# Prevention of data from re-identifiers using map reduce

**Mr. Kolekar Jairam , Mr. Pande Vikas  , Mr. Nitin Khandagale  , Prof. Archana Gaikwad**
**Department of Computer Engineering, D.Y.Patil School of Engineering**

**Abstract — Many data householders area unit required to unleash the data throughout a kind of world application, since it's of necessary importance to discovery valuable knowledge keep behind the data. However, existing re-identification attacks on the AOL and ADULTS data sets have shown that publish such knowledge directly may cause tremendous threads to the individual privacy. so, it's very important to resolve all types of re-identification risks by recommend efficient de-identification policy to confirm every privacy and utility of the data De-identification policies is one in all the models which is able to be used to succeed such desires, however, the amount of de-identification policies is exponentially huge due to the broad domain of quasi-identifier attributes. to higher management the trade off between data utility and data privacy, skyline computation are used to opt for such policies, but it's notwithstanding tough for economical skyline method over sizable quantity of policies. throughout this paper, we tend to tend to propose one parallel algorithmic program called SKY-FILTER-MR, that depends on Map cut back to beat this challenge by computing skylines over huge scale de-identification policies that is drawn by bit-strings. To further improve the performance, a totally distinctive approximate skyline computation theme was projected to prune unqualified policies exploitation the concerning domination relationship. With estimated skyline, the power of filtering at intervals the policy space generation stage was greatly sturdy to effectively decrease the worth of skyline calculation over many policies. full experiments over every universe and artificial datasets demonstrate that our projected SKY-FILTER-MR algorithmic program well outperforms the baseline approach by up to fourfold faster at intervals the most effective case, that indicates wise quality over huge policy sets.**

*Index Terms— Big Data; Access Control; Privacy-preserving Policy; De-identification policies.*

## I. INTRODUCTION

Big information could be a term that refers to information sets or combos of information sets whose size (volume), complexness (variability), and rate of growth (velocity) create them tough to be captured, managed, processed or analyzed by standard technologies and tools, like relative databases and desktop statistics or image packages, inside the time necessary to form them helpful. whereas the dimensions wont to verify whether or not a specific information set is taken into account massive information isn't firmly outlined and continues to vary over time, most analysts and practitioners presently check with information sets from 30-50 terabytes(10 twelve or a thousand gigabytes per terabyte) to multiple petabytes (1015 or a thousand terabytes per petabyte) as massive information.

The analysis of huge information involves multiple distinct phases as shown within the figure below, every of that introduces challenges. many of us sadly focus simply on the analysis/modeling section: whereas that phase is crucial, it's of very little use while not the opposite phases of the information analysis pipeline. Even within the analysis section, that has received a lot of attention, there square measure poorly understood complexities within the context of multi-tenanted clusters wherever many users' programs run at the same time. several vital challenges extend on the far side the analysis section. for instance, massive information should be managed in context, which can be strident, heterogeneous ANd not embody an direct model. Doing therefore raises the requirement to trace cradle and to handle uncertainty and error: topics that square measure crucial to success, and nonetheless seldom mentioned within the same breath as massive information. Similarly, the inquiries to the information analysis pipeline can generally not all be arranged  go into advance. it should ought to discover smart queries supported the information. Doing this can need smarter systems and additionally higher support for user interaction with the analysis pipeline. In fact, there's a serious bottleneck within the variety of individuals authorized  to raise queries of the information and analyze it. It will drastically increase this variety massive information could also be a term that refers to information sets or mixtures of information sets whose size (volume), quality (variability), and rate of growth (velocity) build them hard to be captured, managed, processed or analyzed by normal technologies and tools, like relative databases and desktop statistics or image packages, among the time necessary to make them useful. whereas the size accustomed verify whether or not or not a specific information set is taken into consideration Brobdingnagian information is not firmly made public and continues to vary over time, most analysts and practitioners presently consult with information sets from 30-50 terabytes(10 twelve or one thousand gigabytes per terabyte) to multiple petabytes (1015 or one thousand terabytes per petabyte) as Brobdingnagian information.

## II.LITERATURE SURVEY

### 1) Privacy-Preserving Data Publishing: A Survey of Recent Developments

**Authors**: BENJAMIN C. M. FUNG, KE WANG, RUI CHEN, PHILIP S. YU.

The collection of digital data by governments, companies, and other people has created tremendous opportunities for knowledge- and information-based higher cognitive operation. Driven by mutual benefits, or by laws that require certain data to be written, there is a demand for the exchange and publication of information among varied parties. data in its original kind, however, sometimes contains sensitive data relating to folks, and industrial enterprise such data will violate individual privacy. this apply in data industrial enterprise depends primarily on policies and tips about what kinds of data are written and on agreements on the use of written data. This approach alone may cause excessive data distortion or short protection. Privacy-preserving data industrial enterprise (PPDP) provides ways that and tools for industrial enterprise useful data whereas protecting data privacy.

### 2) APPLET: a privacy-preserving framework for location-aware recommender system

**Authors:** Xindi Ma,HuiLI, Jianfeng MA, Qi JIANG, Sheng GAO, Ning XI &DiLU

Location-aware recommender systems that use location-based ratings to supply recommendations have recently intimate a speedy development and draw important attention from the analysis community. However, current work in the main focused on high-quality recommendations whereas underestimating privacy issues, which can cause problems with privacy. Such problems ar lots of distinguished once service suppliers, World Health Organization have restricted machine and storage resources, leverage on cloud platforms to suit in with the tremendous range of service requirements and users. throughout this paper, we've an inclination to propose a singular framework, specifically applications applied scientist, for shielding user privacy data, beside locations and recommendation results, within a cloud surroundings.

### 3) Efficient Discovery of De-identification Policies Through a Risk-Utility Frontier

**Authors:** Weiyi Xia, Raymond Heatherly, Xiaofeng Ding

Modern information technologies modify organizations to capture big quantities of person-specific information whereas providing routine services. many organizations hope, or unit First State jure required, to share such information for secondary functions (e.g.validation of study findings) throughout a de-identified manner. In previous work, it had been shown de-identification policy alternatives may be modeled on a lattice, which may be explore for policies that met a prespecified risk threshold (e.g., likelihood of re-identification). However, the search was restricted in many ways that. First, its definition of utility was syntactic supported the extent of the lattice - and not linguistics - primarily based

on the actual changes induced at intervals the following information. Second, the sting may not be noted beforehand. The goal of this work is to make the optimum set of policies that trade-off between privacy risk (R) and utility (U), that we've got an inclination to raise as a R-U frontier. To model this downside, we've got an inclination to in-troduce a linguistics definition of utility, supported theory, that is compatible with the lattice illustration of policies. to unravel the matter, we've got an inclination to initially build a bunch of policies that define a frontier. we've got an inclination to then use a probability-guided heuristic to travel wanting the lattice for policies attainable to update the frontier. To demonstrate the effectiveness of our approach, we've got an inclination to perform degree empirical analysis with the Adult dataset of the UCI Machine Learning Repository.

### 4)l -Diversity: Privacy Beyond k-Anonymity

**Authors:** Ashwin Machanavajjhala, Johannes Gehrke, Daniel Kifer

Publishing knowledge regarding people while not revealing sensitive data regarding them is a vital downside. In recent years, a replacement definition of privacy referred to as k-anonymity has gained quality. during a k-anonymized dataset, every record is indistinguishable from a minimum of $k-1$ different records with relation to bound "identifying" attributes. during this paper we have a tendency to show with 2 easy attacks that a k-anonymized dataset has some refined, however severe privacy problems. First, we have a tendency to show that AN offender will discover the values of sensitive attributes once there's very little diversity in those sensitive attributes. Second, attackers typically have background, and that we show that k-anonymity doesn't guarantee privacy against attackers mistreatment background. we have a tendency to provides a careful analysis of those 2 at- tacks and that we propose a unique and powerful privacy definition referred to as -diversity. additionally to assembling a proper foundation for-diversity, we have a tendency to show in AN experimental analysis that -diversity is sensible and may be enforced with efficiency.

### 5) k-ANONYMITY: A MODEL FOR PROTECTING PRIVACY

**Authors:** LATANYA SWEENEY

Consider a knowledge holder, like a hospital or a bank, that encompasses a in camera control assortment of person-specific, field structured knowledge. Suppose the information holder desires to share a version of the information with researchers. however will a knowledge holder unharness a version of its personal knowledge with scientific guarantees that the people United Nations agency square measure the themes of the data can not be re-identified whereas the information stay much useful? the answer provided during this paper includes a proper protection model named face-anonymity and a collection of related policies for preparation. A unharness provides fc-anonymity protection if infofor every person contained within the unharness can not be distinguished from a minimum of k-\ people whose information conjointly seems within the unharness. This paper conjointly examines re-identification attacks which will be accomplished on releases that adhere to k- namelessness unless related policies square measure revered. The fc-anonymity protection model is vital as a result of it forms the premise on that the real-world systems called Datafly, |i-Argus and fc-Similar offer guarantees of privacy protection

## III.EXISTING SYSTEM

Existing re-identification attacks on the AOL and ADULTS information sets have shown that publish such knowledge directly would possibly cause tremendous threads to the individual privacy. Thus, it's pressing to resolve every type of re-identification risks by recommending effective de-identification policies to make sure every privacy and utility of the information. Their work has limitations in many ways that. First, their framework desires a lattice that contains all the selection policies to rearrange with note value. Second, their algorithms area unit approximate approaches that don't have any guarantee of best resolution.

## IV.PROPOSED SYSTEM

In this paper, we tend to propose one parallel rule known as SKY-FILTER-MR that is predicated on Map scale back to beat this challenge by computing skylines over massive scale de-identification policies that's painted by bit-strings. To additional improve the performance, a completely unique approximate skyline computation theme was planned to prune unqualified policies victimization the some domination relationship. With approximate skyline, the ability of filtering within the policy house generation stage was greatly strong to effectively decrease the value of skyline computation over different policies. in depth experiments over each reality and artificial datasets demonstrate that our planned SKY-FILTER-MR rule considerably outperforms the baseline approach by up to fourfold quicker within the optimum case, that indicates sensible measurability over massive policy sets. Our contribution is to see the higher solutions. Keep the most effective solutions, and use them to come up with new doable solutions victimization genetic rule.
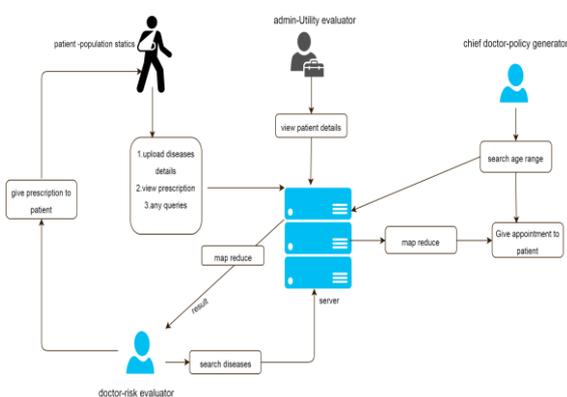
## V.SYSTEM ARCHITECTURE



Fig.1: System Architecture

In our system there are four modules in first multiple user upload our diseases on server then chief doctor will search patient age wise .system perform map reduce operation on patient information and give result to chief doctor then chief doctor will give appointment to patient .then third module is doctor .doctor will search diseases .system will perform map reduce operation on patient diseases and give

result to doctor .then doctor will give prescription to patient. In last module admin will only view all details of patient details

## VI. MODULES

1. **Patient Module:**
   a. Registration
   b. Login
      1. Enter Symptoms & Then Submit.
      2. View Prescription-Give By Doctor
      3. View Appointment –Give By Chief Doctor
      4. Logout.

**Doctor Module:**
   1. Registration- by Specialization
   2. Login
      1. Efficient Recommendation of Patients based on Symptoms & Specialization of Doctor Disease
      2. Search Patient by age, Gender etc.
      3. Give Prescription of medicine to Patient.
      4. Logout.

   **Chief Doctor Module:**
   1. Login
      1. Efficient Recommendation of Patient that have queries & not Performed well prescription from doctor so it will recommend that patient to Chief Doctor by Disease Specialization.
      2. Search Patient By Disease Subtype.
      3. Give Appointment to that Patient for doctor
      4. Logout

   **Admin module**
   1. Login

      1. View All Patient Details

      2. View All Patients Status Details

      3. All Chief Doctor Details

      4. View All Doctor Details

      5. All Symptoms Disease Details

      6. Add Specialist & Disease

      7. Add Symptoms & Key

8. Add Disease & Symptoms

9. Logout.

## VII. GOALS AND   OBJECTIVE

- The main goal of the project is to study, design and implement performance optimizations for big data frameworks.
- This work contributes methods and techniques to build tools for easy and efficient processing of very large data sets. It describes ways to make systems faster, by inventing ways to shorten job completion times.
- To generate faster results.
- It reduces the complexity of data access and retrieval. When we have to dealing with big data.
- The alternative to this is apache Hadoop, which deals with big data with efficiency.
- Hadoop  itself consists of Map Reduce and HDFS.
- Provide security to personal information.
- Protect the user data during transmission.
- We perform a detailed security analysis and performance evaluation of the proposed technique.

## VIII. RESULT



Screenshot 1



Screenshot 2



Screenshot 3



Screenshot 4



Screenshot 5

## IX. CONCLUSION AND FUTURE SCOPE

We study the advice on an excellent variety of Delaware identification policies victimization Map scale back. Firstly, we tend to imply an efficient method of policy generation on the idea of recently projected definition, which may decreases the time of generating policies and therefore the size of other policy set dramatically. Secondly, we tend to propose SKY-FILTER-MR, that may be a three-round Map Reduce-based parallel rule, to answer skyline de-identification policies with efficiency. we tend to use bit-strings to represent one policy within the framework. so as to any improve the performance, a lively approximate skyline theme is projected to decrease the amount of other policy set. By group action the approximate skyline with the minimal  Map Reduce rule, the filtering power within the Map section of initial spherical was optimized while not increasing the transmission price. we tend to perform comprehensive experimental analysis on each real-world and artificial datasets, and therefore the results indicate sensible performance and quantifiability of our projected SKY-FILTER-MR.

### ACKNOWLEDGMENT

### REFERENCES

[1] B. C. M. Fung, K.Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing: A survey of recent

developments," *ACM Comput. Surv.*, vol. 42, no. 4, pp. 14:1–14:53, 2010.

[2] X. MA, H. Li, J. Ma, Q. Jiang, S. Gao, N. Xi, and D. Lu, "Applet: A privacy-preserving framework for location-aware recommender system," *Sci China Inf Sci*, vol. 59, no. 2, pp. 1–15, 2016.

[3] W. Xia, R. Heatherly, X. Ding, J. Li, and B. Malin, "Efficient discovery of de-identification policies through a risk-utility frontier," in *CODASPY*, 2013, pp. 59–70.

[4] K. Benitez, G. Loukides, and B. Malin, "Beyond safe harbor: Automatic discovery of health information de-identification policy alternatives," in *IHI*, 2010, pp. 163–172.

[5] K. E. Emam, "Heuristics for de-identifying health data," *IEEE Security and Privacy*, vol. 6, no. 4, pp. 58–61, 2008.

[6] L. Sweeney, "$k$-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 555–570, 2002.

[7] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam, "$\ell$-diversity: Privacy beyond $k$-anonymity," in *TKDD*, 2007, pp. 1–52.

[8] N. Li, T. Li, and S. Venkatasubramanian, "$t$-closeness: Privacybeyond $k$-anonymity and $\ell$-diversity," in *ICDE*, 2007, pp. 106–115.

[9] J. Brickell and V. Shmatikov, "The cost of privacy: Destruction of data-mining utility in anonymized data publishing," in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008, pp. 70–78.

[10] J. Cao and P. Karras, "Publishing microdata with a robust privacy guarantee," *Proc. VLDB Endow.*, vol. 5, no. 11, pp. 1388–1399, 2012.

[11] W. Xia, R. Heatherly, X. Ding, J. Li, and B. A. Malin, "Ru policy frontiers for health data de-identification," *Journal of the American Medical Informatics Association*, vol. 22, no. 5, pp. 1029–1041, 2015.