

Privacy-Preserving Record Linkage

Dinusha Vatsalan^{a,*}, Dimitrios Karapiperis^b, Vassilios S. Verykios^b

^a*School of Computing, Macquarie University, NSW 2109, Australia*

^b*School of Science and Technology, Hellenic Open University, Patras, Greece*

Definition

Given several databases containing person-specific data held by different organizations, Privacy-Preserving Record Linkage (PPRL) aims to identify and link records that correspond to the same entity/individual across different databases based on the matching of personal identifying attributes, such as name and address, without revealing the actual values in these attributes due to privacy concerns.

1. Synonyms

Private Data Matching, Private Record Linkage, Blind Data Linkage, Private Data Integration

2. Overview

In the current era of Big Data personal data about people, such as customers, patients, tax payers, and clients, are dispersed in multiple different sources collected by different organizations. Several applications have begun to leverage tremendous opportunities and insights provided by linked and integrated data. Examples range from healthcare, businesses, social sciences, to government services and national security. Linking data is also used as a pre-processing step in many data mining and analytics projects in order to clean, enrich, and understand data for quality results [6].

However, the growing concerns of privacy and confidentiality issues about personal data pose serious constraints to share and exchange such data across organizations for linking. Since a unique entity identifier is not available in different data sources, linking of records from different databases needs to rely on available personal identifying attributes, such as names, dates of birth, and addresses. Known as quasi identifiers (QIDs), these values in combination not only allow uniquely identifying individuals but also reveal private and sensitive information about them [44].

Privacy-preserving record linkage (PPRL) aims to identify the same entities from different databases by linking records based on encoded and/or encrypted QIDs, such that no sensitive information of the entities is revealed to any internal (parties involved in the process) and external (external adversaries and eavesdroppers) parties. While a variety of techniques and methodologies have been developed for PPRL over the past two decades, as surveyed in [44, 45], this research field is still open to several challenges, especially with the Big Data revolution.

*Corresponding author.

Email addresses: dinusha.vatsalan@mq.edu.au (Dinusha Vatsalan), dkarapiperis@eap.gr (Dimitrios Karapiperis), verykios@eap.gr (Vassilios S. Verykios)

Chapter published in Encyclopedia of Big Data Technologies

December 13, 2022

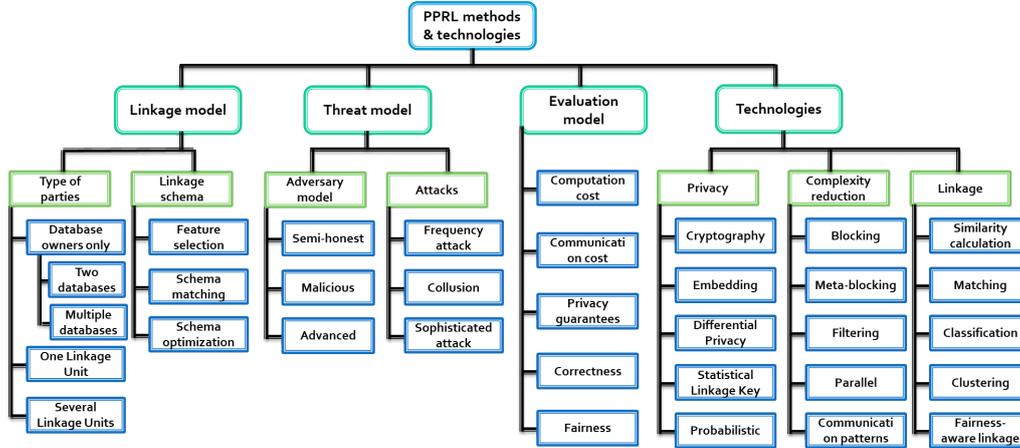


Figure 1: A taxonomy of methodologies and technologies used in PPRL.

A typical PPRL process involves several steps, starting from pre-processing databases, encoding or encrypting records using the privacy masking functions [44], applying blocking or other forms of complexity reduction methods [6], then matching records based on their QIDs using similarity functions [6], and finally clustering or classifying matching records that correspond to the same entity. Additional steps of evaluation of the linkage as well as manual review of certain records are followed in non-PPRL applications. However, these two steps require more research for PPRL as they generally require access to raw data and ground truth which is not possible in a privacy-preserving context.

Despite several challenges, the output of PPRL could certainly bring in enormous potential in the Big Data era for businesses, government agencies, and research organizations. PPRL has been an active area of research in the recent times and it is now being applied in real-world applications [5, 30]. In this chapter, we will describe some of the key research findings of PPRL and example PPRL applications. We will also discuss directions for future research in PPRL.

3. Key research findings

In this section, we present the key research findings in PPRL with regard to the methodologies and technologies used, as characterized in Figure 1.

3.1. Linkage model

Type of parties: Considering the number and type of parties involved, the linkage model can be categorized as two-party or multi-party protocols with or without one or several linkage units (LUs). Two database owners (DOs)/parties involve in two-party protocols in order to identify the matching records in their databases. This linkage model requires more sophisticated techniques to ensure that the parties do not learn any sensitive information about each other’s data. A LU is therefore commonly used in linkage models to conduct/facilitate the linkage. A major drawback of LU-based protocols is that they require a trusted LU in order to avoid possible collusion attacks, where one of the DOs collude with the LU to learn about other DO’s data [44]. In some models several (more than one) LUs are used, for example, one is responsible for secret

keys management, another one for facilitating the complexity reduction step, while matching of records is conducted by a different LU. Separating the tasks across several LUs reduces the amount of information learned by a single party, however collusion could compromise the privacy. In multi-party linkage the aim is to identify cluster of records from multiple (more than two) databases that refer to the same entity. The linkage process becomes complicated with the increase of number of databases. Processing can be distributed among multiple parties to improve efficiency as well as to improve privacy guarantees by reducing the amount of information learned by a single party.

Linkage schema: The linkage schema dimension consists of feature selection, schema matching, and schema optimization. Different types of QIDs have been used for PPRL, with the most commonly used QIDs include name (string), address (string or text), age (numeric), gender (categorical), and date of birth (date). Schema matching identifies the common schema across different databases [32]. The success of linkage depends on the parameter setting used which needs to be tuned appropriately in order to optimize the linkage results. Grid search and random search are two optimization methods, however they tune parameters in an isolated way disregarding past evaluations of parameter combinations [3]. Bayesian optimization can efficiently bring down the time spent to get the optimal set of parameters [36] by taking into account the information on the parameter combinations it has previously seen thus far when choosing the parameter set to evaluate next.

3.2. Threat model

Adversary model: Different adversary models are assumed in PPRL methodologies [44], which are categorized as honest-but-curious/semi-honest, malicious, and advanced models. The semi-honest model is the most commonly used adversary model in existing PPRL techniques [44], where the parties are assumed to honestly follow the steps of the protocol while being curious to learn from the information they received throughout the protocol. This model is not realistic in real applications due to the weak assumption of privacy against adversarial attacks. Malicious model, on the other hand, provides a strong assumption of privacy such that the parties involved in the protocol may not follow the steps of the protocol by deviating from the protocol, sending false input, or behaving arbitrarily. However, more complex and advanced privacy techniques are required to make the protocols resistant against such malicious adversaries. Hybrid models, such as accountable computing and covert models, lie in between the semi-honest model, which is not realistic, and the malicious model, which requires computationally expensive techniques [44].

Attacks: Several attacks have been developed for PPRL techniques to investigate the resistance of such techniques to those attacks. Frequency attacks are most commonly used, where the frequency of encoded values are mapped to the frequency of known unencoded values [41]. Collusion attacks are possible in LU-based and multi-party models where subsets of parties collude to learn another party's data [41, 29]. More sophisticated attacks have been developed against certain privacy techniques. For example, Bloom filters (as described below) are susceptible to cryptanalysis attacks, which allow the iterative mapping of bit patterns back to their original unencoded QID values based on their frequency alignments depending upon the parameter setting used [7, 23].

3.3. Evaluation model

Evaluation of the performance of PPRL consists of five main criteria: **Computation and communication costs** determine the efficiency aspect that are often measured either theoretically

using the big-O notation [10] or empirically using runtime, memory size, number of communication steps, number and size of messages to be communicated, and number of comparisons required [44]. **Privacy guarantees** are either formally proven or empirically measured using metrics such as Information gain and disclosure risk metrics [41] against privacy attacks. **Correctness and fairness** correspond to the linkage quality aspect, where correctness is the accuracy of linkage results measured using precision, recall, area under curve (AUC), and F1-measure [6], and fairness is the accuracy of linkage results with regard to different subgroups of individuals [46].

3.4. Technologies

Privacy technologies: We categorize the key privacy technologies as:

1. **Cryptography** refers to secure multi-party computation techniques, such as homomorphic encryptions, secret sharing, and secure vector operations [25]. These techniques are provably secure and highly accurate, however, they are computationally expensive. An example PPRL technique based on cryptographic techniques is the secure edit distance algorithm for matching two strings or genome sequences [1], which is quadratic in the length of the strings.
2. **Embedding techniques** allow data to be mapped into a multi-dimensional metric space while preserving the distances between original data [32]. It is difficult to determine the appropriate dimensionality for the metric space. A recent work proposed a framework for embedding string and numerical data with theoretical guarantees for rigorous specification of space dimensionality [17].
3. **Differential privacy** is a rigorous definition that provides guarantees of indistinguishability of an individual regardless of the presence or absence of the individual’s record in the data with high probability. It has been used in PPRL to perturb data by adding noise such that every individual in the dataset is indistinguishable [41, 44]. However, adding noise incurs utility loss and volume increase. Output constrained differential privacy is a recently introduced privacy model for PPRL that allows disclosing matching records while being insensitive to the presence or absence of a single non-matching record [14].
4. **Statistical linkage key (SLK)** contains derived values from QIDs which is generated using a combination of components of a set of QIDs. An example is the SLK-581 consisting of the second and third letters of first name, the second, third and fifth letters of surname, full date of birth, and sex, which was developed by the Australian Institute of Health and Welfare to link records from the Home and Community Care datasets [31]. A recent study has shown that SLK-based techniques fall short in providing sufficient privacy protection and sensitivity [31].
5. **Probabilistic methods** are the most widely used techniques for practical PPRL applications due to their efficiency and controllable privacy-utility trade-off. These methods use probabilistic data structures for mapping/encoding data such that the actual distances between original data are preserved depending on the false positive probability of the mapping. A recent study has shown that, if effectively used, probabilistic techniques can achieve high linkage quality comparable to using unencoded records [30].

For example, Bloom filter encoding is one probabilistic method that has been used in several PPRL solutions [12, 30, 33, 34, 38, 39]. A Bloom filter b_i is a bit array of length l bits where all bits are initially set to 0. k independent hash functions, h_j , with $1 \leq j \leq k$, are used to map each of the elements s in a set S into the Bloom filter by setting

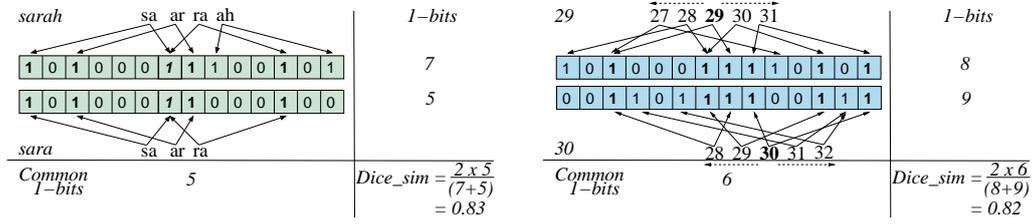


Figure 2: Bloom filter-based matching for string [33] (left) and numerical [40] (right) data.

the bit positions $h_j(s), 1 \leq j \leq k$ to 1. As shown in Figure 2, the set S of q -grams (sub-strings of length q) for string QIDs (left) or neighbouring values for numerical QIDs (right) can be hash-mapped into Bloom filters [33, 40]. The resulting Bloom filters can be matched using a similarity function, such as the Dice-coefficient [6] which is calculated as: $Dice_sim(b_1, \dots, b_p) = \frac{p \times c}{\sum_{i=1}^p x_i}$, where p is the number of Bloom filters compared, c is the number of common bit positions that are set to 1 in all p Bloom filters, and x_i is the number of bit positions set to 1 in $b_i, 1 \leq i \leq p$. The matching can either be done by a LU [12, 33] or collaboratively by the DOs [38, 39].

Complexity reduction technologies: The bottleneck of the PPRL process is the comparison of records across different databases using similarity functions, which is equal to the product of the sizes of the databases. Computational technologies have been used to improve the scalability of PPRL:

1. **Blocking** is defined on selected attributes (blocking keys) and it partitions the records in a database into several blocks or clusters based on the blocking key values such that comparison can be restricted to the records of the same block. A variety of blocking techniques has been developed [44]. Recent examples are randomized blocking methods based on Locality-Sensitive Hashing, which provide theoretical guarantees for identifying similar record pairs in the embedding space with high probability [12, 18].
2. **Meta-blocking** is the process of restructuring a collection of generated blocks to be compared in the next step such that unnecessary comparisons are pruned. Block processing for PPRL only received much attention in the recent years with the aim to improve the scalability in Big Data applications by employing such techniques along with blocking techniques [16, 28].
3. **Filtering** is an optimization technique for a particular similarity function to eliminate pairs/sets of records that cannot meet the similarity threshold for the selected similarity measure [34, 38].
4. **Parallel/distributed processing** for PPRL has only seen limited attention so far. Parallel linkage aims at improving the execution time proportionally to the number of processors [9, 18]. This can be achieved by partitioning the set of all record pairs to be compared, for example using blocking, and conducting the comparison of the different partitions in parallel on different processors.
5. **Advanced communication patterns** can be used to reduce the exponential growth of complexity for multi-party linkage. Such communication patterns include sequential, ring by ring, tree-based, and hierarchical patterns. Some of these patterns have been investigated for PPRL on multiple databases [42].

Linkage technologies: A variety of linkage methods and technologies has been used.

1. **Similarity functions** are required for fuzzy/approximate matching of QIDs in order to account for data errors and variations in the QIDs used for linkage. Different similarity functions have been used for different QID data types and different encoding/masking functions. For example, Bloom filter encoding-based PPRL requires token-based similarity functions, such as Jaccard, Hamming, or Dice coefficient functions (as described above) [33, 12]. Similarly, embedding techniques have used Euclidean or edit distance as similarity functions [32, 17].
2. **Matching** techniques determine how the linkage of records needs to be performed. It is a common practice to first de-duplicate records (internally link) within a single database before linking with records from other databases. This is known as one-to-one linking. If the databases are not de-duplicated (i.e. they contain multiple records corresponding to the same real-world entity), then many-to-many linking is required. Further, in multi-database linking, subset matching is required in certain applications to identify records that match across any subset of databases (for example, patients visited at least three out of five hospitals) [43].
3. **Classification** techniques, ranging from simple threshold-based, rule-based, to probabilistic linkage and machine learning, have been used in the PPRL literature [44]. The aim of these classifiers is to classify the record pairs into ‘matches’ or ‘non-matches’ based on the similarity between their QIDs. While machine learning-based classifiers can provide higher linkage quality, they require training data with ground-truth labels (of ‘matches’ and ‘non-matches’) for supervised techniques.
4. **Clustering** is an unsupervised technique that aims to group matching records corresponding to the same real-world entity into one cluster. A recent work studies incremental clustering techniques for multi-party PPRL [43].
5. **Fairness-aware linkage** is important to ensure fairness in linkage with respect to vulnerable sub-groups of the population. Errors in the linkage will propagate through to the subsequent data analysis. Fairness in the linkage process enables fair upstream data analysis. Fairness-aware classification and clustering algorithms have been developed in the literature [26], however, mitigating fairness-bias specifically in PPRL has not yet been studied.

4. Examples of application

Linking data is increasingly being required in a number of application areas [6]. When databases are linked within a single organization, then generally privacy and confidentiality are not of great concern. However, linking data from several organizations imposes legal and ethical constraints on using personal data for the linkage, as described by the Australian Data-Matching Program Act¹, the EU General Data Protection Regulation², and the HIPAA Act in the USA³. PPRL is therefore required in several real applications, as the following examples illustrate:

¹<https://www.legislation.gov.au/Details/C2016C00755>

²<https://gdpr.eu/>

³<http://www.hhs.gov/ocr/privacy/>

1. **Healthcare applications:** Several healthcare applications ranging from health surveillance, epidemiological studies, and clinical trials, to public health research require PPRL as an efficient building block. For example, a study on surgical treatment received by aboriginal and non-aboriginal people with lung cancer linked data from hospitals and clinical registries with data from central cancer registries and from the Australian Bureau of Statistics using PPRL techniques [8]. Bloom filter-based PPRL was used to link data from several cantonal and national registries in Switzerland to investigate long-term consequences of childhood cancer [20]. In 2016, the Interdisciplinary Committee of the International Rare Diseases Research Consortium launched a task team to explore approaches to PPRL for linking several genomic and clinical data sets [2].
2. **Government services and research:** The traditionally used small-scale survey studies have been replaced by linking databases to develop policies in a more efficient and effective way [19]. The research program ‘Beyond 2011’ established by the Office for National Statistics in the UK, for example aimed to study the options for production of population statistics for England and Wales, by linking anonymous data [35]. Social scientists use PPRL in the field of population informatics to study insights into our society [21].
3. **Business collaboration:** Many businesses take advantage of linking data across different organizations, such as suppliers, retailers, wholesalers, and advertisers, for improving efficiency, targeted marketing, and reducing costs of their supply chains. PPRL can be used for cross-organizational collaborative decision making which involves a great deal of private information that businesses are often reluctant to disclose [47].
4. **National security applications:** PPRL techniques are being used by national security agencies and crime investigators to identify individuals who have committed fraud or crimes [27]. These applications integrate data from various sources, such as law enforcement agencies, Internet service providers, the police, tax agencies, government services, as well as financial institutions to enable the accurate identification of crime and fraud, or of terrorism suspects.

5. Future directions for research

In this section we describe the various open challenges and research directions of PPRL for Big Data applications, as categorized in Figure 3.

5.1. Scalability

The *volume* of data increases both due to the large size of databases and the number of different databases. Challenges of linking large databases have been addressed by using computational techniques to reduce the number of required comparisons between records. However, these techniques do not solve the scalability problem for Big Data completely. For example, even small blocks of records resulting from a blocking technique can still lead to a large number of comparisons between records with the increasing volume of data. Moreover, only limited work has been done to address the challenges of linking multiple databases.

Scalable blocking techniques are required that can generate blocks of suitable sizes across multiple databases, as are advanced filtering techniques that effectively prune potential non-matches even further. With multiple databases, identifying subsets of records that match across only a subset of databases (for example, patients in three out of five hospital databases) is even more challenging due to the large number of combinations of subsets. Advanced communication

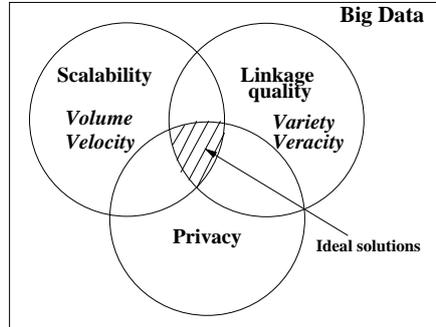


Figure 3: Challenges of PPRL for Big Data.

patterns, distributed computing, and adaptive comparison techniques are required towards this direction for scalable PPRL applications. These techniques are largely orthogonal so that they can be combined to achieve maximal efficiency.

Another major aspect of Big Data is the *velocity*, i.e. the dynamic nature of data. Current techniques are only applicable in batch-mode on static and well defined relational data. Required are approaches for dynamic data, real-time integration, and data streams, for adaptive systems to link data as they arrive at an organization, ideally in (near) real-time.

5.2. Linkage Quality

Big Data are often complex, noisy, and erroneous, which refer to the *veracity* and *variety* aspects. The linkage can be affected by the quality of data, especially in a privacy-preserving context where linkage is done on the masked or encoded version of the data. Masking data reduces the quality of data to account for privacy and hence the quality of data would have an adverse effect on the whole PPRL process leading to low linkage quality. The trade-off between quality and privacy needs to be handled carefully for different privacy masking functions. The affect of poor data quality on the linkage quality becomes worse with increasing number of databases. More advanced classification techniques, such as collective [4] and graph-based [15] techniques, need to be investigated to address the linkage quality problem of PPRL.

Developing classifiers that are fair with respect to a protected/sensitive feature [46], such as gender or race, is an important problem for classification in general and specifically for record linkage. Fairness of a classifier with regard to a certain protected feature determines how much the classifier distorts from producing correct predictions with equal probabilities for individuals across different protected groups/values. There has been increased interest in this field due to the concerns that classifiers may introduce significant bias towards certain minority or vulnerable group with regard to the protected feature, such as race or gender, for example against black people in fraud and crime detection systems [13, 24] or against women in job recommendation systems [11]. However, fairness has not been studied specifically for PPRL so far.

Assessing the linkage quality in a PPRL project is very challenging because it is generally not possible to inspect linked records due to privacy concerns. Knowing the quality of linkage is crucial in many Big Data applications such as in the health or security domains. An initial work has been done on interactive PPRL [22] where parts of sensitive values are iteratively revealed for manual assessment in such a way that the privacy compromise is limited. Implementing such ap-

proaches in real applications is an open challenge that must be solved. Using heuristic measures to approximately evaluate the linkage quality is another option that requires more research.

5.3. Privacy

Another open challenge in PPRL is how resistant the techniques are against different adversarial attacks. Most work in PPRL assume the *semi-honest* adversary model [25] and the trusted LU-based model. Furthermore, these works assume that the parties do not collude with each other [44]. Only few PPRL techniques consider the *malicious* adversary model as it imposes a high complexity [44]. More research is therefore required to develop novel security models that lie between these two models and prevent against collusion risks for PPRL.

Sophisticated attack methods [7, 23] have been recently developed that exploit the information revealed during the PPRL protocols to iteratively gather information about sensitive values. Therefore, existing PPRL techniques need to be hardened to ensure they are not vulnerable to such attacks [33]. Evaluation of PPRL techniques is challenged due to the absence of benchmarks, datasets, and frameworks. Different measurements have been used [12, 41, 44], making the comparison of different PPRL techniques difficult. Synthetic datasets generated with real data characteristics using data generators [37] have been used as an alternative to benchmark datasets. This limits the evaluation of PPRL techniques to assess their application in the real setting.

Further, there has not been much interaction between practitioners and researchers of PPRL to allow better understanding of the whole data life cycle and to evaluate the applicability of PPRL in real applications. A comprehensive privacy strategy is essential including a closely aligned PPRL and privacy-preserving data technologies for strategic Big Data applications.

References

- [1] Atallah, M., Kerschbaum, F., Du, W., 2003. Secure and private sequence comparisons, in: ACM WPES, pp. 39–44.
- [2] Baker, D., Knoppers, B.M., Phillips, M., van Enckevort, D., Kaufmann, P., Lochmuller, H., Taruscio, D., 2018. Privacy-preserving linkage of genomic and clinical data sets. *IEEE Transactions on Computational Biology and Bioinformatics*.
- [3] Bergstra, J., Bengio, Y., 2012. Random search for hyper-parameter optimization. *Journal of machine learning research* 13, 281–305.
- [4] Bhattacharya, I., Getoor, L., 2007. Collective entity resolution in relational data. *TKDD* 1.
- [5] Boyd, J., Randall, S., Ferrante, A., 2015. Application of privacy-preserving techniques in operational record linkage centres, in: *Medical Data Privacy Handbook*, pp. 267–287.
- [6] Christen, P., 2012. *Data matching. Data-Centric Systems and Applications*, Springer.
- [7] Christen, P., Schnell, R., Vatsalan, D., Ranbaduge, T., 2017. Efficient cryptanalysis of bloom filters for PPRL, in: *PAKDD*, Springer. pp. 628–640.
- [8] Condon, J.R., Barnes, T., Cunningham, J., Armstrong, B.K., 2004. Long-term trends in cancer mortality for indigenous australians in the northern territory. *Medical Journal of Australia* 180, 504.
- [9] Dal Bianco, G., Galante, R., Heuser, C.A., 2011. A fast approach for parallel deduplication on multicore processors, in: *SAC*, ACM. pp. 1027–1032.
- [10] Danziger, P., 2010. Big o notation. Source internet: <http://www.scs.ryerson.ca/mth110/Handouts/PD/bigO.pdf>.
- [11] Datta, A., Tschantz, M.C., Datta, A., 2015. Automated experiments on ad privacy settings. *Proceedings on privacy enhancing technologies* 2015, 92–112.
- [12] Durham, E.A., 2012. *A framework for accurate, efficient private record linkage*. Ph.D. thesis. Vanderbilt University, Nashville, TN.
- [13] Flores, A.W., Bechtel, K., Lowenkamp, C.T., 2016. False positives, false negatives, and false analyses: A rejoinder to machine bias: There’s software used across the country to predict future criminals. and it’s biased against blacks. *Fed. Probation* 80, 38.
- [14] He, X., Machanavajjhala, A., Flynn, C.J., Srivastava, D., 2017. Scaling private record linkage using output constrained differential privacy. *CoRR abs/1702.00535*.

- [15] Kalashnikov, D., Mehrotra, S., 2006. Domain-independent data cleaning via analysis of entity-relationship graph. *TODS* 31, 716–767.
- [16] Karakasidis, A., Koloniari, G., Verykios, V.S., 2015. Scalable blocking for PPRL, in: *SIGKDD*, ACM. pp. 527–536.
- [17] Karapiperis, D., Gkoulalas-Divanis, A., Verykios, V.S., 2017. Federal: A framework for distance-aware privacy-preserving record linkage. *TKDE* .
- [18] Karapiperis, D., Verykios, V., 2015. An LSH-based Blocking Approach with a Homomorphic Matching Technique for PPRL. *TKDE* 27, 909–921.
- [19] Kelman, C.W., Bass, J., Holman, D., 2002. Research use of linked health data – A best practice protocol. *ANZJPH* 26, 251–255.
- [20] Kuehni, C.E., Rueegg, C.S., Michel, G., Rebholz, C.E., Strippoli, M.P.F., Niggli, F.K., Egger, M., von der Weid, N.X., (SPOG), S.P.O.G., 2011. Cohort profile: the Swiss childhood cancer survivor study. *International journal of epidemiology* 41, 1553–1564.
- [21] Kum, H., Krishnamurthy, A., Machanavajjhala, A., Ahalt, S., 2013. Population informatics: Tapping the social genome to advance society: A vision for putting” big data” to work for population informatics. *Computer* .
- [22] Kum, H.C., Krishnamurthy, A., Machanavajjhala, A., Reiter, M.K., Ahalt, S., 2014. Privacy preserving interactive record linkage. *JAMIA* 21, 212–220.
- [23] Kuzu, M., Kantarcioglu, M., Durham, E., Malin, B., 2011. A constraint satisfaction cryptanalysis of Bloom filters in private record linkage, in: *PETS*, Springer LNCS, Waterloo, Canada. pp. 226–245.
- [24] Larson, J., Mattu, S., Kirchner, L., Angwin, J., 2016. How we analyzed the compas recidivism algorithm. *ProPublica* (5 2016) 9.
- [25] Lindell, Y., Pinkas, B., 2009. Secure multiparty computation for privacy-preserving data mining. *JPC* 1.
- [26] Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A., 2019. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635* .
- [27] Phua, C., Smith-Miles, K., Lee, V., Gayler, R., 2012. Resilient identity crime detection. *IEEE TKDE* 24.
- [28] Ranbaduge, T., Vatsalan, D., Christen, P., 2016. Scalable block scheduling for efficient multi-database record linkage, in: *ICDM*, IEEE. pp. 1161–1166.
- [29] Ranbaduge, T., Vatsalan, D., Christen, P., 2020. Secure multi-party summation protocols: Are they secure enough under collusion? *Transactions on Data Privacy* 13, 25–60.
- [30] Randall, S.M., Ferrante, A.M., Boyd, J.H., Bauer, J.K., Semmens, J.B., 2014. PPRL on large real world datasets. *JBI* 50, 205–212.
- [31] Randall, S.M., Ferrante, A.M., Boyd, J.H., Brown, A.P., Semmens, J.B., 2016. Limited privacy protection and poor sensitivity: Is it time to move on from the statistical linkage key-581? *HIMJ* 45, 71–79.
- [32] Scannapieco, M., Figotin, I., Bertino, E., Elmagarmid, A., 2007. Privacy preserving schema and data matching, in: *ACM SIGMOD*, pp. 653–664.
- [33] Schnell, R., 2015. Privacy-preserving record linkage, in: *Methodological Developments in Data Linkage*, pp. 201–225.
- [34] Sehili, Z., Kolb, L., Borgs, C., Schnell, R., Rahm, E., 2015. PPRL with PPJoin, in: *BTW*, Hamburg. pp. 85–104.
- [35] Skinner, C., Hollis, J., Murphy, M., 2013. ”matching anonymous data”. *Beyond 2011: Independent review of methodology* , 1–48.
- [36] Snoek, J., Larochelle, H., Adams, R.P., 2012. Practical bayesian optimization of machine learning algorithms, in: *Advances in neural information processing systems*, pp. 2951–2959.
- [37] Tran, K.N., Vatsalan, D., Christen, P., 2013. GeCo: an online personal data generator and corruptor, in: *CIKM*, ACM, San Francisco. pp. 2473–2476.
- [38] Vatsalan, D., Christen, P., 2012. An iterative two-party protocol for scalable PPRL, in: *AusDM, CRPIT*, Sydney. pp. 1–12.
- [39] Vatsalan, D., Christen, P., 2014. Scalable PPRL for multiple databases, in: *CIKM*, ACM, Shanghai. pp. 1795–1798.
- [40] Vatsalan, D., Christen, P., 2016. Privacy-preserving matching of similar patients. *JBI* 59, 285–298.
- [41] Vatsalan, D., Christen, P., O’Keefe, C.M., Verykios, V.S., 2014. An evaluation framework for PPRL. *JPC* 6.
- [42] Vatsalan, D., Christen, P., Rahm, E., 2016. Scalable privacy-preserving linking of multiple databases using counting bloom filters, in: *ICDMW PDDM*, IEEE, Barcelona. pp. 882–889.
- [43] Vatsalan, D., Christen, P., Rahm, E., 2020. Incremental clustering techniques for multi-party privacy-preserving record linkage. *Data & Knowledge Engineering* , 101809.
- [44] Vatsalan, D., Christen, P., Verykios, V.S., 2013. A taxonomy of PPRL techniques. *JIS* 38, 946–969.
- [45] Vatsalan, D., Sehili, Z., Christen, P., Rahm, E., 2017. Privacy-preserving record linkage for big data: Current approaches and research challenges. *Handbook of big data technologies* , 851–895.
- [46] Zafar, M.B., Valera, I., Rodriguez, M.G., Gummadi, K.P., 2017. Fairness constraints: Mechanisms for fair classification, in: *International Conference on Artificial Intelligence and Statistics (AISTATS)*, Florida, USA. pp. 962–970.
- [47] Zhu, H., Liu, H., Ou, C.X., Davison, R.M., Yang, Z., 2017. Privacy preserving mechanisms for optimizing cross-

organizational collaborative decisions based on the karmarkar algorithm. JIS 72, 205–217.