

A systematic review on privacy-preserving distributed data mining

Chang Sun ^{a,*}, Lianne Ippel ^b, Andre Dekker ^c, Michel Dumontier ^d and Johan van Soest ^e

^a *Institute of Data Science, Faculty of Science and Engineering, Maastricht University, Maastricht, The Netherlands*

ORCID: <https://orcid.org/0000-0001-8325-8848>

^b *Institute of Data Science, Faculty of Science and Engineering, Maastricht University, Maastricht, The Netherlands*

ORCID: <https://orcid.org/0000-0001-8314-0305>

^c *Department of Radiation Oncology (MAASTRO), GROW School for Oncology and Developmental Biology, Maastricht University Medical Centre, Maastricht, The Netherlands*

ORCID: <https://orcid.org/0000-0002-0422-7996>

^d *Institute of Data Science, Faculty of Science and Engineering, Maastricht University, Maastricht, The Netherlands*

ORCID: <https://orcid.org/0000-0003-4727-9435>

^e *Department of Radiation Oncology (MAASTRO), GROW School for Oncology and Developmental Biology, Maastricht University Medical Centre, Maastricht, The Netherlands and Brightlands Institute of Smart Society (BISS), Faculty of Science and Engineering, Maastricht University, Maastricht/Heerlen, The Netherlands*

ORCID: <https://orcid.org/0000-0003-2548-0330>

Editor: Karin Verspoor (<https://orcid.org/0000-0002-8661-1544>)

Solicited reviews: Abdur Rahim (<https://orcid.org/0000-0003-0237-170>); Dayana Spagnuolo (<https://orcid.org/0000-0001-6882-6480>)

Received 19 March 2021

Accepted 16 August 2021

Abstract. Combining and analysing sensitive data from multiple sources offers considerable potential for knowledge discovery. However, there are a number of issues that pose problems for such analyses, including technical barriers, privacy restrictions, security concerns, and trust issues. Privacy-preserving distributed data mining techniques (PPDDM) aim to overcome these challenges by extracting knowledge from partitioned data while minimizing the release of sensitive information. This paper reports the results and findings of a systematic review of PPDDM techniques from 231 scientific articles published in the past 20 years. We summarize the state of the art, compare the problems they address, and identify the outstanding challenges in the field. This review identifies the consequence of the lack of standard criteria to evaluate new PPDDM methods and proposes comprehensive evaluation criteria with 10 key factors. We discuss the ambiguous definitions of privacy and confusion between privacy and security in the field, and provide suggestions of how to make a clear and applicable privacy description for new PPDDM techniques. The findings from our review enhance the understanding of the challenges of applying theoretical PPDDM methods to real-life use cases, and the importance of involving legal-ethical and social experts in implementing

*Corresponding author. E-mail: chang.sun@maastrichtuniversity.nl.

PPDDM methods. This comprehensive review will serve as a helpful guide to past research and future opportunities in the area of PPDDM.

Keywords: Survey, data mining, privacy preserving, distributed learning

1. Introduction

Mining distributed, sensitive data offers tantalising potential for new insights and a wide variety of applications, but is generally fraught with concerns of model accuracy and data privacy. Consider the case of analyzing patient data in the healthcare domain: hospitals have used patient data to improve diagnostic accuracy and efficiency [29,31] and to fuel the transition to preventive [17] and precision medicine [6,27,95]. However, learning patient data from a single hospital might cause limited model performance and incomplete knowledge discovery [59]. Patients' health are not only affected by genetic and biological factors, but also by individual behaviour and social circumstances [19]. Combining various patient data from multiple sources offers one pathway to obtain more accurate and reliable analytical models for health outcomes [3,97]. However, combining distributed sensitive data faces a number of challenges including: data protection compliance to one or more legal jurisdictions, privacy concerns, security, and trust issues. Beyond the healthcare domain, this also applies to applications in many other fields, such as finance and law [82,114]. Conventional centralised data mining techniques are challenged in this environment and require viable alternatives.

Privacy-preserving distributed data mining (PPDDM), which focuses on the analysis of decentralised data without leaking sensitive information from any party to the other parties, offers one way forward for multiple data parties to overcome the challenges posed by centralising the data for analysis [72]. PPDDM techniques, whether data mining or machine learning, aim to make it technically or mathematically infeasible to deduce the original data from a communication message, and certainly from the final analysis result. To make use of PPDDM in practical applications, we should consider the data problems (e.g., classification, regression), the adversarial concerns the involving data parties have (e.g., malicious, honest), and the balance between data privacy and model performance. PPDDM is sometimes referred to privacy-preserving federated learning after Google first proposed the concept in 2016 [66,76]. However, privacy-preserving federated learning can be regarded as a specific category of PPDDM, in which there is a federation of autonomous organisations that express an interest to contribute to a joint analysis [92].

A number of PPDDM methods have been reported in the last 20 years. The existing survey papers have compared the theoretical backgrounds, strengths, and limitations. However, the analysis of distributed data has been poorly addressed as only one special case of privacy-preserving data mining [1,9,89,110]. The distributed data problem has been addressed to a limited extent in the survey of Hina Vaghashia [102] and Suchitra Shelke [91]. Vassilios S. et al. [110] presented five dimensions of state-of-the-art privacy-preserving data mining algorithms where the problem of analysing distributed data was merely considered to be addressed by cryptography-based techniques and only the association rule mining problem and decision tree induction were presented in this survey. Several surveys summarized the evaluation parameters to assess privacy-preserving techniques including privacy level, hiding failure, data quality, complexity, efficiency, and resistance of different data mining algorithms [9,10,36,110]. Others have a major focus on the definition and construction of Secure Multiparty Computation (SMC) and how SMC can be combined with data mining algorithms [18,72,103]. In a recent survey [77], privacy-preserving

approaches were summarized for data collection, data publishing, data mining output, and distributed learning. The majority of the published surveys have typically treated PPDDM as a specialised subtopic of either distributed data mining or privacy-preserving data mining. As an emerging field, PPDDM is under-reported in the existing surveys and now requires a more comprehensive and complete analysis.

Accordingly, the main aim of this systematic review is to provide an overview of existing approaches and identify outstanding challenges in the field of PPDDM. This paper reports the results and findings of a comprehensive review of PPDDM techniques from 231 scientific articles published in the past 20 years. We present the characteristics of the 18 most cited studies and analyze their influence on other studies in the field. The results show a wide range of privacy-preserving methods and data mining algorithms have been well-studied. We highlight the findings showing a lack of standard evaluation criteria in the field, the ambiguous definition of privacy, and insufficient experimental information in some studies. These findings enhance the understanding of the challenges of applying the theoretical PPDDM methods to real-life use cases, and the importance of involving legal-ethical and social experts in implementing PPDDM methods.

The main contributions of this work to the literature in the PPDDM field are:

- (1) to propose comprehensive criteria with 10 key factors to evaluate the new PPDDM techniques. The evaluation criteria include adversarial behaviour of data parties, data partitioning, experiment datasets, privacy/security analysis, privacy-preserving methods, data mining problems, analysis algorithms, complexity and cost, performance measures, and scalability.
- (2) to present different definitions of privacy, distinguish information privacy from information security in the PPDDM field, and provide suggestions of how to make clear and applicable privacy descriptions to propose new PPDDM techniques.
- (3) to identify the most cited PPDDM articles, analyze their characteristics and how these articles influence other studies in the field, and
- (4) to provide a guideline based on the proposed evaluation criteria for researchers to conduct future research and publications in the PPDDM field.

This systematic review offers new insights into the important factors that should be considered to propose and evaluate new PPDDM techniques and how to bridge the gap between theoretical methods and practical applications in the field. We present this review paper as a helpful guide to past research and future opportunities in the area of PPDDM.

The outline of this paper is as follows. In the next section, we present existing privacy-preserving methods and define terms related to PPDDM. In Section 3, we describe the approach in conducting this systematic review. In Section 4, we provide the results of our review, including evaluation criteria. In Section 5, we compare the key influential papers. In the last section, we summarize our main findings, present a list of recommendations, and discuss future directions.

2. Privacy-preserving methods

Privacy-preserving methods, as the major component of PPDDM techniques, are used to minimize the release of information during data mining model training and communication among multiple parties. Various privacy-preserving methods have been proposed from different communities such as statistics, cryptography, data mining, and secure data transfer. In this section, we summarize the most commonly-used privacy-preserving methods in PPDDM.

2.1. Secure multiparty computation (SMC)

Secure multiparty computation protocols are designed for multiple parties to jointly compute some function over their own data without revealing the original data to any other parties [72]. The foundation for SMC started from cryptography. In addition to protect the participants from being attacked by external parties (who are outside of the system or protocol), SMC also protects the participants from each other. For example, some SMC protocols are implemented to prevent participants from learning private information from other parties or deliberately sending incorrect computation results to other parties. The following sub-sections describe some well-known protocols in SMC.

2.1.1. Building blocks (primitives) SMC of protocols

Secure protocols that are deployed as building blocks of secure computation are used to prevent data being revealed or deduced from the communication and/or computation between data parties [72]. Commonly used encryption protocols include oblivious transfer and homomorphic encryption. Oblivious transfer, first developed by Even et al. [33], considers two data parties, a requester and a sender, where the requester obtains exactly one instance without the sender knowing which element was queried, and without the requester knowing about the other instances that were not retrieved. Oblivious transfer protocols iteratively pass over the data many times during training, and as a result are computationally expensive. Another technique, homomorphic encryption, was introduced by Rivest [86]. This technique supports certain algebraic operations such as additions and multiplications on encrypted text (i.e., ciphertext). The decrypted result from the operations on ciphertext matches the result of the operations performed on the plain text. Homomorphic encryption systems are grouped into fully homomorphic encryption (FHE) or partial homomorphic encryption (PHE) [81]. As the initial scheme of a homomorphic cryptosystem, PHE can only perform a specific algebra operation such as addition or multiplication in each iteration. This limits the usability for data mining algorithms, as the algorithms consist of several complex operations. On the contrary, FHE supports any desirable operation and functionality that can run on the ciphertext. Since the ciphertext is never decrypted, the input from each data party is not revealed. The first generation of FHE system was proposed by Gentry in 2009 [42]. However, FHE systems are not sufficiently efficient due to the high computational cost of performing iterative operations over encrypted data during the training epochs.

2.1.2. Generic SMC protocols

Generic SMC protocols were implemented for any probabilistic polynomial-time function [72]. Unlike homomorphic encryption systems, these generic protocols are sensitive to the number of data parties. The commonly-used protocol of secure two-party computation is Yao's garbled circuit protocol [124]. The protocol is based on evaluating the function that needs to be computed by two data parties as a combinatorial circuit with a collection of gates (e.g., AND, XOR gate). These gates connect with circuit-input wires, circuit-output wires and intermediate wires. Each gate has two input wires and one single output wire. The required communication of the protocol depends on the size of the circuit, while the computation cost depends on the number of input wires. Extensions to more than two data parties, i.e. the cases of multiparty computation, have been developed by Micali et al. [79], Beaver et al. [5], and Ben-Or et al. [7]. Following Yao's theory, these protocols are based on designing the function as a circuit and applying a secure computation protocol to the circuit [72]. Beside computational complexity, communication cost is a considerable factor in these protocols. All protocols need a one-to-one communication channel between every pair of parties. Some require a broadcast channel for all parties.

2.1.3. Specialized SMC protocols

Specialized SMC protocols are commonly used as primitives to the data mining algorithms including secure sum, secure set union, secure size of intersection, and secure scalar product protocols. These protocols allow certain operations without revealing any inputs from any of the participating data parties.

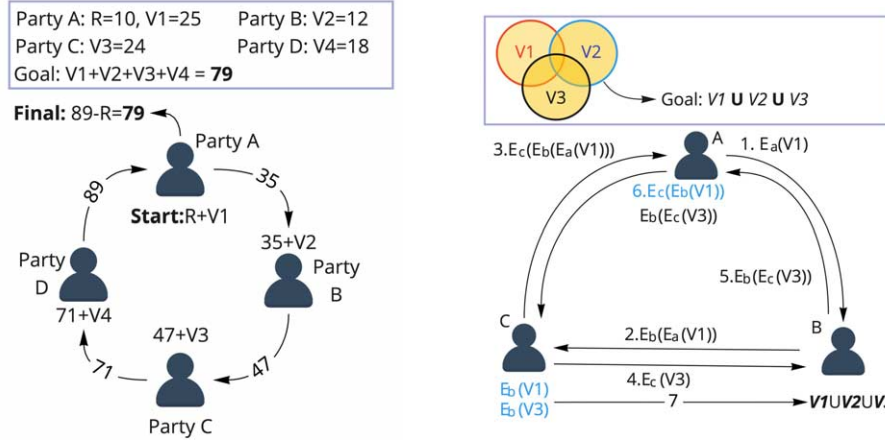
Secure sum as a basic and simple example of secure multiparty computation was introduced by Clifton et al. to obtain the sum of the inputs [18]. The protocol is as follows: data party A has $V1$ local value. Party A generates a random number R and calculates $(R + V1)$ and sends this result to data party B (PB). Then, Party B adds their local value to the received value and sends it $(R + V1 + V2)$ to the next party. In the end, to obtain the final result, the last sum value will be sent back to party A to subtract R . The protocol ends with sending this final result to all participating parties. An example of securely computing a sum among 4 four parties is shown in Fig. 1(a).

Secure set union has been applied to the case where data parties want to jointly create unions of sets from rules and itemsets shared by multiple parties but not leaking the owner of each set. To guarantee a secure computation, one approach is to apply a commutative encryption system in computing the set union [18,85]. A commutative encryption system can encrypt original data multiple times using different users' public keys. The final encrypted data can be decrypted without considering the order of the public keys in the encryption process [51]. In the secure set union protocol, one data party encrypts its own itemsets using commutative encryption and transfers them to other parties. The receiver party encrypts both its own sets and the received encrypted sets and passes it to the next party. Once the data is encrypted by all parties, decryption can start at each party in any order. The permutation of the encryption order prevents the participating parties from tracking the ownership of itemsets. However, if one item is present at multiple data parties, then the number of the item will be exposed because of duplication. Figure 1(b) presents an example of securely computing a set union among three data parties.

Secure size of set intersection is solving the problem that multiple data parties want to obtain the size of set intersection of their local datasets without revealing the ownership. Similar to secure set union, each data party encrypts its own item sets by using commutative encryption and sends it to another data party. The receiver encrypts these items, arbitrarily permutes the order, and sends it to the next data party. This process ends when all item sets are encrypted by all data parties. Due to the commutative encryption, if and only if the original inputs are the same, then the final outcomes from two different item sets can be equal. Therefore, the number of values that occur in all encrypted item sets is the size of the set intersection. No input will get exposed since only encryption (no decryption) is required. Figure 1(c) demonstrates the protocol of securely computing the size of set intersection.

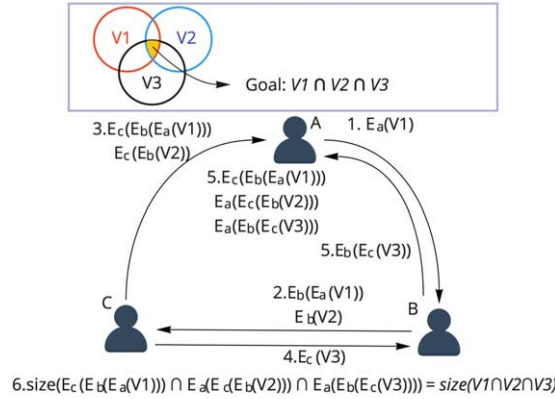
Secure scalar product protocols are essential and powerful. It has been widely applied in many data mining algorithms which can be decomposed to the calculation of scalar products. As a notable example, Vaidya and Clifton extended a secure scalar products protocol to solve association rule mining problems between two parties [104]. The general idea is as follows:

- (1) Data party A has $X = \{x_1, \dots, x_n\}$, while data party B has $Y = \{y_1, \dots, y_n\}$. The goal is to calculate $X * Y = \sum_{i=1}^n (x_i * y_i)$ without revealing inputs to the other party. Both parties share a matrix C which is generated by random numbers.
- (2) The protocol starts at Party A who generates n random numbers $Ra = \{r_1, \dots, r_n\}$. Then, party A calculates $X' = X + C * Ra$ and send to party B.
- (3) Party B generates $m (< n)$ random numbers Rb and calculate $Y' = C_1 * Y + Rb_1, \dots, C_{n/m} * Y + Rb_1, \dots, C_{2n/m} * Y + Rb_2, \dots, C_n * Y + Rb_n$ and $S' = \sum_{i=1}^n (x'_i * y_i)$. Y' and S' are sent to party A.



(a) An example of secure computation of a sum among four parties. R is a random number generated by Party A. $V1$, $V2$, $V3$, and $V4$ presents the private data from party A to party D.

(b) An example of secure computation of a set union among three parties. Party A, B, C encrypts their private data using a commutative encryption scheme respectively (E_a , E_b , E_c). Text in blue is decrypted text.



(c) An example of secure computation of a size of set intersection among three parties. Party A, B, C encrypts their private data using a commutative encryption scheme respectively (E_a , E_b , E_c).

Fig. 1. Examples of three secure multiparty computation protocols - Fig. 1(a). Secure sum protocol, Fig. 1(b). Secure set union protocol, Fig. 1(c). Secure size of set intersection protocol.

- (4) Party A calculates $S'' = S' - \sum_{i=1}^n (Ra * Y')$ and m sets of sum of Ra which is $Ra' = Ra_1 + Ra_2 + \dots + Ra_{n/m} + Ra_{n/m+1} + \dots + Ra_{2n/m}, \dots, Ra_{((m-1)n/m)+1} + Ra_{((m-1)n/m)+2} + \dots + Ra_n$. Party A sends S'' and Ra' for final result calculation.
- (5) Party B computes the final scalar product as $S = S'' + Ra' * Rb$.

The security of this secure scalar product protocol is guaranteed by the inability of either side to deduce k equations with more than k unknowns. As with many other existing scalar product protocols [4,

54], it is limited to the collaboration between only two parties because of the lack of efficiency in practice [18].

2.2. Data perturbation

Data Perturbation preserves data privacy by adding ‘noise’ to the individual records but still keeps the key summary information about the data [116]. One major approach of data perturbation is to use statistical techniques to replace the original data with synthetic values which have the same or comparable statistical information (e.g., distributions) as the original values. The synthetic data can be generated by a statistical model which learns from the original data. The other main approach is to distort the values by applying additive noise, multiplicative noise, or other randomization procedures [2]. Data swapping, another method of data perturbation, switches a set of (sensitive) attributes between different data entities to prevent the linkage of records to identities [23,35]. The major drawback of these methods is the decrease of data quality and accuracy of the learning model. Data perturbation techniques are more commonly used to protect privacy in data publishing problems [77].

2.3. Local learning and global integration

The method that integrates local models to one global model uses the foundation of ensemble learning that trains a set of models in order to enhance the performance of one single model [84,112]. Each data party can train their own local data miners independently. Then, these local data miners are sequentially or parallelly integrated to compose a center or global data miner which can generate the final results. Consequently, the original data of each party is never transferred to other data parties. A majority of data mining algorithms have been theoretically developed to this approach including Support Vector Machine [40,74,100,108], Decision Tree [34,98,106], Neural Networks [21,28,100,128] and so forth. A few of them have been successfully implemented, applied and evaluated in practical use cases such as [59] and [118].

3. Methodology

This paper follows the systematic review procedures described by Kitchenham [65]. In this section, we will detail the workflow. First, we discuss the inclusion and exclusion criteria of study selection, followed by the search strategies, and evaluation criteria for reviewing selected studies.

3.1. Eligibility criteria

We selected papers that are peer-reviewed publications in English between 2000 and 2020 (August) working on data mining and machine learning techniques that solve problems of classification, regression, clustering, or association rule mining. The eligible papers must take privacy preservation into account when data mining and machine learning models are executed on partitioned data. Partitioned data includes horizontally partitioned/homogeneous data, vertically partitioned/heterogeneous data, and arbitrarily partitioned data (The definitions of different partitioned data are presented in Section 3.3). Furthermore, included papers must 1) propose and/or implement a new approach and/or; 2) apply existing approaches to a practical case and/or; 3) improve the performance of existing approaches.

To narrow down the number of publications, we excluded poster and workshop abstracts, survey papers, and articles that only contain discussions on current concerns and future research directions. To set the scope of this survey, the authors screened titles, keywords, and abstracts to exclude the papers that 1) only focus on privacy-preserving data mining/machine learning on centralised data, 2) solve problems of parallel computing, cloud computing, grid computing, edge computing, and fog computing to improve computational performance rather than the complexity of the data analysis problem, 3) solve privacy issues in data collecting, data publishing, data storage, and data querying, and 4) focus on Blockchain, web attacks detection, intrusion detection, data privacy focusing on mobile devices, geographic data privacy, and differential privacy. If the papers could not be identified based on its title, keywords, and abstract, the authors reviewed the full text of the paper.

3.2. Search strategy

According to the eligibility criteria above, we used the following search engines and digital libraries: IEEE Xplore Digital Library,¹ ACM Digital Library,² Science Direct,³ ISI Web of Science,⁴ Springer Link,⁵ PubMed.⁶ Based on the inclusion criteria, we formulated the following terms to search in the title, abstract, and keywords of papers. The entire workflow for selecting relevant studies is presented with search results in Fig. 3 in Section 4.1.

- (1) *privacy* and (*distributed* or *de-centralized* or *de-centralised* or *partitioned*) and “*machine learning*” (**PPDML**)
- (2) *privacy* and (*distributed* or *de-centralized* or *de-centralised* or *partitioned*) and “*data mining*” (**PPDDM**)
- (3) *privacy* and (*vertically* or *heterogeneous*) and “*machine learning*” (**PPVML**)
- (4) *privacy* and (*vertically* or *heterogeneous*) and “*data mining*” (**PPVDM**)
- (5) *privacy* and (*horizontally* or *homogeneous*) and “*machine learning*” (**PPHML**)
- (6) *privacy* and (*horizontally* or *homogeneous*) and “*data mining*” (**PPHDM**)

3.3. Evaluation criteria for reviewing papers

To evaluate the paper on PPDDM techniques, conventional data mining evaluation criteria are not adequate [84]. Beside conventional evaluation methods, additional factors such as communication costs, data partitioning, adversary behavior, privacy measures should be considered. To the best of our knowledge, there are no standard criteria for evaluating new PPDDM approaches. Consequently, studies selected a various set of evaluation methods which they think are necessary for their approaches. In this review, we assessed selected papers considering the following 10 factors including adversarial behavior of data party, data partitioning, experimented datasets, privacy/security analysis, privacy-preserving methods, data mining problems, analysis algorithms, complexity and cost, performance measures, and scalability. The authors initially generated and modified these evaluation criteria by reviewing 10% of

¹ IEEE Xplore: <https://ieeexplore.ieee.org/Xplore/home.jsp/>.

² ACM Digital Library: <https://dl.acm.org/>.

³ ScienceDirect: <https://www.sciencedirect.com/>.

⁴ Web of Science – Clarivate: <https://clarivate.com/products/web-of-science/>.

⁵ Springer Link: <https://link.springer.com/>.

⁶ PubMed: <https://www.ncbi.nlm.nih.gov/pubmed/>.

the included articles. Then, the evaluation criteria have been discussed by the co-authors in several iterations of reviewing until an agreement has been made on these 10-factor evaluation criteria. Afterwards, all selected papers have been reviewed and assessed again using the criteria.

1) Adversarial behavior of data parties covers the assumed adversarial behavior that involved data parties have. In this review, we consider two types of adversarial behavior of involved parties – semi-honest and malicious. A semi-honest (also called passive, or honest-but-curious) party follows the protocol properly, however is also curious about other parties' data [72]. The semi-honest party will attempt to learn or deduce data from other parties. A malicious (or active) party will arbitrarily deviate from the protocol and will make deliberate attacks to obtain access to data from other parties [44]. For example, possible malicious behavior might be not starting the execution of protocols at all or suspending (or aborting) the execution at any desired point in time. Papers that use ambiguous expressions such as 'untrusted' or 'non-trusting' or 'non-collaborative' are not classified into any category, because they did not clearly indicate the adversarial property of data parties, nor did they provide any privacy or security proof of their methods. In addition, we include the situation where a third party was involved. A third party, as another independent entity, can combine data from multiple parties, execute analysis on the joint datasets, or do the final computation based on information from data parties. A third party can be fully-honest, semi-honest, and malicious.

2) Data partitioning Fig. 2 shows three scenarios of data partitioning which are considered in this review: 1) Horizontally partitioned data which contains the same attributes from different data instances (see Fig. 2(a)). For example, different hospitals see different patients, though they collect the same patient attributes; 2) Vertically partitioned data which contains the same data instances but with different attributes (see Fig. 2(b)). For example, a hospital has data on the same individuals as the tax office, while the attributes collected differs per data party; 3) Arbitrarily partitioned data, the hybrid situation of horizontally and vertically partitioned data. In this scenario, the data providing institutes hold different attributes for different data instances (see Fig. 2(c)).

3) Dataset information factor indicates whether the study provides adequate information about the applied datasets in their experiments. Basic information of datasets including sources, names, numbers of features and instances, categorical or numeric type (if available) were recorded. Considering the readability, collected information is composed into five categories for this factor:

- (1) Datasets that are publicly available (e.g., UCI repository) [101]
- (2) Datasets from practical cases such as real patients data from a clinic
- (3) Synthetic datasets and datasets which were generated by authors
- (4) Experiments are presented in the paper but information about datasets is missing
- (5) No experiments are presented in the paper

4) Privacy definition or measurement describes whether the study gave an explicit privacy definition, analyses, or measurements. Due to a lack of a universally accepted standard definition, there are many different definitions of privacy from various aspects such as law and philosophical point of view covering personal information, body, communications, and territory [24,57]. This review only focuses on information privacy which concerns the control of collection, use, retention, and distribution of personal information. During reviewing, we do not assess if the privacy definitions are correct and the levels of privacy these studies can preserve though whether they gave a sufficient description, measurement, or analysis of privacy.

5) Privacy-preserving methods are classified into 5 categories: 1) secure multiparty computation – building blocks, 2) secure multiparty computation – generic and specialized construction protocols, 3)

ID	Age	Gender	Education	Wellbeing	Diabetes
1	56	Male	University	Good	YES
2	25	Female	University	Medium	NO
3	31	Female	High School	Good	NO
4	45	Male	Primary School	Poor	YES
5	32	Male	Primary School	Good	No
6	60	Female	High School	Poor	YES
7	55	Male	University	Medium	NO

(a) An example of horizontally partitioned data.

ID	Age	Gender	Education	Wellbeing	Diabetes
1	56	Male	University	Good	YES
2	25	Female	University	Medium	NO
3	31	Female	High School	Good	NO
4	45	Male	Primary School	Poor	YES
5	32	Male	Primary School	Good	No
6	60	Female	High School	Poor	YES
7	55	Male	University	Medium	NO

(b) An example of vertically partitioned data.

ID	Age	Gender	Education	Wellbeing	Diabetes
1	56	Male	University	Good	YES
2	25	Female	University	Medium	NO
3	31	Female	High School	Good	NO
4	45	Male	Primary School	Poor	YES
5	32	Male	Primary School	Good	No
6	60	Female	High School	Poor	YES
7	55	Male	University	Medium	NO

(c) An example of arbitrarily partitioned data.

Fig. 2. Examples of three different partitioned data. Figure 2(a) shows horizontally partitioned data which contains the same attributes/features from different data instances. Figure 2(b) shows vertically partitioned data which contains the same data instances but with different attributes/features. Figure 2(c) shows arbitrarily partitioned data which is a hybrid situation of horizontally and vertically partitioned data.

data modification, 4) local learning and global integration, and 5) others. First 4 categories have been explained in detail in the Privacy-Preserving Method Section. The papers which did not use any method from above are categorized to “others”.

6) Types of problems covers four main data mining areas: i.e., classification, regression, clustering, and association rule mining. Classification predicts a class with categorical labels. These categorical labels can be represented by discrete values, where the ordering among values has no meaning. In contrast, regression is to predict continuous-valued function or ordered value. Clustering is to group a set

of data objects into multiple groups (clusters) so that objects within a cluster have high similarity, but are very dissimilar to objects in other clusters. Association rule mining is to discover interesting associations and correlations between itemsets in transactional and relational databases [48]. Additionally, we labeled the studies as “general” that solved some mathematical or statistical problems which are applied to classification, regression, and clustering. The studies which worked on outlier detection, record linkage, recommendation system, attribute/dimension reduction, feature selection, and probabilistic graph are categorized into “others”.

7) Data mining algorithms present the algorithms which have been developed in a privacy-preserving manner and which ones lack attention. There are plenty of algorithms across the data mining and statistics domain [9,38]. In this review, the top eight algorithms are listed in the result table including decision tree, K-nearest neighbor, bayesian networks, support vector machine, neural networks, K-means, linear/logistic regressions, and A-priori algorithms.

8) Complexity and cost indicates whether the study explicitly measures computational complexity, time cost, and communication cost. The papers which did not present any experiments but only briefly discussed computation, time, and communication costs are counted as “No Measurement”.

9) Performance measures covers whether the study compared the performance of their approaches with 1) other published PPDDM methods, 2) centralised data mining methods, and 3) distributed without preserving privacy methods. The performance measures include accuracy, precision, recall, F1 score, AUC (Area Under the Curve), mean squared error, mean absolute error, and other standard evaluation criteria in the data mining domain [13,39,48,49,83]. Owing to the high degree of heterogeneity in the reporting of performance measures across the reviewed papers, we determine whether any performance measure was applied to evaluate the methods rather than comparing different performance measures. The papers which contained experiments but did not compare their results with other methods are categorized into “No comparison (with experiment)”. The studies which did not provide any experiments are classified to “No experiments”.

10) Scalability covers whether the study presented a scalability analysis or the experiments prove the scalability of their approach. The scalability in this review means if the approach can tackle large-size datasets which contain a large number of either features or instances. It is noteworthy that only discussing scalability or mentioning their approaches are scalable were not included.

4. Results

In this section, we first describe the number and distribution of search results retrieved from the six search engines in the last 20 years. Detailed reviews of selected papers based on the evaluation criteria are elaborated in Section 4.2. The analysis of the relations among selected papers is described in Section 4.3.

4.1. Search results

In Fig. 3, we present the workflow of this systematic review with the number of papers included in each step. Following the inclusion criteria, 4222 publications including duplicates were retrieved from six search engines. Most papers were from IEEE and Springer Link followed by ACM Digital Library. To remove the duplicates, we used Digital Object Identifiers (DOI) to keep the unique papers. The number of publications was reduced from 4222 to 2424. Furthermore, we filtered out irrelevant papers by screening the titles and abstracts of the retrieved papers. Papers that focused on parallel computing, cloud computing, edge computing, network security, intrusion detection, web attack detection, privacy

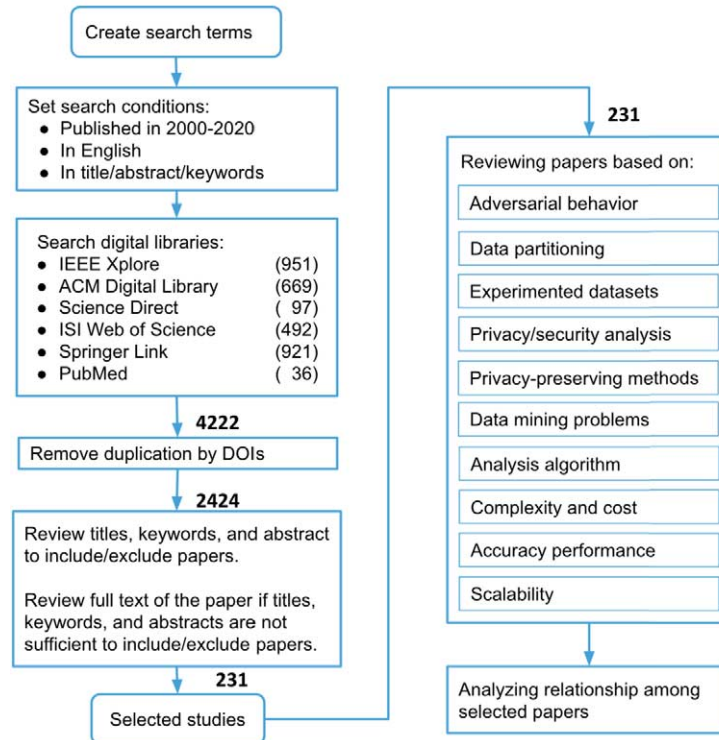


Fig. 3. Workflow of conducting this systematic review.

in mobile data and geographic data, differential privacy, privacy in data collecting, data publishing, data storing, data querying were excluded. In the end, 231 papers were selected to be preliminarily reviewed.

To improve the insight of the search result, we map the selected papers into graphs by using the Gephi visualization tool [43]. In Fig. 4, the distribution of 231 selected papers using different search terms is presented. Papers are presented as nodes and clustered by the search terms. For instance, 182 selected papers were found by using the search term – PPDDM, while 38 of them were findable in PPHDM category and 50 of them were findable in PPVDM. It is obvious that data mining papers are the majority of the search outcomes. It is reasonable as data mining covers a larger scope than machine learning. Privacy issues should be considered in the entire data processing procedure instead of only the part of analysis and building machine learning models. Moreover, a large number of papers (71 papers from PPDDM, 22 papers from PPDML) did not indicated what exact data partitioning problems (vertical, horizontal, or arbitrary) their method can solve in their titles, abstracts, and keywords. This increases difficulties for other researchers and practitioners to find the correct papers based on their needs.

4.2. Review results

In Fig. 5, we summarize the review results of 231 papers using the 10 evaluation factors we discussed previously. The full review results of 231 papers are publicly available in the data repository: <https://doi.org/10.6084/m9.figshare.14239937.v4>. (DOI: 10.6084/m9.figshare.14239937). The following subsection elaborates on the review result of each factor.

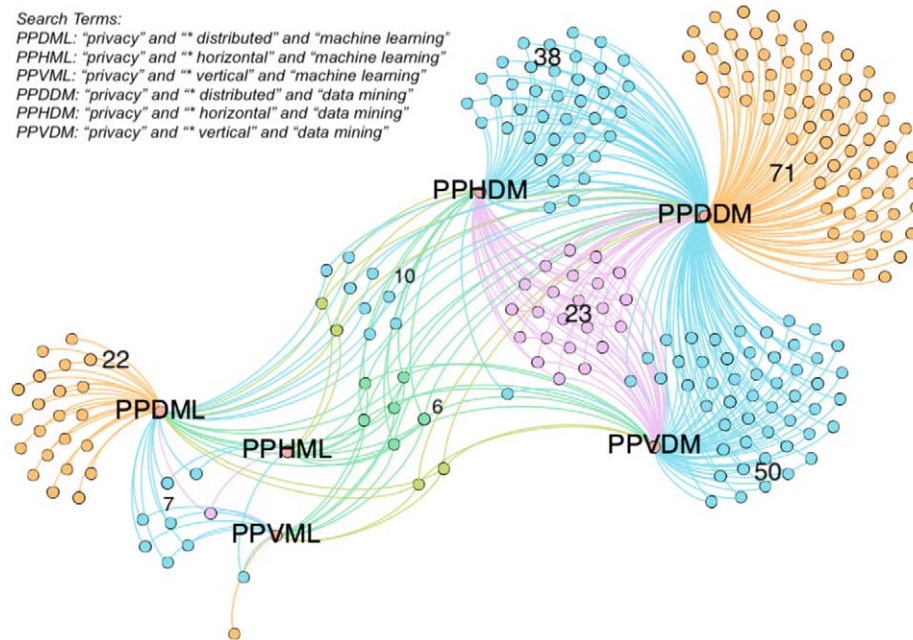


Fig. 4. Numbers and clusters of selected papers from different search terms. Papers are presented as nodes and clustered by the search terms. The number of papers in each cluster is labeled in the figure. The edges show which search terms were used to find the papers. For example, the 23 nodes in the purple cluster were found from using search terms PPDDM, PPHDM, and PPVDM.

Adversarial behavior of data parties. About half of the reviewed studies assuming their approaches are applicable for the data parties with semi-honest adversary behavior. In contrast, only 17 reviewed studies developed their methods against malicious parties. Third party constructions were applied in the method of 47 studies. More than half of them handled semi-honest behavior data parties together with employing the third party. However, it is worth noting that over 30% of selected papers did not state a clear assumption that which adversarial behavior their approach can deal with.

Data partitioning. Horizontally partitioned data (105 reviewed papers) and vertically partitioned data (112 reviewed papers) seem to be represented equally in the selected literature. There are 35 papers handling both horizontally partitioned data and vertically partitioned data. However, only 9 reviewed studies developed PPDDM methods on arbitrarily partitioned data which can work with semi-honest data parties. Additionally, 20% of selected studies did not indicate in which data partitioning situation their methods can be applied.

Privacy is one of the most important evaluation parameters for PPDDM techniques. However, only one fifth of selected studies describe an explicit definition of privacy and mathematical analysis of how much information is leaked by the proposed method. There are 81 papers proving the security of their approaches rather than a privacy analysis. The difference between security and privacy will be discussed in the next section. The majority of studies describe "privacy preservation" very briefly in their own understanding. These descriptions are heterogeneous: e.g., "not revealing privacy of any database", "not compromising the privacy of the data owners", "preserving the confidentiality of datasets", and "no important information leakage". The remaining 30 papers proposed new PPDDM methods without indicating any definition or description about privacy.

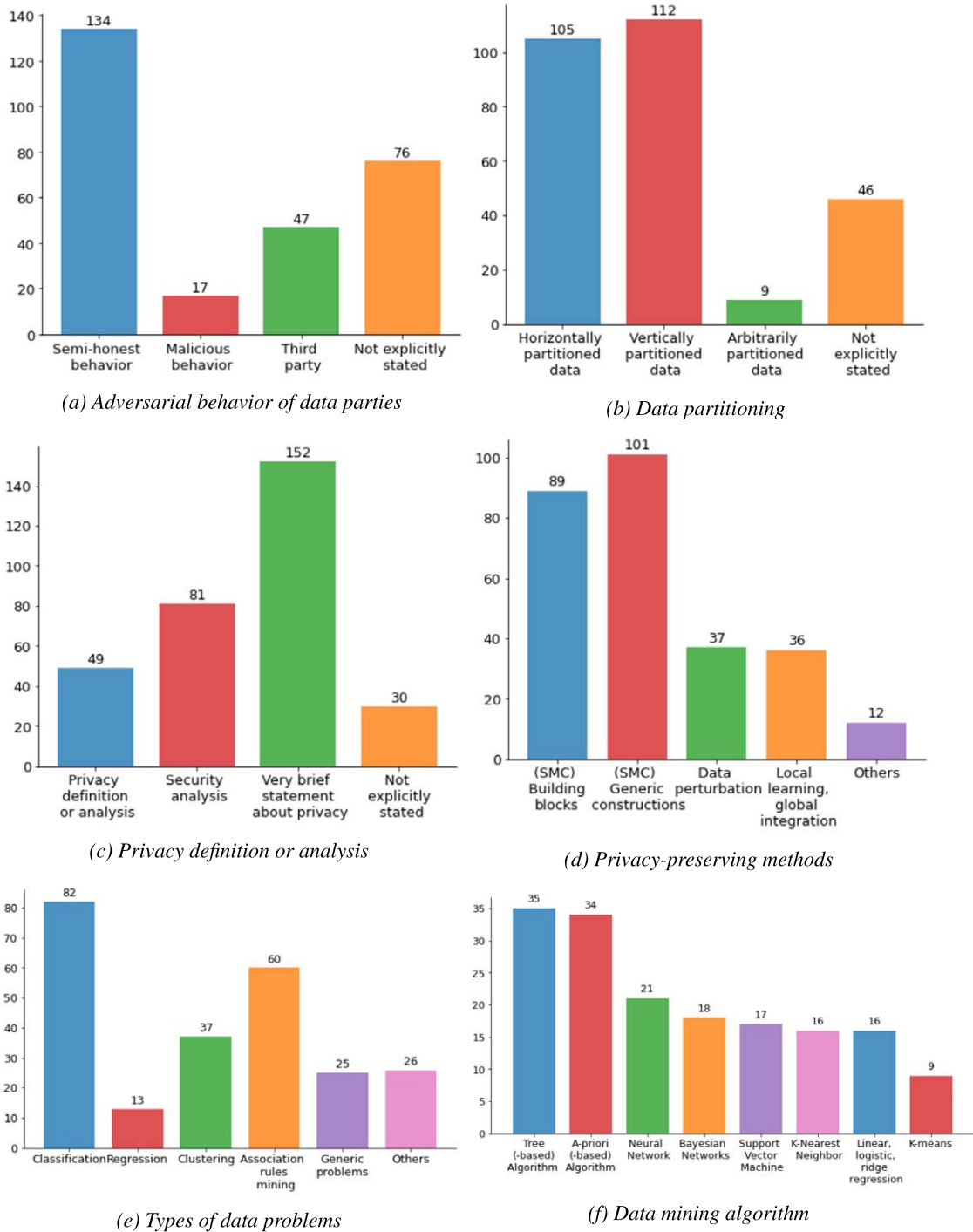


Fig. 5. Bar charts of presenting review results using 10-factor evaluation criteria. Papers can cover one or more items in the factors except privacy definition/analysis and scalability.

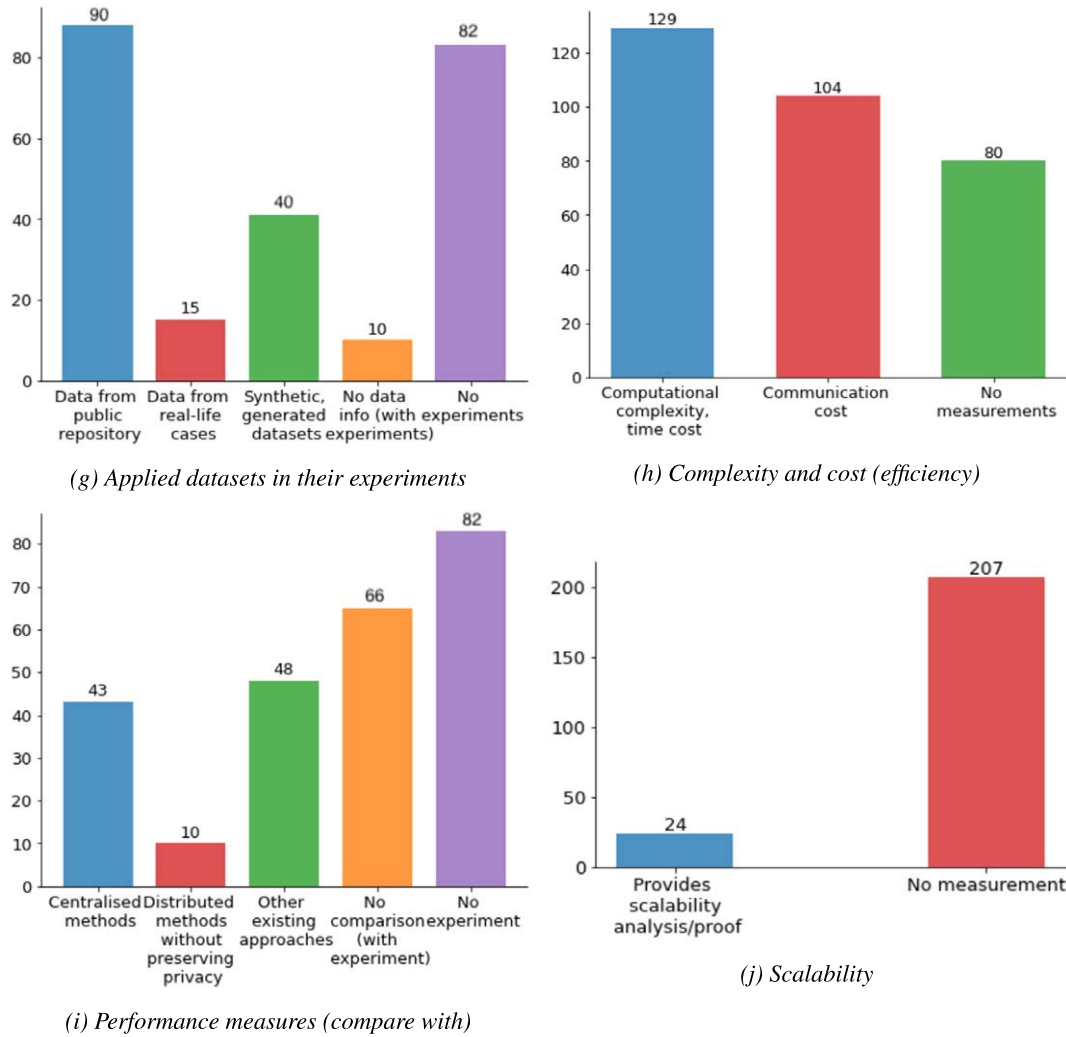


Fig. 5. (Continued.)

Privacy-preserving methods. Secure multiparty computation techniques are the most encountered solutions in the PPDDM domain. The generic and specialized protocols were applied in 101 papers, while 89 studies employed homomorphic encryption or oblivious transfer protocols. A minority of reviewed studies used data modification, or methodologies to train local models and combine these local models into a global model. A combination of techniques such as combining data modification and homomorphic encryption protocols has been applied by 41 studies.

Types of data problems and data mining algorithms. Classification problems attracted the most attention from researchers in the PPDDM domain, followed by association rule mining and clustering. By contrast, a minority of studies deal with regression modeling. The most implemented data mining algorithms tackling these data problems are: Tree-based algorithms such as decision tree, random forest (35 papers), A-priori-based algorithms (34 papers), Neural Networks (21), Bayesian Networks (18), Support Vector Machine (17), K-Nearest Neighbor (16), Linear/Logistic/Ridge Regression (16), and K-

means (9). There are over 10% of reviewed papers studied on generic algorithms that can be applied to multiple data mining techniques such as gradient descent. About 12% of reviewed papers worked on solving privacy problems in outlier detection, record linkage, recommendation system approaches, attribute/dimension reduction, feature selection, and probabilistic graphs.

Applied datasets in their experiments. From the selected studies, we identified the datasets that were applied in their experiments, measurement of complexity and cost, and performance on accuracy and scalability. We found 90 studies used datasets from public repositories, while 40 studies generated synthetic datasets to conduct their experiments. It is noteworthy that only 15 papers applied real-world datasets in practical use cases. Furthermore, it is remarkable to find that 82 papers proposed new methods by only presenting mathematical theories without any experiments, while 10 papers conducted experiments but did not provide any information about the datasets.

Complexity and cost. To prove the efficiency of proposed methods, 129 papers calculated computational complexity and/or time cost, while 104 papers reported communication cost of their approaches. Among them, 85 papers measured both computational complexity/time cost and communication cost. However, one third of (80) reviewed papers did not have any measurement of computation, running time, or communication cost.

Accuracy performance. We found 82 reviewed papers were lacking in evaluating accuracy performance of their methods because no experiments were conducted in these studies. In the rest of the papers, 43 papers proved their PPDDM methods can achieve comparable accuracy as the centralised data mining methods, while 48 studies proved their methods exceeded other existing PPDDM methods or achieved the same accuracy with higher efficiency. A small proportion of (10) studies proved their privacy-preserving models have comparable performance on learning partitioned data as the non-privacy-preserving models. Lastly, 66 papers conducted experiments but did not compare with any other methods or situations.

Scalability. The last factor – scalability – shows 10% papers proved or analyzed the scalability of their proposed methods. The majority of papers either only provided very brief statements in the discussion and future work section of the paper, or did not consider the scalability challenge.

4.3. Result of referencing relationship among selected papers

We investigated how selected papers influence each other based on their references and citations. We extracted text from reference sections of all selected studies and recognized titles and authors from the text. As DOIs are not available in the reference section of all papers, only titles and authors were used to recognize different studies. Figure 6 illustrates the citation network, where papers are represented as nodes, and citing relations are represented as edges. The size of nodes are proportional to the number of citations among the 231 papers. Papers [80,104,105] are most cited, with 1354, 1320, and 875 citations respectively (until 2021 Feb).

Table 1 lists the attributes of the most cited articles. Semi-honest behavior is the most common assumption, while none of these influential papers addressed malicious adversarial behavior. 3 out of 18 studies considered a third party. Two papers [55,108] took all possible data distribution situations (horizontally, vertically, and arbitrarily partitioned data) into account. Horizontally and vertically partitioned data problems have been covered with a good balance. Although the vertically partitioned data problem is more complicated than the horizontally one [107,113], our review indicates that they have been developed at the same pace.

A similar balance is apparent in the types of problems as well. Seven papers focused on solving a classification problem by using SVM, decision tree, bayesian networks, while 8 papers looked at

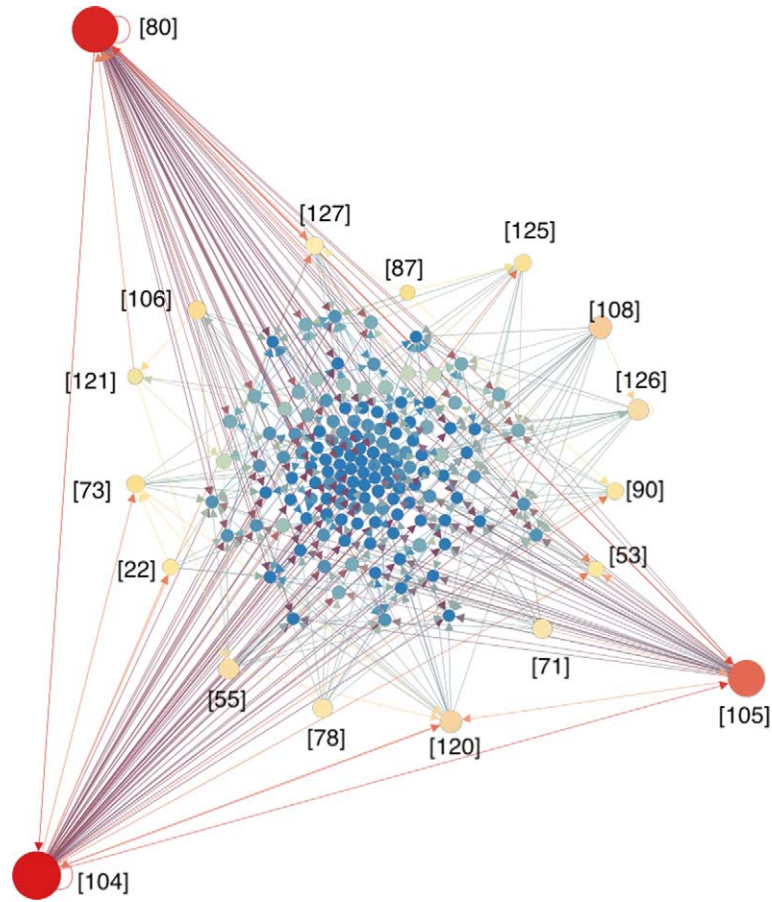


Fig. 6. Citation network among the selected papers. Papers are presented as nodes, while the citing relations are presented as edges. The size of nodes are proportional to the number of citations among the 231 papers.

clustering problems particularly at K-means, Expectation Maximization algorithms (EM), Local Outlier Factor (LOF) algorithm. Association rule mining problem has fewer influential papers, but the top 2 influential papers [80,104] both focused on this problem. In contrast to the balance in the types of problems, privacy-preserving solutions from the influential papers are completely dominated by SMC. 16 out of 18 influential papers covered SMC [121,127] combined SMC with homomorphic encryption, while [108,125,126] combined it with structuring local and global data miners. More than half of existing studies in our review applied SMC as the major privacy-preserving method.

It is notable that 12 out of 18 studies did not conduct experiments, but they provided explicit privacy/security analyses and costs measurements instead. These privacy/security analyses have been presented in different ways, but the main objectives were similar. All influential papers described what information their approaches can protect, what information have to be disclosed, and what potential risks, problems or troubles might exist. Moreover, their computational complexity and communication costs of their approaches were clearly presented as one of the evaluation parameters. Hence, the described performance evaluation on privacy and efficiency may be the reasons why these papers are often cited.

Table 1

Review results for the 18 most cited papers in this review. (PP method: privacy-preserving methods; local-global: local learning and global integration; ARM: association rule mining)

Ref	User scenario		Data distribution			Privacy/ security analysis	PP method*		Type of problems			Experiment	Cost	
	Semihonest	Third party	Horizontal	Vertical	Arbitrary		SMC	Local global*	Classification	Clustering	ARM*		Computation	Communication
[104]	✓			✓		✓	✓				✓			✓
[80]	✓		✓			✓	✓				✓		✓	✓
[105]	✓			✓		✓	✓			✓			✓	✓
[108]	✓	✓	✓	✓	✓	✓	✓	✓	✓			✓	✓	✓
[120]	✓			✓			✓		✓				✓	✓
[126]				✓		✓	✓	✓	✓					
[55]	✓	✓	✓	✓	✓	✓	✓			✓			✓	✓
[71]	✓		✓			✓	✓			✓				
[78]			✓			✓		✓		✓		✓		✓
[106]	✓			✓		✓		✓	✓				✓	✓
[127]				✓		✓	✓			✓			✓	✓
[73]	✓		✓	✓		✓	data perturbation		✓	✓		✓	✓	✓
[125]	✓		✓				✓	✓	✓			✓	✓	✓
[90]	✓		✓			✓	✓			✓				
[53]	✓	✓	✓			✓	✓			✓		✓	✓	✓
[22]				✓		✓	✓		Probabilistic graph			✓		✓
[121]	✓		✓				✓		✓				✓	✓
[87]	✓			✓		✓	✓		✓		✓		✓	✓

5. Discussion

PPDDM has been rapidly developing through active research programs across different scientific communities including data mining and machine learning, mathematics and statistics, cryptography, and data management. The total number of publications in this domain has dramatically increased in the last 20 years. Many of the studies included promising results in the efficiency and accuracy of their models in an experimental environment. These promising experimental results helped move the field forward towards practical applications. In the past five years, use cases have been developed in healthcare [25,59,63,69], finance [15], and technology companies [11,50,66] to examine different PPDDM methods. Participation of industry partners accelerates the transformation of PPDDM theoretical methods to practical applications. The existing PPDDM methods have been well-developed to solve a wide range of data problems (e.g., classification, clustering, association rule mining) using various data mining algorithms. To achieve the goal of PPDDM methods in practical studies, methods that will preserve privacy require legal, ethical, and social scholars in addition to scientific and technical experts. Successful implementation of PPDDM needs a joint effort from researchers with diverse backgrounds.

5.1. Inadequate definition and measurement of privacy

There are some challenges hindering PPDDM methods to be further developed and widely applied in practice. One of the key issues is the lack of the definition and measurement of (information) privacy. The meaning and operational definition of privacy is commonly ambiguous and subjective in the selected papers. It is not sufficiently expressed by the papers what privacy means to them, and what their proposed approaches can preserve. The three most common definitions of privacy preservation in the selected papers are 1) not revealing sensitive information; 2) not revealing private information; 3) not revealing raw data. However, it is unclear if “sensitive information” or “private information” or “raw data” is equal to personal information privacy. To understand personal information privacy from a legal and ethical perspective, it is the right of an individual or group to seclude themselves, or information about themselves, and thereby express themselves selectively [8,20,93]. Similarly, privacy is seen as the claim of individuals, groups, or institutions to determine for themselves when, how, and to what extent information about them is communicated to others [115]. In relation to controlling and protecting privacy, two definitions from legal literature state “Privacy, as a whole or in part, represents the control of transactions between person(s) and other(s), the ultimate aim of which is to enhance autonomy and/or to minimize vulnerability” [75] and “Privacy is to protect personal data and information related to a communication entity to be collected from other entities that are not authorized” [26].

According to privacy definitions above, any information about a person can be considered as privacy regardless of its sensitivity, originality, and transformation. It is the data subject that determines what data is private. For instance, a data subject might consider their state of mental health more private than their date of birth. However, existing PPDDM methods have not yet addressed different privacy requirements from each data subject. All data elements have equal treatment for all data subjects. This might cause insufficient privacy preservation for some data elements and data subjects, while over-protection for the others. To personalize the privacy preservation, Xiao and Tao [122] proposed a new generalization framework using personalized anonymity that data subjects can specify the degree of privacy protection for her/his data elements. In the study, Xiao and Tao [122] assume: 1) data subjects can easily set/change their privacy requirements with data parties, 2) data subjects are knowledgeable about the benefits and consequences of setting different degrees of privacy. This method is only applicable when

the data is centralized. In the partitioned data scenario, there is no platform yet facilitating data subjects to customize privacy requirements for each data element across multiple parties. Second, privacy requirements can be satisfied when using one single data source. However, analyzing an amount of partitioned data from multiple sources increases risk of privacy violation. As indicated by the 2020 European Commission White Paper on Artificial Intelligence [32], data about persons can be re-identified through the analysis of large amounts of other non-private data.

5.2. Ambiguity between privacy and security

Another ambiguity lies in the difference between (information) privacy and (information) security. Different from privacy, security has an explicit definition and measurement from the cryptography domain, separating the problem into semantic security and technical security [44]. Semantic security is a computational-complexity analogue of Shannon's definition of perfect privacy (which requires that the ciphertext yield no information regarding the plaintext). Technical security is the infeasibility of distinguishing between encryptions of a given pair of messages. Generally speaking, security focuses on maximally protecting information/data from malicious attacks and stealing data. Satisfying security requirements is not always sufficient for addressing privacy issues [56]. However, in the majority of the reviewed papers, the difference between security and privacy is not clearly stated. For example, some studies defined the data privacy but evaluated the methods by conducting security analysis [46,58,70]. Certain approaches guarantee that the data used for the analyses remain unknown to other parties through secure computation. However, this does not mean that the resulting output from the analyses is equally privacy-preserving [56,60,72]. The output can reveal information about the person so that the privacy is still not preserved according to the privacy definition we discussed above. For instance, the outcome of the analysis might portray a harmful profile for individuals sharing certain characteristics. Some essential problems are not taken into consideration, such as how much data or information will be revealed by the output although the output is computed securely [69], whether the models and algorithms are harmless to the data party or individuals, does the purpose of formula or function satisfy the legal and ethical concerns [96,99]. A typical example is building a decision tree on vertically partitioned data in a privacy-preserving way. The decision tree model can be securely and correctly built up. However, to some extent, the decision tree, as an output, leaks information about the input data [37]. Decision tree algorithm splits nodes based on attributes or features, while the splitting decision is dictated by the data. When the final decision tree is completed, the leaf nodes in the tree might reveal some information about the input data such as class counts. Therefore, releasing the final decision tree to all participating parties could potentially breach privacy.

Providing an applicable privacy description is significant to any PPDDM studies. What data or information should be preserved from mining can be influenced by different legal restrictions, ethical concerns, organizational regulations, personal preference, and application domains. Instead of generalizing the solution of a specific scheme to all situations, it is more reasonable to make a precise statement on the specific scenario to address. Therefore, the authors could provide a clear description to readers about what privacy means to them, and in which situation the proposed approach is privacy preserving by answering the following questions:

- (1) *What is the operational definition of privacy-preservation for the work?*
- (2) *Which data are deemed sensitive or require protection, and why?*
- (3) *What computational operation is intended to preserve privacy, and where does it fail?*

- (4) *What is the role or responsibility of each actor (e.g., data collector, data holder, data publisher, data analyst) in the scenario?*

5.3. *Inadequate experiments and practical use cases*

Our review result shows half of the reviewed papers did not provide any experiments to evaluate their methods, and as such there were no reports of accuracy, efficiency, and scalability in these papers. This is probably one of the gaps between the theoretical research and practical use cases in this domain. Solutions based on theory might not solve real world problems. In our review, only a few papers applied real-world use cases to evaluate their methods. It reflects a fact in this domain that many solutions have been proposed by researchers, but only a few of them were implemented in practice. Without experimenting on real data, the proposed approaches might neglect essential problems such as sparse or biased datasets [21,52], or record linkage problems in vertically partitioned data [61,94,109]. Future research in PPDDM should consider conducting experiments using real-world datasets and provide adequate information about the experiments. Meanwhile, we observed most real-life use cases to examine existing PPDDM approaches from the healthcare domain [25,69,99]. We suggest researchers apply the PPDDM methods to practical cases also in other research domains such as social science and finance. In addition to developing new theories, implementing and improving existing approaches in practice can also make a meaningful contribution to the PPDDM domain.

Nevertheless, these findings were observed in the light of limitations in our search strategy, which are elaborated in Section 5.6. This review did not specifically search for follow-up studies of reviewed papers. A possible effect is that papers which lack experiments might present their experiments in the follow-up studies, and might introduce selection bias towards the low number of practical experiments. However, we would argue that our search strategy would have found these papers if proper terminology was used.

5.4. *Challenge of linking data in vertically partitioned data scenario*

The accurate linking of entities across distributed datasets is of crucial importance in vertically partitioned data mining. Data parties must link their data and/or order them in an identical manner prior to data analysis. However, most papers assume this correspondence between data entities (records) exist by default. Matching data entities from multiple datasets can be error-prone particularly where the use of direct identifiers – even encrypted – are prevented by law, as is the case in the use of the national Citizen Service Number (“Burgerservicenummer”) in the Netherlands [12]. Sharing such identifiers compromises privacy as the sole information that a data subject is known to another data entity might be sensitive. Furthermore, one often assumes that records can be linked by doing exact matching on this unique identifier. However, exact matching can be very difficult due to the unstable and incorrect identifiers. Winkler and Schnell showed that 25% true matches would have been missed by exact matching in a census operation [88,117]. In another case, two data parties do not share the unique identifiers but have some features in common. As an alternative solution, two parties can match the data entities based on their common features. The matching accuracy will be affected by the correctness, completeness, and updating promptness of these common features from both data parties. In addition, privacy needs to be preserved in the matching procedure. Some efficient and privacy-compliant algorithms for the field of privacy-preserving entity matching have been developed [16,41,47,111] in the past 10 years.

5.5. A recommendation list of key parameters for PPDDM studies

It is challenging to compare similar PPDDM methods where there is a lack of key parameters presented. For instance, approaches which are designed for semi-honest parties might not be comparable with the approaches aiming to handle malicious behavior. The privacy-preserving methods for semi-honest parties will fail if involved parties show malicious behavior such as manipulating the input or output or completely aborting the protocol. Thus, the allowed adversarial behavior of participating parties is essential to be explicitly stated in the PPDDM papers. To consider all key parameters in PPDDM techniques, we provide a list of recommendations for the reporting of studies proposing new PPDDM methods or improving existing PPDDM methods as Table 2 shows. The recommendations detail the key parameters that should be described in each section of the paper of PPDDM. The factors in Table 2 refer to the 10 factors in the evaluation criteria which were discussed in the Methodology Section.

5.6. Potential limitations

The findings of this review have to be seen in light of some potential limitations. First, the 231 reviewed studies were searched from only 6 digital bibliographic databases (IEEE Xplore Digital Library, ACM Digital Library, Science Direct, ISI Web of Science, SpringerLink, and PubMed) and must be peer-reviewed publications. Some relevant studies may be missed in this review because they were not findable in these 6 bibliographic databases during searching. Studies that have not been peer-reviewed such as relevant articles published on arXiv.org⁷ were excluded.

Second, we did not apply an iterative “snowballing” approach to further identify more relevant studies [45]. “Snowballing” searching includes 1) reference tracking which identifies relevant studies from the reference lists of the primarily selected papers, 2) citation tracking which identifies relevant articles that cite primarily selected papers. We decided not to apply “snowballing” approach is because it may introduce a bias in favour of what authors think is relevant to their narrative [30]. Contrary, omitting the “snowballing” approach results in omitting follow-up studies of the reviewed papers. We decided to choose the latter approach, as we deemed our search criteria to be broad enough to cover follow-up studies. We have found several follow-up papers, where these papers present an extension of their existing methods to: 1) solve other data partitioning problems [125,126]; 2) apply to more advanced data analysis algorithms [62,64]; 3) to include more complicated user scenarios [67,68]; 4) to conduct more experiments by using real-life datasets [25,59,96,109].

Moreover, due to the scope of this review (providing a general overview of existing PPDDM methods and identifying outstanding challenges), more details of some privacy-preserving methods were not extensively discussed. For instance, in the category of ‘local learning and global integration’, multiple different methods can be applied to integrate the local miner (model) into a global miner (model) such as stacked generalization [119] and meta-learning [14]. In our belief this field warrants a separate in-depth review. Additionally, it has been well-recognized that there is an important trade-off between leakage of information and effectiveness or efficiency of learning in PPDDM technologies [15,66,77,123]. In practice, it is crucial to balance this trade-off depending on the specific use cases, the purposes of the data analysis, and the urgency of the problems. Although we included the privacy and efficiency factors in our review, we did not further investigate how each method weights the trade-off between them. For example, we did not measure how much and in which way information loss was tolerated to increase

⁷arXiv – a free distribution service and an open-access archive: <https://arxiv.org/>.

Table 2

A list of recommendations for reporting PPDDM studies

Section	Factor	Recommendations
Title and abstract		
Title and keywords	2,7	Identify the study as developing new or improving existing PPDDM algorithms to solve which data problem by using which type of partitioned data in a privacy-preserving manner
Abstract	1,2,4,6,7	Summarize the problems, objectives covering assumed adversarial behavior of data parties, data partitioning, brief description about privacy-preserving method, data mining algorithms, and applied dataset in the experiments.
Introduction		
Problem statement and background	2,3,5,6	Describe how data partitioned in which domain are considered by this study, what privacy issues are involved in that domain, which data mining algorithm is studied to solve what problems. Additionally, the number of participating parties and if all parties or only some parties have the target class should be also covered by this section.
Objectives and study design	1,3,4,7	Specify the objectives and study design include what level of privacy (or information leakage) is preserved against what adversarial behavior, applied privacy-preserving methods, evaluation criteria (for accuracy, efficiency, and privacy level), applied datasets in the experiments.
Methods		
Method design	4,5,6	Clearly explain which privacy-preserving methods are applied including the specific protocols/structures, proofs of preserving information leakage. Then, describe how certain data mining algorithms are adapted to combine with privacy-preserving methods, what information is communicated among parties, and complexity in different scenarios such as using categorical or numerical data, or involving different numbers of data parties. Lastly, make the code publicly available so that other researchers can reproduce the work.
Data	7	Describe data sources (and where and how other researchers can request the same dataset), the type and size of the datasets, basic description about data, what the target features/attributes are, missing values, and other basic information about the datasets.
Data analysis design	5,6	If real-life datasets are applied in the study, this subsection should describe the pre-processing of features/attributes (such as normalization, re-sampling), data analysis algorithms, parameter setting, and so on with reference to other comparable studies.
Experiment design	7	Describe how the datasets are partitioned (both feature-wise and instances-wise), how data parties communicate/transfer files, what validation is used, and what machine(such as CPU, memory) and software(versions) are used to do the experiments. In addition, experiments should be set up to compare with other existing PPDDM methods, or compare with privacy-preserving centralised data mining methods, or compare with distributed data mining methods without preserving privacy.
Evaluation design	8,9,10	Describe the evaluations of accuracy, efficiency (computational complexity, time cost on computation and communication among parties), privacy/security (such as information disclosure measurement)
Result		
Discovery from datasets	7	If real-life datasets are applied in the study, this subsection should describe what new knowledge was obtained from their analysis
Model performance	8	Present the performance measures such as accuracy scores of the proposed models in comparison with other existing PPDDM methods, or privacy-preserving centralised data mining methods, or distributed data mining methods without preserving privacy. Performance will be presented based on the evaluation criteria which was described in the methods section.
Privacy and/or security analysis	9	Provide sufficient privacy/security analysis based on the assumed adversarial behavior (semi-honest or malicious). Describe what information is exchanged among parties, what can be learnt from the exchanged information, if the models as a final outcome can cause information leakage, what the potential risks exist during the training process or in the final model.

Table 2
(Continued)

Section	Factor	Recommendations
Scalability analysis	10	Present the computation complexity and time consumption of the methods and describe what the volume (number of instances) and variety (number of features/attributes) of data can be handled by the proposed methods
Discussion		
Limitations	/	Discuss any limitations of proposed methods such as special cases where the methods are not applicable or certain assumptions which are not common in practice.
Interpretation	/	If real-life datasets are applied in the study, this subsection should discuss the findings with reference to any other validation data from other studies. Then, interpret the model performance on accuracy, efficiency, feasibility in practice, strengths and weaknesses with reference to other existing PPDDM methods.
Implementation	/	Discuss what other resources, paperwork, or supports are needed to implement the proposed methods, what potential challenges or risks will appear if apply the methods on real-life data.

efficiency. We believe this specific trade-off issue between privacy (information leakage) and learning performance (effectiveness or efficiency) deserves further investigation.

6. Conclusion

Privacy-preserving distributed data mining (PPDDM) techniques consider the issue of executing data mining algorithms on private, sensitive, and/or confidential data from multiple data parties while maintaining privacy. This review presented a comprehensive overview of current PPDDM methods to help researchers better understand the development of this domain and assist practitioners to select the suitable solutions for their practical cases.

In this review, we discovered there is a lack of standard criteria for evaluating new PPDDM techniques. The previous studies applied a variety of different evaluation methods, which brings challenges to objectively comparing existing PPDDM techniques. Therefore, an comprehensive evaluation criteria was proposed in this review including 10 key factors – adversarial behavior of data parties, data partitioning, experiment datasets, privacy/security analysis, privacy-preserving methods, data mining problems, analysis algorithms, complexity and cost, performance measures, and scalability to assess 231 recent studies published between 2000 to 2020 (August). We highlighted the characteristics of the 18 most cited studies and analyzed their influence on other studies in the field. Furthermore, a variety of definitions of privacy and distinguishment between information privacy and information security in the PPDDM field were discussed in this review, followed by some suggestions of making applicable privacy descriptions for new PPDDM methods. Finally, we also provided a list of recommendations for future research such as explicitly describing the privacy aspect under consideration, and evaluating new approaches using real-life data to narrow the gap between theoretical solutions and practical applications.

Acknowledgements

Financial support for this study was provided by a grant from the Dutch National Research Agenda (NWA; project number: 400.17.605). We gratefully acknowledge the time and effort devoted by two reviewers Dr. Abdur Rahim and Dr. Dayana Spagnuolo for their valuable comments. We would like

to thank Dr. Leto Peel for his generous feedback and suggestions to help us improve the quality of the manuscript. Special thanks are given to Dr. Amrapali Zaveri for her constructive suggestions on preparing the manuscript.

References

- [1] Y. Abdul Alsaheb, S. Aldeen, M. Salleh and M. Abdur Razzaque, A comprehensive review on privacy preserving data mining, *SpringerPlus* **4**(1) (2015), 1–36. doi:[10.1186/2193-1801-4-1](https://doi.org/10.1186/2193-1801-4-1).
- [2] N.R. Adam and J.C. Worthmann, Security-control methods for statistical databases: A comparative study, *ACM Computing Surveys (CSUR)* **21**(4) (1989), 515–556. doi:[10.1145/76894.76895](https://doi.org/10.1145/76894.76895).
- [3] J.S. Ancker, M.-H. Kim, Y. Zhang, Y. Zhang and J. Pathak, The potential value of social determinants of health in predicting health outcomes, *Journal of the American Medical Informatics Association* **25**(8) (2018), 1109–1110. doi:[10.1093/jamia/ocy061](https://doi.org/10.1093/jamia/ocy061).
- [4] M.J. Atallah and W. Du, Secure multi-party computational geometry, in: *Workshop on Algorithms and Data Structures*, Springer, 2001, pp. 165–179. doi:[10.1007/3-540-44634-6_16](https://doi.org/10.1007/3-540-44634-6_16).
- [5] D. Beaver, S. Micali and P. Rogaway, The round complexity of secure protocols, in: *Proceedings of the Twenty-Second Annual ACM Symposium on Theory of Computing*, Association for Computing Machinery, 1990, pp. 503–513. doi:[10.1145/100216.100287](https://doi.org/10.1145/100216.100287).
- [6] J.S. Beckmann and D. Lew, Reconciling evidence-based medicine and precision medicine in the era of big data: Challenges and opportunities, *Genome Medicine* **8**(1) (2016), 1–11. doi:[10.1186/s13073-016-0388-7](https://doi.org/10.1186/s13073-016-0388-7).
- [7] M. Ben-Or, S. Goldwasser and A. Wigderson, Completeness theorems for non-cryptographic fault-tolerant distributed computation, in: *Providing Sound Foundations for Cryptography: On the Work of Shafi Goldwasser and Silvio Micali*, Association for Computing Machinery, 2019, pp. 351–371. doi:[10.1145/3335741.3335756](https://doi.org/10.1145/3335741.3335756).
- [8] F.-Z. Benjelloun and A.A. Lahcen, Big data security: Challenges, recommendations and solutions, in: *Web Services: Concepts, Methodologies, Tools, and Applications*, IGI Global, 2019, pp. 25–38. doi:[10.4018/978-1-4666-8387-7.CH014](https://doi.org/10.4018/978-1-4666-8387-7.CH014).
- [9] E. Bertino, D. Lin and W. Jiang, A survey of quantification of privacy preserving data mining algorithms, in: *Privacy-Preserving Data Mining*, Springer, 2008, pp. 183–205. doi:[10.1007/978-0-387-70992-5_8](https://doi.org/10.1007/978-0-387-70992-5_8).
- [10] E. Bertino and I. Nai Fovino, Information driven evaluation of data hiding algorithms, in: *International Conference on Data Warehousing and Knowledge Discovery*, Springer, 2005, pp. 418–427. doi:[10.1007/11546849_41](https://doi.org/10.1007/11546849_41).
- [11] A. Bhowmick, J. Duchi, J. Freudiger, G. Kapoor and R. Rogers, Protection against reconstruction and its applications in private federated learning, 2018, arXiv preprint [arXiv:1812.00984](https://arxiv.org/abs/1812.00984).
- [12] Binnenlandse Zaken en Koninkrijksrelaties, Wet van 21 juli 2007, houdende algemene bepalingen betreffende de toekenning, het beheer en het gebruik van het burgerservicenummer (wet algemene bepalingen burgerservicenummer), 2018-07-28. <https://wetten.overheid.nl/jci1.3:c:BWBR0022428&z=2018-07-28&g=2018-07-28>.
- [13] A. Botchkarev, Performance metrics (error measures) in machine learning regression, forecasting and prognostics: Properties and typology, 2018, arXiv preprint [arXiv:1809.03006](https://arxiv.org/abs/1809.03006).
- [14] P.K. Chan, S.J. Stolfo et al., Toward parallel and distributed learning by meta-learning, in: *AAAI Workshop in Knowledge Discovery in Databases*, 1993, pp. 227–240. <https://dl.acm.org/doi/10.5555/3000767.3000789#d49627527e1>.
- [15] Y. Cheng, Y. Liu, T. Chen and Q. Yang, Federated learning for privacy-preserving ai, *Communications of the ACM* **63**(12) (2020), 33–36. doi:[10.1145/3387107](https://doi.org/10.1145/3387107).
- [16] P. Christen, *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*, Springer, 2012. doi:[10.1007/978-3-642-31164-2](https://doi.org/10.1007/978-3-642-31164-2).
- [17] E.A. Clarke, What is preventive medicine?, *Canadian Family Physician* **20**(11) (1974), 65. PMID:[20469128](https://pubmed.ncbi.nlm.nih.gov/20469128/).
- [18] C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin and M.Y. Zhu, Tools for privacy preserving distributed data mining, *ACM Sigkdd Explorations Newsletter* **4**(2) (2002), 28–34. doi:[10.1145/772862.772867](https://doi.org/10.1145/772862.772867).
- [19] Commission on Social Determinants of Health et al. *Closing the gap in a generation: health equity through action on the social determinants of health: final report of the commission on social determinants of health*, World Health Organization, 2008. https://www.who.int/social_determinants/final_report/csdh_finalreport_2008.pdf.
- [20] M.M. Cruz-Cunha, *Handbook of Research on Digital Crime, Cyberspace Security, and Information Assurance*, IGI Global, 2014. doi:[10.4018/978-1-4666-6324-4](https://doi.org/10.4018/978-1-4666-6324-4).
- [21] E. Czeizler, W. Wiessler, T. Koester, M. Hakala, S. Basiri, P. Jordan and E. Kuusela, Using federated data sources and varian learning portal framework to train a neural network model for automatic organ segmentation, *Physica Medica* **72** (2020), 39–45. doi:[10.1016/j.ejmp.2020.03.011](https://doi.org/10.1016/j.ejmp.2020.03.011).
- [22] Da Meng, K. Sivakumar and H. Kargupta, Privacy-sensitive Bayesian network parameter learning, in: *Fourth IEEE International Conference on Data Mining (ICDM'04)*, IEEE, 2004, pp. 487–490. doi:[10.1109/ICDM.2004.10076](https://doi.org/10.1109/ICDM.2004.10076).

- [23] T. Dalenius and S.P. Reiss, Data-swapping: A technique for disclosure control, *Journal of Statistical Planning and Inference* **6**(1) (1982), 73–85. doi:10.1016/0378-3758(82)90058-1.
- [24] J.W. DeCew, *In Pursuit of Privacy: Law, Ethics, and the Rise of Technology*, Cornell University Press, 1997. <https://www.jstor.org/stable/10.7591/j.ctv75d3zc>.
- [25] T.M. Deist, F.J.W.M. Dankers, P. Ojha, M. Scott Marshall, T. Janssen, C. Faivre-Finn, C. Masciocchi, V. Valentini, J. Wang, J. Chen et al., Distributed learning on 20000+ lung cancer patients – the personal health train, *Radiotherapy and Oncology* **144** (2020), 189–200. doi:10.1016/j.radonc.2019.11.019.
- [26] Y. Djemaiel, S. Rekhis and N. Boudriga, Trustworthy networks, authentication, privacy, and security models, in: *Handbook of Research on Wireless Security*, IGI Global, 2008, pp. 189–209. doi:10.4018/978-1-59904-899-4.ch014.
- [27] S. Dolley, Big data's role in precision public health, *Frontiers in Public Health* **6** (2018), 68. doi:10.3389/fpubh.2018.00068.
- [28] Y. Dong, X. Chen, L. Shen and D.W. Eastfly, Efficient and secure ternary federated learning, *Computers & Security* **94** (2020), 101824. doi:10.1016/j.cose.2020.101824.
- [29] G. Dougherty, *Digital Image Processing for Medical Applications*, Cambridge University Press, 2009. doi:10.1017/CBO9780511609657.
- [30] M. Egger, G. Davey-Smith and D. Altman, *Systematic Reviews in Health Care: Meta-Analysis in Context*, John Wiley & Sons, 2008. doi:10.1002/9780470693926.
- [31] H. Elshazly, A.T. Azar, A. El-Korany and A.E. Hassanien, Hybrid system for lymphatic diseases diagnosis, in: *2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, IEEE, 2013, pp. 343–347. doi:10.1109/ICACCI.2013.6637195.
- [32] European Commission, White paper on artificial intelligence: A european approach to excellence and trust. Technical report, European Commission, 2020. https://ec.europa.eu/info/sites/default/files/commission-white-paper-artificial-intelligence-feb2020_en.pdf.
- [33] S. Even O. Goldreich and A. Lempel, A randomized protocol for signing contracts, *Communications of the ACM* **28**(6) (1985), 637–647. doi:10.1145/3812.3818.
- [34] W. Fang and B. Yang, Privacy preserving decision tree learning over vertically partitioned data, in: *2008 International Conference on Computer Science and Software Engineering*, Vol. 3, IEEE, 2008, pp. 1049–1052. doi:10.1109/CSSE.2008.731.
- [35] S.E. Fienberg and J. McIntyre, Data swapping: Variations on a theme by dalenius and reiss, in: *International Workshop on Privacy in Statistical Databases*, Springer, 2004, pp. 14–29. doi:10.1007/978-3-540-25955-8_2.
- [36] S. Fletcher and M.Z. Islam, Measuring information quality for privacy preserving data mining, *International Journal of Computer Theory and Engineering* **7**(1) (2015), 21. doi:10.7763/IJCTE.2015.V7.924.
- [37] S. Fletcher and M.Z. Islam, Decision tree classification with differential privacy: A survey, *ACM Computing Surveys (CSUR)* **52**(4) (2019), 1–33. doi:10.1145/3337064.
- [38] A.A. Freitas, A survey of evolutionary algorithms for data mining and knowledge discovery, in: *Advances in Evolutionary Computing*, Springer, 2003, pp. 819–845. doi:10.1007/978-3-642-18965-4_33.
- [39] J. Fürnkranz and P.A. Flach, An analysis of rule evaluation metrics, in: *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, 2003, pp. 202–209. <https://www.aaai.org/Papers/ICML/2003/ICML03-029.pdf>.
- [40] D. Gao, Y. Liu, A. Huang, C. Ju, H. Yu and Q. Yang, Privacy-preserving heterogeneous federated transfer learning, in: *2019 IEEE International Conference on Big Data (Big Data)*, IEEE, 2019, pp. 2552–2559. doi:10.1109/BigData47090.2019.9005992.
- [41] A. Gascón, P. Schoppmann, B. Balle, M. Raykova, J. Doerner, S. Zahur and D. Evans, Privacy-preserving distributed linear regression on high-dimensional data, in: *Proceedings on Privacy Enhancing Technologies*, Vol. 2017, 2017, pp. 345–364. doi:10.1007/978-3-540-71701-0.
- [42] C. Gentry et al., *A Fully Homomorphic Encryption Scheme*, Vol. 20, Stanford University, Stanford, 2009. <https://crypto.stanford.edu/craig/craig-thesis.pdf>.
- [43] Gephi – The Open Graph Viz Platform, <https://gephi.org/>.
- [44] O. Goldreich, *Foundations of Cryptography: Volume 2, Basic Applications*, Cambridge University Press, 2009. doi:10.1017/CBO9780511721656.
- [45] T. Greenhalgh and R. Peacock, Effectiveness and efficiency of search methods in systematic reviews of complex evidence: Audit of primary sources, *Bmj* **331**(7524) (2005), 1064–1065. doi:10.1136/bmj.38636.593461.68.
- [46] B. Gu, Z. Dang, X. Li and H. Huang, Federated doubly stochastic kernel learning for vertically partitioned data, in: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Association for Computing Machinery, 2020, pp. 2483–2493. doi:10.1145/3394486.3403298.
- [47] R. Hall and S.E. Fienberg, Privacy-preserving record linkage, in: *International Conference on Privacy in Statistical Databases*, Springer, 2010, pp. 269–283. doi:10.1007/978-3-642-15838-4_24.
- [48] J. Han, M. Kamber and J. Pei, *Data Mining Concepts and Techniques*, 3rd edn, Elsevier, 2011. doi:10.1016/C2009-0-61819-5.

- [49] M. Hossin and M.N. Sulaiman, A review on evaluation metrics for data classification evaluations, *International Journal of Data Mining & Knowledge Management Process* **5**(2) (2015), 1. doi:[10.5121/ijdkp.2015.5201](https://doi.org/10.5121/ijdkp.2015.5201).
- [50] Y. Hu, D. Niu, J. Yang and S.Z. Fdml, A collaborative machine learning framework for distributed features, in: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Association for Computing Machinery, 2019, pp. 2232–2240. doi:[10.1145/3292500.3330765](https://doi.org/10.1145/3292500.3330765).
- [51] K. Huang and R. Tso, A commutative encryption scheme based on elgamal encryption, in: *2012 International Conference on Information Security and Intelligent Control*, IEEE, 2012, pp. 156–159. doi:[10.1109/ISIC.2012.6449730](https://doi.org/10.1109/ISIC.2012.6449730).
- [52] L. Huang, A.L. Shea, H. Qian, A. Masurkar, H. Deng and D. Liu, Patient clustering improves efficiency of federated machine learning to predict mortality and hospital stay time using distributed electronic medical records, *Journal of Biomedical Informatics* **99** (2019), 103291. doi:[10.1016/j.jbi.2019.103291](https://doi.org/10.1016/j.jbi.2019.103291).
- [53] A. Inan, S.V. Kaya, Y. Saygin, E. Savaş, A.A. Hintoğlu and A. Levi, Privacy preserving clustering on horizontally partitioned data, *Data & Knowledge Engineering* **63**(3) (2007), 646–666. doi:[10.1016/j.datak.2007.03.015](https://doi.org/10.1016/j.datak.2007.03.015).
- [54] I. Ioannidis, A. Grama and M. Atallah, A secure protocol for computing dot-products in clustered and distributed environments, in: *Proceedings International Conference on Parallel Processing*, IEEE, 2002, pp. 379–384. doi:[10.1109/ICPP.2002.1040894](https://doi.org/10.1109/ICPP.2002.1040894).
- [55] G. Jagannathan and R.N. Wright, Privacy-preserving distributed k-means clustering over arbitrarily partitioned data, in: *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, Association for Computing Machinery, 2005, pp. 593–599. doi:[10.1145/1081870.1081942](https://doi.org/10.1145/1081870.1081942).
- [56] P. Jain, M. Gyanchandani and N. Khare, Big data privacy: A technological perspective and review, *Journal of Big Data* **3**(1) (2016), 1–25. doi:[10.1186/s40537-016-0059-y](https://doi.org/10.1186/s40537-016-0059-y).
- [57] H. Jeff Smith, T. Dinev and H. Xu, Information privacy research: An interdisciplinary review, *MIS quarterly* (2011), 989–1015. doi:[10.2307/41409970](https://doi.org/10.2307/41409970).
- [58] W. Jiang and M. Atzori, Secure distributed k-anonymous pattern mining, in: *Sixth International Conference on Data Mining (ICDM'06)*, IEEE, 2006, pp. 319–329. doi:[10.1109/ICDM.2006.140](https://doi.org/10.1109/ICDM.2006.140).
- [59] A. Jochems, T.M. Deist, J. Van Soest, M. Eble, P. Bulens, P. Coucke, W. Dries, P. Lambin and A. Dekker, Distributed learning: Developing a predictive model based on data from multiple hospitals without data leaving the hospital – a real life proof of concept, *Radiotherapy and Oncology* **121**(3) (2016), 459–467. doi:[10.1016/j.radonc.2016.10.002](https://doi.org/10.1016/j.radonc.2016.10.002).
- [60] G.A. Kaissis, M.R. Makowski, D. Rückert and R.F. Braren, Secure, privacy-preserving and federated machine learning in medical imaging, *Nature Machine Intelligence* **2**(6) (2020), 305–311. doi:[10.1038/s42256-020-0186-1](https://doi.org/10.1038/s42256-020-0186-1).
- [61] A. Karakasidis and V.S. Verykios, A sorted neighborhood approach to multidimensional privacy preserving blocking, in: *2012 IEEE 12th International Conference on Data Mining Workshops*, IEEE, 2012, pp. 937–944. doi:[10.1109/ICDMW.2012.70](https://doi.org/10.1109/ICDMW.2012.70).
- [62] H. Kikuchi, C. Hamanaga, H. Yasunaga, H. Matsui and H. Hashimoto, Privacy-preserving multiple linear regression of vertically partitioned real medical datasets, in: *2017 IEEE 31st International Conference on Advanced Information Networking and Applications (AINA)*, 2017, pp. 1042–1049. doi:[10.1109/AINA.2017.52](https://doi.org/10.1109/AINA.2017.52).
- [63] H. Kikuchi, C. Hamanaga, H. Yasunaga, H. Matsui, H. Hashimoto and C.-I. Fan, Privacy-preserving multiple linear regression of vertically partitioned real medical datasets, *Journal of Information Processing* **26** (2018), 638–647. doi:[10.2197/ipsjip.26.638](https://doi.org/10.2197/ipsjip.26.638).
- [64] H. Kikuchi, H. Hashimoto, H. Yasunaga and T. Saito, Scalability of privacy-preserving linear regression in epidemiological studies, in: *2015 IEEE 29th International Conference on Advanced Information Networking and Applications*, 2015, pp. 510–514. doi:[10.1109/AINA.2015.229](https://doi.org/10.1109/AINA.2015.229).
- [65] B. Kitchenham, Procedures for performing systematic reviews, *Keele, UK, Keele University* **33**(2004) (2004), 1–26. http://artemisa.unicauca.edu.co/~ecalton/docs/spi/kitchenham_2004.pdf.
- [66] J. Konečný, H.B. McMahan, F.X. Yu and P. Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency, 2016, arXiv preprint [arXiv:1610.05492](https://arxiv.org/abs/1610.05492).
- [67] S.X. Lee, K.L. Leemaqz and G.J. McLachlan, Ppem: Privacy-preserving em learning for mixture models, *Concurrency and Computation: Practice and Experience* **31**(24) (2019), e5208. doi:[10.1002/cpe.5208](https://doi.org/10.1002/cpe.5208).
- [68] K.L. Leemaqz, S.X. Lee and G.J. McLachlan, Corruption-resistant privacy preserving distributed em algorithm for model-based clustering, in: *2017 IEEE Trustcom/BigDataSE/ICSS*, 2017, pp. 1082–1089. doi:[10.1109/Trustcom/BigDataSE/ICSS.2017.356](https://doi.org/10.1109/Trustcom/BigDataSE/ICSS.2017.356).
- [69] J. Li, Y. Tian, Y. Zhu, T. Zhou, J. Li, K. Ding and J. Li, A multicenter random forest model for effective prognosis prediction in collaborative clinical research network, *Artificial Intelligence in Medicine* **103** (2020), 101814. doi:[10.1016/j.artmed.2020.101814](https://doi.org/10.1016/j.artmed.2020.101814).
- [70] L. Li, L. Huang, W. Yang, X. Yao and A. Liu, Privacy-preserving lof outlier detection, *Knowledge and Information Systems* **42**(3) (2015), 579–597. doi:[10.1007/s10115-013-0692-0](https://doi.org/10.1007/s10115-013-0692-0).
- [71] X. Lin, C. Clifton and M. Zhu, Privacy-preserving clustering with distributed em mixture modeling, *Knowledge and Information Systems* **8**(1) (2005), 68–81. doi:[10.1007/s10115-004-0148-7](https://doi.org/10.1007/s10115-004-0148-7).

- [72] Y. Lindell, Secure multiparty computation for privacy preserving data mining, in: *Encyclopedia of Data Warehousing and Mining*, IGI Global, 2005, pp. 1005–1009. doi:[10.4018/978-1-59140-557-3.ch189](https://doi.org/10.4018/978-1-59140-557-3.ch189).
- [73] K. Liu, H. Kargupta and J. Ryan, Random projection-based multiplicative data perturbation for privacy preserving distributed data mining, *IEEE Transactions on Knowledge and Data Engineering* **18**(1) (2005), 92–106. doi:[10.1109/TKDE.2006.14](https://doi.org/10.1109/TKDE.2006.14).
- [74] Y. Lu, P. Phoungphol and Y. Zhang, Privacy aware non-linear support vector machine for multi-source big data, in: *2014 IEEE 13th International Conference on Trust, Security and Privacy in Computing and Communications*, IEEE, 2014, pp. 783–789. doi:[10.1109/TrustCom.2014.103](https://doi.org/10.1109/TrustCom.2014.103).
- [75] S.T. Margulis, Conceptions of privacy: Current status and next steps, *Journal of Social Issues* **33**(3) (1977), 5–21. doi:[10.1111/j.1540-4560.1977.tb01879.x](https://doi.org/10.1111/j.1540-4560.1977.tb01879.x).
- [76] B. McMahan, E. Moore, D. Ramage, S. Hampson and B.A. y Arcas, Communication-efficient learning of deep networks from decentralized data, in: *Artificial Intelligence and Statistics*, PMLR, 2017, pp. 1273–1282. <http://proceedings.mlr.press/v54/mcmahan17a/mcmahan17a.pdf>.
- [77] R. Mendes and J.P. Vilela, Privacy-preserving data mining: Methods, metrics, and applications, *IEEE Access* **5** (2017), 10562–10582. doi:[10.1109/ACCESS.2017.2706947](https://doi.org/10.1109/ACCESS.2017.2706947).
- [78] S. Merugu and J. Ghosh, Privacy-preserving distributed clustering using generative models, in: *Third IEEE International Conference on Data Mining*, IEEE, 2003, pp. 211–218. doi:[10.1109/ICDM.2003.1250922](https://doi.org/10.1109/ICDM.2003.1250922).
- [79] S. Micali, O. Goldreich and A. Wigderson, How to play any mental game, in: *Proceedings of the Nineteenth ACM Symp. on Theory of Computing, STOC*, Association for Computing Machinery, 1987, pp. 218–229. doi:[10.1145/28395.28420](https://doi.org/10.1145/28395.28420).
- [80] K. Murat and C. Clifton, Privacy-preserving distributed mining of association rules on horizontally partitioned data, *IEEE Transactions on Knowledge and Data Engineering* **16**(9) (2004), 1026–1037. doi:[10.1109/TKDE.2004.45](https://doi.org/10.1109/TKDE.2004.45).
- [81] M. Ogburn, C. Turner and P. Dahal, Homomorphic encryption, *Procedia Computer Science* **20** (2013), 502–509. doi:[10.1016/j.procs.2013.09.310](https://doi.org/10.1016/j.procs.2013.09.310).
- [82] O. Ohrimenko, F. Schuster, C. Fournet, A. Mehta, S. Nowozin, K. Vaswani and M. Costa, Oblivious multi-party machine learning on trusted processors, in: *25th {USENIX} Security Symposium ({USENIX} Security 16)*, USENIX Association, 2016, pp. 619–636. <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/07/paper.pdf>.
- [83] J.-O. Palacio-Niño and F. Berzal, Evaluation metrics for unsupervised learning algorithms, 2019, arXiv preprint [arXiv:1905.05667](https://arxiv.org/abs/1905.05667).
- [84] D. Peteiro-Barral and B. Guijarro-Berdiñas, A survey of methods for distributed machine learning, *Progress in Artificial Intelligence* **2**(1) (2013), 1–11. doi:[10.1007/s13748-012-0035-5](https://doi.org/10.1007/s13748-012-0035-5).
- [85] S. Pohlig and M. Hellman, An improved algorithm for computing logarithms over $GF(p)$ and its cryptographic significance (corresp.), *IEEE Transactions on Information Theory* **24**(1) (1978), 106–110. doi:[10.1109/TIT.1978.1055817](https://doi.org/10.1109/TIT.1978.1055817).
- [86] R.L. Rivest, L. Adleman, M.L. Dertouzos et al., On data banks and privacy homomorphisms, *Foundations of Secure Computation* **4**(11) (1978), 169–180. <http://people.csail.mit.edu/rivest/RivestAdlemanDertouzos-OnDataBanksAndPrivacyHomomorphisms.pdf>.
- [87] B. Rozenberg and E. Gudes, Association rules mining in vertically partitioned databases, *Data & Knowledge Engineering* **59**(2) (2006), 378–396. doi:[10.1016/j.datak.2005.09.001](https://doi.org/10.1016/j.datak.2005.09.001).
- [88] R. Schnell, Efficient private record linkage of very large datasets, in: *59th World Statistics Congress of the International Statistical Institute*, International Statistical Institute, 2013. <https://openaccess.city.ac.uk/id/eprint/14652/>.
- [89] A. Shah and R. Gulati, Privacy preserving data mining: Techniques, classification and implications – a survey, *Int. J. Comput. Appl* **137**(12) (2016), 40–46. doi:[10.5120/IJCA2016909006](https://doi.org/10.5120/IJCA2016909006).
- [90] M. Shaneck, Y. Kim and V. Kumar, Privacy preserving nearest neighbor search, in: *Machine Learning in Cyber Trust*, Springer, 2009, pp. 247–276. doi:[10.1007/978-0-387-88735-7_10](https://doi.org/10.1007/978-0-387-88735-7_10).
- [91] S. Shelke and B. Bhagat, Techniques for privacy preservation in data mining, *International Journal of Engineering Research* **4**(10) (2015). doi:[10.17577/ijertv4is100473](https://doi.org/10.17577/ijertv4is100473).
- [92] S. Shen, T. Zhu, D. Wu, W. Wang and W. Zhou, From distributed machine learning to federated learning: In the view of data privacy and security. *Concurrency and Computation: Practice and Experience*, 2020. doi:[10.1002/cpe.6002](https://doi.org/10.1002/cpe.6002).
- [93] C.A. Shoniregun, K. Dube and F. Mtenzi, *Electronic Healthcare Information Security*, Vol. 53, Springer Science & Business Media, 2010. doi:[10.1007/978-0-387-84919-5](https://doi.org/10.1007/978-0-387-84919-5).
- [94] A.B. Slavkovic, Y. Nardi and M.M. Tibbits, “Secure” logistic regression of horizontally and vertically partitioned distributed databases, in: *Seventh IEEE International Conference on Data Mining Workshops (ICDMW 2007)*, IEEE, 2007, pp. 723–728. doi:[10.1109/ICDMW.2007.114](https://doi.org/10.1109/ICDMW.2007.114).
- [95] R.B. Stricker and L. Johnson, Lyme disease: The promise of big data, companion diagnostics and precision medicine, *Infection and Drug Resistance* **9** (2016), 215. doi:[10.2147/IDR.S114770](https://doi.org/10.2147/IDR.S114770).
- [96] C. Sun, L. Ippel, J. Van Soest, B. Wouters, A. Malic, O. Adekunle, B. van den Berg, O. Musmann, A. Koster, C. van der Kallen et al., A privacy-preserving infrastructure for analyzing personal health data in a vertically partitioned scenario, in: *MEDINFO 2019: Health and Wellbeing E-Networks for All: Proceedings of the 17th World Congress on Medical and Health Informatics*, Vol. 264, IOS Press, 2019, pp. 373–377. doi:[10.3233/SHTI190246](https://doi.org/10.3233/SHTI190246).

- [97] N. Suranga, Kasthurirathne, J.R. Vest, N. Menachemi, P.K. Halverson and S.J. Grannis, Assessing the capacity of social determinants of health data to augment predictive models identifying patients in need of wraparound social services, *Journal of the American Medical Informatics Association* **25**(1) (2018), 47–53. doi:10.1093/jamia/ocx130.
- [98] E. Suthampan and S. Maneewongvatana, Privacy preserving decision tree in multi party environment, in: *Asia Information Retrieval Symposium*, Springer, 2005, pp. 727–732. doi:10.1007/11562382_75.
- [99] M. Timo, Deist, A. Jochems, J. van Soest, G. Nalbantov, C. Oberije, S. Walsh, M. Eble, P. Bulens, P. Coucke, W. Dries et al., Infrastructure and distributed learning methodology for privacy-preserving multi-centric rapid learning health care: Eurocat, *Clinical and Translational Radiation Oncology* **4** (2017), 24–31. doi:10.1016/j.ctro.2016.12.004.
- [100] A. Tuladhar, S. Gill, Z. Ismail and N.D. Forkert, Alzheimer’s Disease Neuroimaging Initiative et al. Building machine learning models without sharing patient data: A simulation-based analysis of distributed learning by ensembling, *Journal of Biomedical Informatics* **106** (2020), 103424. doi:10.1016/j.jbi.2020.103424.
- [101] UCI Machine Learning Repository, <https://archive.ics.uci.edu/ml/index.php>.
- [102] H. Vaghashia and A. Ganatra, A survey: Privacy preservation techniques in data mining, *International Journal of Computer Applications* **119**(4) (2015). doi:10.5120/21056-3704.
- [103] J. Vaidya, A survey of privacy-preserving methods across vertically partitioned data, in: *Privacy-Preserving Data Mining*, Springer, 2008, pp. 337–358. doi:10.1007/978-0-387-70992-5_14.
- [104] J. Vaidya and C. Clifton, Privacy preserving association rule mining in vertically partitioned data, in: *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, 2002, pp. 639–644. doi:10.1145/775047.775142.
- [105] J. Vaidya and C. Clifton, Privacy-preserving k-means clustering over vertically partitioned data, in: *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, 2003, pp. 206–215. doi:10.1145/956750.956776.
- [106] J. Vaidya, C. Clifton, M. Kantarcioglu and A.S. Patterson, Privacy-preserving decision trees over vertically partitioned data, *ACM Transactions on Knowledge Discovery from Data (TKDD)* **2**(3) (2008), 1–27. doi:10.1145/1409620.1409624.
- [107] J. Vaidya, C.W. Clifton and Y.M. Zhu, *Privacy Preserving Data Mining*, Vol. 19, Springer Science & Business Media, 2006. doi:10.1007/978-0-387-29489-6.
- [108] J. Vaidya, H. Yu and X. Jiang, Privacy-preserving svm classification, *Knowledge and Information Systems* **14**(2) (2008), 161–178. doi:10.1007/s10115-007-0073-7.
- [109] J. Van Soest, C. Sun, O. Mussmann, M. Puts, B. van den Berg, A. Malic, C. van Oppen, D. Townend, A. Dekker and M. Dumontier, Using the personal health train for automated and privacy-preserving analytics on vertically partitioned data, in: *Building Continents of Knowledge in Oceans of Data: The Future of Co-Created eHealth*, Vol. 247, IOS Press, 2018, pp. 581–585. doi:10.3233/978-1-61499-852-5-581.
- [110] S. Vassiliou, V. Verykios, E. Bertino, I. Nai Fovino, L. Parasiliti Provenza, Y. Saygin and Y. Theodoridis, State-of-the-art in privacy preserving data mining, *ACM Sigmod Record* **33**(1) (2004), 50–57. doi:10.1145/974121.974131.
- [111] D. Vatsalan, P. Christen and V.S. Verykios, A taxonomy of privacy-preserving record linkage techniques, *Information Systems* **38**(6) (2013), 946–969. doi:10.1016/j.is.2012.11.005.
- [112] J. Verbraeken, M. Wolting, J. Katzy, J. Kloppenburg, T. Verbelen and J.S. Rellermeyer, A survey on distributed machine learning, *ACM Computing Surveys (CSUR)* **53**(2) (2020), 1–33. doi:DOI:10.1145/3377454.
- [113] J. Wang, *Encyclopedia of Data Warehousing and Mining*, IGI Global, 2005. doi:10.1108/14684520610675852.
- [114] R. Wang, W. Ji, M. Liu, X. Wang, J. Weng, S. Deng, S. Gao and C. Yuan, Review on mining data from multiple data sources, *Pattern Recognition Letters* **109** (2018), 120–128. doi:10.1016/j.patrec.2018.01.013.
- [115] A.F. Westin, Privacy and freedom, *Washington and Lee Law Review* **25**(1) (1968), 166. doi:10.2307/3102188.
- [116] R.L. Wilson and P.A. Rosen, Protecting data through perturbation techniques: The impact on knowledge discovery in databases, *Journal of Database Management (JDM)* **14**(2) (2003), 14–26. doi:10.4018/jdm.2003040102.
- [117] W.E. Winkler, Record linkage, in: *Handbook of Statistics*, Vol. 29, Elsevier, 2009, pp. 351–380. doi:10.1016/S0169-7161(08)00014-X.
- [118] M. Wolfson, S.E. Wallace, N. Masca, G. Rowe, N.A. Sheehan, V. Ferretti, P. LaFlamme, M.D. Tobin, J. Macleod, J. Little et al., Datashield: Resolving a conflict in contemporary bioscience – performing a pooled analysis of individual-level data without sharing the data, *International Journal of Epidemiology* **39**(5) (2010), 1372–1382. doi:10.1093/ije/dyq111.
- [119] D.H. Wolpert, Stacked generalization, *Neural Networks* **5**(2) (1992), 241–259. doi:10.1016/S0893-6080(05)80023-1.
- [120] R. Wright and Z. Yang, Privacy-preserving Bayesian network structure computation on distributed heterogeneous data, in: *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, 2004, pp. 713–718. doi:10.1145/1014052.1014145.
- [121] M.-J. Xiao, L.-S. Huang, Y.-L. Luo and H. Shen, Privacy preserving id3 algorithm over horizontally partitioned data, in: *Sixth International Conference on Parallel and Distributed Computing Applications and Technologies (PDCAT’05)*, IEEE, 2005, pp. 239–243. doi:10.1109/PDCAT.2005.191.

- [122] X. Xiao and Y. Tao, Personalized privacy preservation, in: *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*, Association for Computing Machinery, 2006, pp. 229–240. doi:[10.1145/1142473.1142500](https://doi.org/10.1145/1142473.1142500).
- [123] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen and H. Yu, Federated learning, *Synthesis Lectures on Artificial Intelligence and Machine Learning* **13**(3) (2019), 1–207. doi:[10.2200/S00960ED2V01Y201910AIM043](https://doi.org/10.2200/S00960ED2V01Y201910AIM043).
- [124] A.C.-C. Yao, How to generate and exchange secrets, in: *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*, IEEE, 1986, pp. 162–167. doi:[10.1109/SFCS.1986.25](https://doi.org/10.1109/SFCS.1986.25).
- [125] H. Yu, X. Jiang and J. Vaidya, Privacy-preserving svm using nonlinear kernels on horizontally partitioned data, in: *Proceedings of the 2006 ACM Symposium on Applied Computing*, Association for Computing Machinery, 2006, pp. 603–610. doi:[10.1145/1141277.1141415](https://doi.org/10.1145/1141277.1141415).
- [126] H. Yu, J. Vaidya and X. Jiang, Privacy-preserving svm classification on vertically partitioned data, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer, 2006, pp. 647–656. doi:[10.1007/11731139_74](https://doi.org/10.1007/11731139_74).
- [127] J. Zhan, S. Matwin and L. Chang, Privacy-preserving collaborative association rule mining, *Journal of Network and Computer Applications* **30**(3) (2007), 1216–1227. doi:[10.1016/j.jnca.2006.04.010](https://doi.org/10.1016/j.jnca.2006.04.010).
- [128] Q. Zhao, C. Zhao, S. Cui, S. Jing and Z.C. Privatedl, Privacy-preserving collaborative deep learning against leakage from gradient sharing, *International Journal of Intelligent Systems* **35**(8) (2020), 1262–1279. doi:[10.1002/int.22241](https://doi.org/10.1002/int.22241).