm which takes efficient candidate generation and statistics-based pruning methods. HCA-Gordian performs Gordian algorithm on a smaller sample of table to discover non-uniques and non-uniques will be used as pruning candidates when executing HCA on the entire table.

In the experiment setup, the sample size for the preprocessing step of non-unique discover is always 10,000 tuple sample. Especially when the amount of data is large and the number of unique is small, the runtime of HCA-Gordian is lower than Gordian. For example, when using real world tables for experiments, search speed of HCA-Gordian is four times faster than Gordian. And as the data set grows larger, e.g., the National file contains 1,394,725 tuples, Gordian takes too long to run, while HCA-Gordian only takes 115 seconds to complete. In addition, When the number of detected non-uniques is high, the discovery effect of HCA-Gordian is better than Gordian.

#### B. Sampling for Functional Dependencies

A functional dependency refers to a set of attributes in a relationship that determines another set of attributes. For example, there is such a functional dependency  $A \rightarrow B$ , which means that any two records in the relationship, when their values on the attribute set  $A$  are equal, the values on the attribute set  $B$  must be equal. Bleifuss et al. [74] propose an approximate discovery strategy AID-FD (Approximate Iterative Discovery of FDs) which sacrifices a certain correct rate in exchange for performance improvement. AID-FD uses an incremental, focused sampling of tuple pairs to deduce non-FDs until user-configured termination criterion is met. The authors have demonstrated in experiments that the AID-FD

method uses only 2%-40% of the time of the exact algorithm when processing the same data set, but finds more than 99% of the functional dependencies.

Papenbrock and Naumann [75] mention that today's various functional dependencies discovery algorithms do not have the ability to process more than 50 columns and 1 million rows of data. Thus, they propose the sampling-based FD discovery algorithm HYFD. And there are three properties in sampling-based FD discovery algorithms: Completeness, Minimality, Proximity, which are important for HYFD. HYFD combines column-efficient FD induction techniques with row-efficient FD search techniques in two phases. In Phase 1, they apply focused sampling techniques to select samples with a possibly large impact on the results precision and produce a set of FD candidates based on samples. In Phase 2, the algorithm applies row-efficient FD search techniques to validate the FD candidates produced in Phase 1. The sampling method allows functional dependencies discovery algorithms to be extended to large data sets.

In experiments, when the data set is not very large, the runtime of HYFD is almost all lower than other algorithms. When the data set exceeds 50 columns and 10 million rows, HYFD can get the result through a few days of calculation. However, other algorithms cannot complete the calculation, because the time complexity for these algorithms is exponential. This again demonstrates that sampling is important for data profiling, e.g., FD discovery.

In the above, we mention that using focused sampling to find functional dependencies. In this section, we will mention the use of random sampling to find soft functional dependencies. The so-called "soft" functional dependency is relative to the "hard" functional dependency. A "hard" functional dependency means that the entire relationship satisfies the functional dependency, while a "soft" functional dependency means that the entire relationship is almost satisfied, or that there is a high probability of satisfying the functional dependency.

Ilyas et al. [76] propose sampling-based CORDS, which means that automatic discovery of correlations and soft functional dependencies between columns, to find approximate dependencies. Among them, correlation refers to the general statistical dependence, while soft functional dependence refers to that value of attribute C1 determines the value of attribute C2 with high probability. CORDS use enumeration to generate pairs of columns that may be associated, and heuristically cuts out those unrelated column pairs with high probable. CORDS apply random sampling with replacement to generate sample. In the implementation of CORDS, they only use a few hundred rows of sample data, and the sample size is independent of the data size. In the experiment to evaluate the advantages of applying CORDS, where run a workload of 300 queries on the Accidents database, the median query execution time and worst query execution time with CORDS applied were better than those without CORDS. Hence, CORDS is efficient and scalable when it encounters large-scale dataset.

Approximate functional dependence is similar to the meaning of soft functional dependency. Approximate functional dependence requires the normal functional dependency to be satisfied by most tuples of relation  $R$  [103], [104]. Of course,

approximation functional dependencies contain exact functional dependencies that are satisfied throughout the relationship. As mentioned in [103], when the amount of data is large, the time for discovery of functional dependency will increase exponentially. Therefore, Kivinen and Mannila [103] propose to discover approximate dependencies by random sampling. In fact, sampling can be used not only to find approximate functional dependencies, but also to verify exact functional dependencies [104]. If the exact functional dependency does not satisfy all the sample data, then the whole relationship is definitely not satisfied, hence such functional dependencies can be removed.

Functional dependencies are satisfied for all tuples in the relation, while conditional functional dependencies (CFDs) is to hold on the subset of tuples that satisfies some patterns [105]. And CFDs can be used for data cleaning [105], [106]. Fan et al. [107] propose three methods for conditional functional dependencies discovery. However, when the size of data set is large, no dependency discovery algorithms scale very well to discover minimal conditional functional dependencies.

When mining CFDs on big data, the volume issue of big data has to be solved. Li et al. [77] develop the sampling algorithms to obtain a small representative training set from large and low-quality datasets and discover CFDs on the samples. They use sampling technology for two reasons. One is that finding CFD needs to scan the data set multiple times, and sampling helps reduce the amount of data. The second is to use the sampling method to help them filter those dirty items on the low-quality data set and choose clean items as the training set. They define criteria for misleading tuples, which are dirty, incomplete or very similar to popular tuples. And then they design a Representative and Random Sampling for CFDs (BRRSC), which is similar to reservoir sampling [43]. The difference is that they combine the criteria defined above during the sampling process. Furthermore, they propose fault-tolerant CFDs discovery and conflict-resolution algorithms to find CFDs. Finally, experimental results show that their sampling-based CFD discovery algorithms can find valid CFD rules for billions of data in a reasonable time.

### C. Sampling-based Test for Inclusion Dependency Candidates

The definition of inclusion dependencies (INDs) is that the combination of values that appear in a set of attribute columns must also appear in another set of attribute columns [108]. Therefore, inclusion dependencies are often used to discover foreign keys [98]. However, discovery of inclusion dependencies is computationally expensive. One of the reasons is that the existing algorithms need to shuffle huge amounts of data to test inclusion dependencies candidates, which puts pressure on both computing and memory [78].

Under these circumstances, Kruse et al. [78] propose fast approximate discovery of inclusion dependencies (FAIDA). FAIDA can guarantee to find all INDs and only false positives with a low probability in order to balance efficiency and correctness. FAIDA uses algorithms [109], [110] of Apriori-style to generate inclusion dependencies candidates. The inverted index values and operates on a small sample of the input data.

The sampling algorithm is applied to each table to get each sample. Rather than use random sampling to get sample, they assure that sample table contains  $\min\{s, d_A\}$  distinct values for each column  $A$ , where  $s$  represents sample size and  $d_A$  represents number of distinct values in column  $A$ .

In their experiments, they set sample size to a default of 500. In order to verify the efficiency of FAIDA, Kruse et al. [78] compare FAIDA's runtime with the state-of-the-art algorithm for exact IND discovery BINDER [110] on multiple datasets. On four datasets, FAIDA is steadily 5 to 6 times faster than BINDER, and they generate and test almost the same number of IND candidates. Especially when one of the datasets reaches 79.4GB, BINDER takes 9 hours and 32 minutes to complete, while FAIDA only takes 1 hour and 47 minutes. Their evaluation shows that sampling-based FAIDA outperforms the state-of-the-art algorithm by a factor of up to six in terms of runtime without reporting any false positives.

## VII. SUMMARY AND FUTURE WORKS

Data in various fields are increasing on a large scale. Big data brings us new opportunities and challenges. Through data analysis and data mining of big data, we can get a lot of potential value. However, due to the large amount of data, it brings great challenges to the processing and storage. Therefore, data analysis, data mining or data profiling on large data sets have to face the pressure of calculation and time. Increasing computing power by using clusters of computers is one solution, but many times this is not the case, and designing distributed computing is often difficult. Hence, the application of data reduction techniques like sampling is very important. There are some mature research articles on data profiling and sampling, while little attention is paid to sampling and profiling over big data, therefore this article focuses on researching sampling and profiling in big data context. We first give a brief introduction of data profiling and introduce some important factors of sampling in detail. Then, according to the classification of data profiling in [16], we introduce the application of sampling in single column data profiling, multiple columns data profiling and dependency discovery. In conclusion, Table III summarizes the sampling for data profiling tasks investigated in survey, indicating the widespread use of sampling in data profiling.

The above survey on "sampling and profiling over big data" is mainly about relational databases, and rarely involves graph data or time series data. Since there is less research on sampling-based data profiling for graph data or time series data, we provide some future directions as follows.

### A. Sampling for Profiling Time Series Data

Many tasks on time series data need data profiling, e.g., matching heterogeneous events in a sequence [111], [112] with profiled patterns [113], [114], cleaning time series data [115] under speed constraints [22], or repairing timestamps according to the given temporal constraints [116] such as sequential dependencies [117]. All these studies use data profiling to detect and repair erroneous temporal data. The computational cost and time cost in large-scale temporal data

streams can be high. Therefore, sampling for profiling time series data is valuable and necessary.

In the time series data stream, we do not need to get exact results, e.g., when calculating the quantiles or probability distributions of speeds. Approximate results are valuable in time-series data streams, for example approximate probability distributions of speeds can also help us perform effective anomaly detection. In the sampling of time series data, statistical probability distributions of speeds are different from discovering quantiles. The speed of time series data depends on the adjacent time-series data points, which means that sampling for calculating speed of time series requires a set of data points in a window. Therefore, how to apply the sampling technology to the aforesaid data profiling task of time series data needs further experimental analysis and research.

### B. Sampling for Profiling Graph Data

Data profiling is also heavily used in graph data, e.g., using Petri Nets in process mining to recover missing events [118], [119] and clean event data [120], discovering keys for graphs and applying keys to study entity matching [121], or defining functional dependencies for graphs [25] and discovering them [122]. However, the above studies still seem to be difficult when encountering large graphs. Fan et al. [121] prove that entity matching is NP-complete for graphs and recursively defined keys for graphs bring more challenges. In this case, one has to design two parallel scalable algorithms, in MapReduce and a vertex-centric asynchronous model. In order to find Graph Functional Dependencies, Fan et al. [122] have to deal with large-scale graphs by designing effective pruning strategies, using parallel algorithms, and adding processors. As mentioned earlier, designing parallel algorithms is difficult.

Equivalently, profiling for graph data has to face the pressure of computing and memory when data profiling encounters large graphs. Therefore, it is necessary and worth researching to sample the graph data and carry out the tasks of data profiling based on the sample. But sampling graph data is more difficult than sampling relational data. Leskovec and Faloutsos [41] did practical experiments on sampling from large graphs. They concluded that best performing methods are the ones based on random-walks and "forest fire", with sample sizes as low as 15% of the original graph. However, how to apply these graph sampling methods to the above-mentioned graph data-based data profiling tasks is waiting for further experiments and exploration.

### C. Sampling for Profiling Heterogeneous Data

Data profiling is also widely used for heterogeneous data, e.g., discovering matching dependencies (MDs) [123], [124], reasoning about matching rules [125], [126], discovering a concise set of matching keys [127] and conditional matching dependencies (CMDs) [128]. However, these profiling tasks also have to face computational pressure in a big data context.

In fact, MDs, DDs and data dependencies are all based on differential functions. When calculating the measures for differential dependencies, performing sampling of pairwise comparison is more difficult. Given an instance of relation

R with N data tuples, pairwise comparison M will increase the total number to  $\frac{N*(N-1)}{2}$ , which will greatly increase the number of populations. However, many pairs in M are meaningless when calculating support for DDs [129], which means that the proportion of pairs we want is very small. Therefore, we must increase the sampling rate to expect to include these pairs in the sample, so as to get the approximate results as close as possible.

## REFERENCES

- [1] Min Chen, Shiwen Mao, and Yunhao Liu. Big data: A survey. *MONET*, 19(2):171–209, 2014.
- [2] James Manyika. Big data: The next frontier for innovation, competition, and productivity. [http://www.mckinsey.com/Insights/MGI/Research/Technology\\_and\\_Innovation](http://www.mckinsey.com/Insights/MGI/Research/Technology_and_Innovation), 2011.
- [3] J Anuradha et al. A brief introduction on big data 5vs characteristics and hadoop technology. *Procedia computer science*, 48:319–324, 2015.
- [4] Bin Chen, Peter J. Haas, and Peter Scheuermann. A new two-phase sampling based algorithm for discovering association rules. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, July 23-26, 2002, Edmonton, Alberta, Canada*, pages 462–468, 2002.
- [5] Albert Bifet. Mining big data in real time. *Informatica (Slovenia)*, 37(1):15–20, 2013.
- [6] Xindong Wu, Xingquan Zhu, Gong-Qing Wu, and Wei Ding. Data mining with big data. *IEEE Trans. Knowl. Data Eng.*, 26(1):97–107, 2014.
- [7] M. S. Mahmud, J. Z. Huang, S. Salloum, T. Z. Emar, and K. Sadat-diyov. A survey of data partitioning and sampling methods to support big data analysis. *Big Data Mining and Analytics*, 3(2):85–101, June 2020.
- [8] Sandra González-Bailón. Social science in the era of big data. *Policy & Internet*, 5(2):147–160, 2013.
- [9] Waleed Albattah. The role of sampling in big data analysis. In *Proceedings of the International Conference on Big Data and Advanced Wireless Technologies, BDAW 2016, Blagoevgrad, Bulgaria, November 10-11, 2016*, pages 28:1–28:5, 2016.
- [10] Ajay S Singh and Micah B Masuku. Sampling techniques & determination of sample size in applied statistics research: An overview. *International Journal of Economics, Commerce and Management*, 2(11):1–22, 2014.
- [11] Ashwin Satyanarayana. Intelligent sampling for big data using bootstrap sampling and chebyshev inequality. In *IEEE 27th Canadian Conference on Electrical and Computer Engineering, CCECE 2014, Toronto, ON, Canada, May 4-7, 2014*, pages 1–6, 2014.
- [12] George H. John and Pat Langley. Static versus dynamic sampling for data mining. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA*, pages 367–370, 1996.
- [13] Prashant Singh, Joachim Van Der Hertten, Dirk Deschrijver, Ivo Couckuyt, and Tom Dhaene. A sequential sampling strategy for adaptive classification of computationally expensive data. *Structural & Multidisciplinary Optimization*, 55(4):1–14, 2017.
- [14] Mohammed Javeed Zaki, Srinivasan Parthasarathy, Wei Li, and Mitsunori Ogihara. Evaluation of sampling for data mining of association rules. In *7th International Workshop on Research Issues in Data Engineering (RIDE '97) High Performance Database Management for Large-Scale Applications, Birmingham, UK, April 7-8, 1997*, pages 42–50, 1997.
- [15] Dorian Pyle. *Data Preparation for Data Mining*. Morgan Kaufmann, 1999.
- [16] Ziawasch Abedjan, Lukasz Golab, and Felix Naumann. Profiling relational data: a survey. *VLDB J.*, 24(4):557–581, 2015.
- [17] Shaoxu Song, Aoqian Zhang, Lei Chen, and Jianmin Wang. Enriching data imputation with extensive similarity neighbors. *PVLDB*, 8(11):1286–1297, 2015.
- [18] Shaoxu Song, Yu Sun, Aoqian Zhang, Lei Chen, and Jianmin Wang. Enriching data imputation under similarity rule constraints. *IEEE Trans. Knowl. Data Eng.*, 32(2):275–287, 2020.
- [19] Shaoxu Song, Han Zhu, and Jianmin Wang. Constraint-variance tolerant data repairing. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 877–892, 2016.
- [20] Felix Naumann. Data profiling revisited. *SIGMOD Record*, 42(4):40–49, 2013.
- [21] Ziawasch Abedjan, Lukasz Golab, and Felix Naumann. Data profiling. In *32nd IEEE International Conference on Data Engineering, ICDE 2016, Helsinki, Finland, May 16-20, 2016*, pages 1432–1435, 2016.
- [22] Shaoxu Song, Aoqian Zhang, Jianmin Wang, and Philip S. Yu. SCREEN: stream data cleaning under speed constraints. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 827–841, 2015.
- [23] Aoqian Zhang, Shaoxu Song, and Jianmin Wang. Sequential data cleaning: A statistical approach. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 909–924, 2016.
- [24] Fang Wang, Menggang Li, Yiduo Mei, and Wenrui Li. Time series data mining: A case study with big data analytics approach. *IEEE Access*, 8:14322–14328, 2020.
- [25] Wenfei Fan, Yinghui Wu, and Jingbo Xu. Functional dependencies for graphs. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 1843–1857, 2016.
- [26] Shaoxu Song, Hong Cheng, Jeffrey Xu Yu, and Lei Chen. Repairing vertex labels under neighborhood constraints. *PVLDB*, 7(11):987–998, 2014.
- [27] Shaoxu Song, Boge Liu, Hong Cheng, Jeffrey Xu Yu, and Lei Chen. Graph repairing under neighborhood constraints. *VLDB J.*, 26(5):611–635, 2017.
- [28] David Maier, Alon Y. Halevy, and Michael J. Franklin. Dataspaces: Co-existence with heterogeneity. In *Proceedings, Tenth International Conference on Principles of Knowledge Representation and Reasoning, Lake District of the United Kingdom, June 2-5, 2006*, page 3, 2006.
- [29] Shaoxu Song, Lei Chen, and Philip S. Yu. On data dependencies in dataspace. In *Proceedings of the 27th International Conference on Data Engineering, ICDE 2011, April 11-16, 2011, Hannover, Germany*, pages 470–481, 2011.
- [30] Shaoxu Song, Lei Chen, and Philip S. Yu. Comparable dependencies over heterogeneous data. *VLDB J.*, 22(2):253–274, 2013.
- [31] Ziawasch Abedjan, Lukasz Golab, and Felix Naumann. Data profiling: A tutorial. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14-19, 2017*, pages 1747–1751, 2017.
- [32] Hans T Schreuder, Timothy G Gregoire, and Johann P Weyer. For what applications can probability and non-probability sampling be used? *Environmental Monitoring and Assessment*, 66(3):281–291, 2001.
- [33] Ronaldo Iachan, Lewis Berman, Tonja M Kyle, Kelly J Martin, Yangyang Deng, Davia N Moysse, Deirdre Middleton, and Audie A Atienza. Weighting nonprobability and probability sample surveys in describing cancer catchment areas, 2019.
- [34] Konstantinos Slavakis, Georgios B. Giannakis, and Gonzalo Mateos. Modeling and optimization for big data analytics: (statistical) learning tools for our era of data deluge. *IEEE Signal Process. Mag.*, 31(5):18–31, 2014.
- [35] Rajeev Agrawal, Anirudh Kadadi, Xiangfeng Dai, and Frédéric Andrès. Challenges and opportunities with big data visualization. In *Proceedings of the 7th International Conference on Management of computational and collective intelligence in Digital EcoSystems, Caraguatatuba, Brazil, October 25 - 29, 2015*, pages 169–173, 2015.
- [36] Lina Zhou, Shimei Pan, Jianwu Wang, and Athanasios V. Vasilakos. Machine learning on big data: Opportunities and challenges. *Neuro-computing*, 237:350–361, 2017.
- [37] Cem Kadir and Hulya Cingi. Ratio estimators in simple random sampling. *Applied Mathematics and Computation*, 151(3):893–902, 2004.
- [38] Peter J Bickel and David A Freedman. Asymptotic normality and the bootstrap in stratified sampling. *The annals of statistics*, pages 470–482, 1984.
- [39] HJG Gundersen and EB Jensen. The efficiency of systematic sampling in stereology and its prediction. *Journal of microscopy*, 147(3):229–263, 1987.
- [40] Ralph H Henderson and Thalanayar Sundaresan. Cluster sampling to assess immunization coverage: a review of experience with a simplified sampling method. *Bulletin of the World Health Organization*, 60(2):253, 1982.
- [41] Jure Leskovec and Christos Faloutsos. Sampling from large graphs. In *Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadelphia, PA, USA, August 20-23, 2006*, pages 631–636, 2006.

- [42] Bee Wah Yap, Khatijahusna Abd Rani, Hezlin Aryani Abd Rahman, Simon Fong, Zuraida Khairudin, and Nik Nik Abdullah. An application of oversampling, undersampling, bagging and boosting in handling imbalanced datasets. In *Proceedings of the First International Conference on Advanced Data and Information Engineering, DaEng 2013, Kuala Lumpur, Malaysia, December 16-18, 2013*, pages 13–22, 2013.
- [43] Jeffrey Scott Vitter. Random sampling with a reservoir. *ACM Trans. Math. Softw.*, 11(1):37–57, 1985.
- [44] Qing He, Haocheng Wang, Fuzhen Zhuang, Tianfeng Shang, and Zhongzhi Shi. Parallel sampling from big data with uncertainty distribution. *Fuzzy Sets and Systems*, 258:117–133, 2015.
- [45] Yu-Lin He, Joshua Zhexue Huang, Hao Long, Qiang Wang, and Chenghao Wei. I-sampling: A new block-based sampling method for large-scale dataset. In *2017 IEEE International Congress on Big Data, BigData Congress 2017, Honolulu, HI, USA, June 25-30, 2017*, pages 360–367, 2017.
- [46] Zhicheng Liu, Biye Jiang, and Jeffrey Heer. *imMens*: Real-time visual querying of big data. *Comput. Graph. Forum*, 32(3):421–430, 2013.
- [47] Salman Salloum, Joshua Zhexue Huang, and Yu-Lin He. Empirical analysis of asymptotic ensemble learning for big data. In *Proceedings of the 3rd IEEE/ACM International Conference on Big Data Computing, Applications and Technologies, BDCAT 2016, Shanghai, China, December 6-9, 2016*, pages 8–17, 2016.
- [48] Qing He, Zhong-Zhi Shi, Li-An Ren, and ES Lee. A novel classification method based on hypersurface. *Mathematical and computer modelling*, 38(3-4):395–407, 2003.
- [49] Chun-Wei Tsai, Chin-Feng Lai, Han-Chieh Chao, and Athanasios V. Vasilakos. Big data analytics: a survey. *J. Big Data*, 2:21, 2015.
- [50] Foster J. Provost, David D. Jensen, and Tim Oates. Efficient progressive sampling. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Diego, CA, USA, August 15-18, 1999*, pages 23–32, 1999.
- [51] Tim Harford. Big data: A big mistake? *Significance*, 11(5):14–19, 2014.
- [52] Jonathan Nagler and Joshua A Tucker. Drawing inferences and testing theories with big data. *PS: Political Science & Politics*, 48(1):84–88, 2015.
- [53] Meir J Stampfer, Graham A Colditz, Walter C Willett, JoAnn E Manson, Bernard Rosner, Frank E Speizer, and Charles H Hennekens. Postmenopausal estrogen therapy and cardiovascular disease: ten-year follow-up from the nurses’ health study. *New England Journal of Medicine*, 325(11):756–762, 1991.
- [54] Robert M Kaplan, David A Chambers, and Russell E Glasgow. Big data and large sample size: a cautionary note on the potential for bias. *Clinical and translational science*, 7(4):342–346, 2014.
- [55] Jae Kwang Kim and Zhonglei Wang. Sampling techniques for big data analysis in finite population inference. *International Statistical Review*, 2018.
- [56] Siu-Ming Tam and Jae-Kwang Kim. Big data ethics and selection-bias: An official statistician’s perspective. *Statistical Journal of the IAOS*, 34(4):577–588, 2018.
- [57] Maisel and R. Inference from nonprobability samples. *Public Opinion Quarterly*, 36(3):407–407.
- [58] Xiao-Li Meng et al. Statistical paradises and paradoxes in big data (i): Law of large populations, big data paradox, and the 2016 us presidential election. *The Annals of Applied Statistics*, 12(2):685–726, 2018.
- [59] Peter J. Haas, Jeffrey F. Naughton, S. Seshadri, and Lynne Stokes. Sampling-based estimation of the number of distinct values of an attribute. In *VLDB’95, Proceedings of 21th International Conference on Very Large Data Bases, September 11-15, 1995, Zurich, Switzerland*, pages 311–322, 1995.
- [60] Moses Charikar, Surajit Chaudhuri, Rajeev Motwani, and Vivek R. Narasayya. Towards estimation error guarantees for distinct values. In *Proceedings of the Nineteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, May 15-17, 2000, Dallas, Texas, USA*, pages 268–279, 2000.
- [61] Phillip B. Gibbons. Distinct sampling for highly-accurate answers to distinct values queries and event reports. In *VLDB 2001, Proceedings of 27th International Conference on Very Large Data Bases, September 11-14, 2001, Roma, Italy*, pages 541–550, 2001.
- [62] Surajit Chaudhuri, Rajeev Motwani, and Vivek R. Narasayya. Random sampling for histogram construction: How much is enough? In *SIGMOD 1998, Proceedings ACM SIGMOD International Conference on Management of Data, June 2-4, 1998, Seattle, Washington, USA*, pages 436–447, 1998.
- [63] Phillip B. Gibbons, Yossi Matias, and Viswanath Poosala. Fast incremental maintenance of approximate histograms. *ACM Trans. Database Syst.*, 27(3):261–298, 2002.
- [64] Gurmeet Singh Manku, Sridhar Rajagopalan, and Bruce G. Lindsay. Random sampling techniques for space efficient online computation of order statistics of large datasets. In *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA*, pages 251–262, 1999.
- [65] Zengfeng Huang, Lu Wang, Ke Yi, and Yunhao Liu. Sampling based algorithms for quantile computation in sensor networks. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2011, Athens, Greece, June 12-16, 2011*, pages 745–756, 2011.
- [66] Xian Wu, Wei Fan, Jing Peng, Kun Zhang, and Yong Yu. Iterative sampling based frequent itemset mining for big data. *Int. J. Machine Learning & Cybernetics*, 6(6):875–882, 2015.
- [67] George Kollios, Dimitrios Gunopulos, Nick Koudas, and Stefan Berchtold. Efficient biased sampling for approximate clustering and outlier detection in large data sets. *IEEE Trans. Knowl. Data Eng.*, 15(5):1170–1187, 2003.
- [68] Ying Yan, Liang Jeff Chen, and Zheng Zhang. Error-bounded sampling for analytics on big sparse data. *PVLDB*, 7(13):1508–1519, 2014.
- [69] Julian Ramos Rojas, Mary Beth Kery, Stephanie Rosenthal, and Anind K. Dey. Sampling techniques to improve big data exploration. In *7th IEEE Symposium on Large Data Analysis and Visualization, LDVA 2017, Phoenix, AZ, USA, October 2, 2017*, pages 26–35, 2017.
- [70] HaiYing Wang, Min Yang, and John Stufken. Information-based optimal subdata selection for big data linear regression. *Journal of the American Statistical Association*, 114(525):393–405, 2019.
- [71] Sunghae Jun, Seung-Joo Lee, and Jea-Bok Ryu. A divided regression analysis for big data. *International Journal of Software Engineering and Its Applications*, 9(5):21–32, 2015.
- [72] Yannis Sismanis, Paul Brown, Peter J. Haas, and Berthold Reinwald. GORDIAN: efficient and scalable discovery of composite keys. In *Proceedings of the 32nd International Conference on Very Large Data Bases, Seoul, Korea, September 12-15, 2006*, pages 691–702, 2006.
- [73] Ziawasch Abedjan and Felix Naumann. Advancing the discovery of unique column combinations. In *Proceedings of the 20th ACM Conference on Information and Knowledge Management, CIKM 2011, Glasgow, United Kingdom, October 24-28, 2011*, pages 1565–1570, 2011.
- [74] Tobias Bleifuß, Susanne Bülow, Johannes Frohnhofen, Julian Risch, Georg Wiese, Sebastian Kruse, Thorsten Papenbrock, and Felix Naumann. Approximate discovery of functional dependencies for large datasets. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016*, pages 1803–1812, 2016.
- [75] Thorsten Papenbrock and Felix Naumann. A hybrid approach to functional dependency discovery. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 821–833, 2016.
- [76] Ihab F. Ilyas, Volker Markl, Peter J. Haas, Paul Brown, and Ashraf Aboulnaga. CORDS: automatic discovery of correlations and soft functional dependencies. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, Paris, France, June 13-18, 2004*, pages 647–658, 2004.
- [77] M. Li, H. Wang, and J. Li. Mining conditional functional dependency rules on big data. *Big Data Mining and Analytics*, 3(1):68–84, March 2020.
- [78] Sebastian Kruse, Thorsten Papenbrock, Christian Dullweber, Moritz Finke, Manuel Hegner, Martin Zabel, Christian Zöllner, and Felix Naumann. Fast approximate discovery of inclusion dependencies. In *Datenbanksysteme für Business, Technologie und Web (BTW 2017), 17. Fachtagung des GI-Fachbereichs „Datenbanken und Informationssysteme“ (DBIS), 6.-10. März 2017, Stuttgart, Germany, Proceedings*, pages 207–226, 2017.
- [79] Ruihong Huang, Zhiwei Chen, Zhicheng Liu, Shaoxu Song, and Jianmin Wang. Tsoutlier: Explaining outliers with uniform profiles over iot data. In *2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, December 9-12, 2019*, pages 2024–2027. IEEE, 2019.
- [80] Ziawasch Abedjan, Lukasz Golab, Felix Naumann, and Thorsten Papenbrock. *Data Profiling*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2018.
- [81] Shaoxu Song, Han Zhu, and Lei Chen. Probabilistic correlation-based similarity measure on text records. *Inf. Sci.*, 289:8–24, 2014.

- [82] Hazar Harmouch and Felix Naumann. Cardinality estimation: An experimental survey. *PVLDB*, 11(4):499–512, 2017.
- [83] Robert Kooi. *The Optimization of Queries in Relational Databases*. PhD thesis, Case Western Reserve University, September 1980.
- [84] Stephen Macke, Yiming Zhang, Silu Huang, and Aditya G. Parameswaran. Adaptive sampling for rapidly matching histograms. *Proc. VLDB Endow.*, 11(10):1262–1275, 2018.
- [85] Patricia G. Selinger, Morton M. Astrahan, Donald D. Chamberlin, Raymond A. Lorie, and Thomas G. Price. Access path selection in a relational database management system. In *Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data, Boston, Massachusetts, USA, May 30 - June 1*, pages 23–34, 1979.
- [86] Ira Pohl. *A minimum storage algorithm for computing the median*. IBM TJ Watson Research Center, 1969.
- [87] Lili Zhang, Wenjie Wang, and Yuqing Zhang. Privacy preserving association rule mining: Taxonomy, techniques, and metrics. *IEEE Access*, 7:45032–45047, 2019.
- [88] Shaoxu Song, Chunping Li, and Xiaoquan Zhang. Turn waste into wealth: On simultaneous clustering and cleaning over dirty data. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pages 1115–1124, 2015.
- [89] Aojian Zhang, Shaoxu Song, Yu Sun, and Jianmin Wang. Learning individual models for imputation. In *35th IEEE International Conference on Data Engineering, ICDE 2019, Macao, China, April 8-11, 2019*, pages 160–171. IEEE, 2019.
- [90] Rakesh Agrawal and John C. Shafer. Parallel mining of association rules. *IEEE Trans. Knowl. Data Eng.*, 8(6):962–969, 1996.
- [91] David Wai-Lok Cheung, Jiawei Han, Vincent T. Y. Ng, Ada Wai-Chee Fu, and Yongjian Fu. A fast distributed algorithm for mining association rules. In *Proceedings of the Fourth International Conference on Parallel and Distributed Information Systems, December 18-20, 1996, Miami Beach, Florida, USA*, pages 31–42, 1996.
- [92] Ali Seyed Shirkhorshidi, Saeed Reza Aghabozorgi, Ying Wah Teh, and Tutut Herawan. Big data clustering: A review. In *Computational Science and Its Applications - ICCSA 2014 - 14th International Conference, Guimarães, Portugal, June 30 - July 3, 2014, Proceedings, Part V*, pages 707–720, 2014.
- [93] Publishing Models/article Dates Explained. Sampling: Design and analysis. *Technometrics*, 42(2):223–224, 2009.
- [94] Roger Sauter. Introduction to probability and statistics for engineers and scientists. *Technometrics*, 47(3):378, 2005.
- [95] Shaoxu Song and Lei Chen. Differential dependencies: Reasoning and discovery. *ACM Trans. Database Syst.*, 36(3):16:1–16:41, 2011.
- [96] Thorsten Papenbrock and Felix Naumann. Data-driven schema normalization. In *Proceedings of the 20th International Conference on Extending Database Technology, EDBT 2017, Venice, Italy, March 21-24, 2017*, pages 342–353, 2017.
- [97] Xiang Lian, Lei Chen, and Shaoxu Song. Consistent query answers in inconsistent probabilistic databases. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2010, Indianapolis, Indiana, USA, June 6-10, 2010*, pages 303–314, 2010.
- [98] Stéphane Lopes, Jean-Marc Petit, and Farouk Toumani. Discovering interesting inclusion dependencies: application to logical database tuning. *Inf. Syst.*, 27(1):1–19, 2002.
- [99] Shaoxu Song, Lei Chen, and Hong Cheng. Parameter-free determination of distance thresholds for metric distance constraints. In *IEEE 28th International Conference on Data Engineering (ICDE 2012), Washington, DC, USA (Arlington, Virginia), 1-5 April, 2012*, pages 846–857, 2012.
- [100] Shaoxu Song, Lei Chen, and Hong Cheng. Efficient determination of distance thresholds for differential dependencies. *IEEE Trans. Knowl. Data Eng.*, 26(9):2179–2192, 2014.
- [101] Arvid Heise, Jorge-Arnulfo Quiané-Ruiz, Ziawasch Abedjan, Anja Jentzsch, and Felix Naumann. Scalable discovery of unique column combinations. *PVLDB*, 7(4):301–312, 2013.
- [102] Jyrki Kivinen and Heikki Mannila. The power of sampling in knowledge discovery. In *Proceedings of the Thirteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, May 24-26, 1994, Minneapolis, Minnesota, USA*, pages 77–85, 1994.
- [103] Jyrki Kivinen and Heikki Mannila. Approximate dependency inference from relations. In *Database Theory - ICDT'92, 4th International Conference, Berlin, Germany, October 14-16, 1992, Proceedings*, pages 86–98, 1992.
- [104] Jixue Liu, Jiuyong Li, Chengfei Liu, and Yongfeng Chen. Discover dependencies from data - A review. *IEEE Trans. Knowl. Data Eng.*, 24(2):251–264, 2012.
- [105] Philip Bohannon, Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. Conditional functional dependencies for data cleaning. In *Proceedings of the 23rd International Conference on Data Engineering, ICDE 2007, The Marmara Hotel, Istanbul, Turkey, April 15-20, 2007*, pages 746–755, 2007.
- [106] Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *ACM Trans. Database Syst.*, 33(2):6:1–6:48, 2008.
- [107] Wenfei Fan, Floris Geerts, Jianzhong Li, and Ming Xiong. Discovering conditional functional dependencies. *IEEE Trans. Knowl. Data Eng.*, 23(5):683–698, 2011.
- [108] Fabien De Marchi, Stéphane Lopes, and Jean-Marc Petit. Efficient algorithms for mining inclusion dependencies. In *Advances in Database Technology - EDBT 2002, 8th International Conference on Extending Database Technology, Prague, Czech Republic, March 25-27, Proceedings*, pages 464–476, 2002.
- [109] Fabien De Marchi, Stéphane Lopes, and Jean-Marc Petit. Unary and n-ary inclusion dependency discovery in relational databases. *J. Intell. Inf. Syst.*, 32(1):53–73, 2009.
- [110] Thorsten Papenbrock, Sebastian Kruse, Jorge-Arnulfo Quiané-Ruiz, and Felix Naumann. Divide & conquer-based inclusion dependency discovery. *PVLDB*, 8(7):774–785, 2015.
- [111] Xiaochen Zhu, Shaoxu Song, Xiang Lian, Jianmin Wang, and Lei Zou. Matching heterogeneous event data. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014*, pages 1211–1222, 2014.
- [112] Yu Gao, Shaoxu Song, Xiaochen Zhu, Jianmin Wang, Xiang Lian, and Lei Zou. Matching heterogeneous event data. *IEEE Trans. Knowl. Data Eng.*, 30(11):2157–2170, 2018.
- [113] Xiaochen Zhu, Shaoxu Song, Jianmin Wang, Philip S. Yu, and Jianguang Sun. Matching heterogeneous events with patterns. In *IEEE 30th International Conference on Data Engineering, Chicago, ICDE 2014, IL, USA, March 31 - April 4, 2014*, pages 376–387, 2014.
- [114] Shaoxu Song, Yu Gao, Chaokun Wang, Xiaochen Zhu, Jianmin Wang, and Philip S. Yu. Matching heterogeneous events with patterns. *IEEE Trans. Knowl. Data Eng.*, 29(8):1695–1708, 2017.
- [115] Aojian Zhang, Shaoxu Song, Jianmin Wang, and Philip S. Yu. Time series data cleaning: From anomaly detection to anomaly repairing. *PVLDB*, 10(10):1046–1057, 2017.
- [116] Shaoxu Song, Yue Cao, and Jianmin Wang. Cleaning timestamps with temporal constraints. *PVLDB*, 9(10):708–719, 2016.
- [117] Lukasz Golab, Howard J. Karloff, Flip Korn, Avishek Saha, and Divesh Srivastava. Sequential dependencies. *PVLDB*, 2(1):574–585, 2009.
- [118] Jianmin Wang, Shaoxu Song, Xiaochen Zhu, and Xuemin Lin. Efficient recovery of missing events. *PVLDB*, 6(10):841–852, 2013.
- [119] Jianmin Wang, Shaoxu Song, Xiaochen Zhu, Xuemin Lin, and Jianguang Sun. Efficient recovery of missing events. *IEEE Trans. Knowl. Data Eng.*, 28(11):2943–2957, 2016.
- [120] Jianmin Wang, Shaoxu Song, Xuemin Lin, Xiaochen Zhu, and Jian Pei. Cleaning structured event logs: A graph repair approach. In *31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015*, pages 30–41, 2015.
- [121] Wenfei Fan, Zhe Fan, Chao Tian, and Xin Luna Dong. Keys for graphs. *PVLDB*, 8(12):1590–1601, 2015.
- [122] Wenfei Fan, Chunming Hu, Xueli Liu, and Ping Lu. Discovering graph functional dependencies. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, pages 427–439. ACM, 2018.
- [123] Shaoxu Song and Lei Chen. Discovering matching dependencies. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009, Hong Kong, China, November 2-6, 2009*, pages 1421–1424, 2009.
- [124] Shaoxu Song and Lei Chen. Efficient discovery of similarity constraints for matching dependencies. *Data Knowl. Eng.*, 87:146–166, 2013.
- [125] Wenfei Fan, Xibei Jia, Jianzhong Li, and Shuai Ma. Reasoning about record matching rules. *PVLDB*, 2(1):407–418, 2009.
- [126] Wenfei Fan, Hong Gao, Xibei Jia, Jianzhong Li, and Shuai Ma. Dynamic constraints for record matching. *VLDB J.*, 20(4):495–520, 2011.
- [127] Shaoxu Song, Lei Chen, and Hong Cheng. On concise set of relative candidate keys. *PVLDB*, 7(12):1179–1190, 2014.
- [128] Yihan Wang, Shaoxu Song, Lei Chen, Jeffrey Xu Yu, and Hong Cheng. Discovering conditional matching rules. *TKDD*, 11(4):46:1–46:38, 2017.

- [129] Selasi Kwashie, Jixue Liu, Jiuyong Li, and Feiyue Ye. Efficient discovery of differential dependencies through association rules mining. In *Databases Theory and Applications - 26th Australasian Database Conference, ADC 2015, Melbourne, VIC, Australia, June 4-7, 2015. Proceedings*, pages 3–15, 2015.