Chapter 60 Privacy Preserving Big Data Publishing: Challenges, Techniques, and Architectures

Nancy Victor VIT University, India

Daphne Lopez

VIT University, India

ABSTRACT

Data privacy plays a noteworthy part in today's digital world where information is gathered at exceptional rates from different sources. Privacy preserving data publishing refers to the process of publishing personal data without questioning the privacy of individuals in any manner. A variety of approaches have been devised to forfend consumer privacy by applying traditional anonymization mechanisms. But these mechanisms are not well suited for Big Data, as the data which is generated nowadays is not just structured in manner. The data which is generated at very high velocities from various sources includes unstructured and semi-structured information, and thus becomes very difficult to process using traditional mechanisms. This chapter focuses on the various challenges with Big Data, PPDM and PPDP techniques for Big Data and how well it can be scaled for processing both historical and real-time data together using Lambda architecture. A distributed framework for privacy preservation in Big Data by combining Natural language processing techniques is also proposed in this chapter.

INTRODUCTION

"Data is the new oil", declared Clive Humby, a Sheffield mathematician ('Tech giants may be huge, but nothing matches big data', 2013). Michael Palmer expanded the quote as: "Data is just like crude. It's valuable, but if unrefined it cannot really be used. It has to be changed into gas, plastic, chemicals etc. to create a valuable entity that drives profitable activity; so must data be broken down, analyzed for it to have value". This is true in the case of Big Data. Data is the natural resource growing bigger and bigger

DOI: 10.4018/978-1-7998-8954-0.ch060

each and every second. Big Data is so large amount of data that cannot be processed using traditional systems. Big Data analytics is the process of generalizing values from large data sets through which hidden patterns, unknown correlations and other useful information can be uncovered ('The state of the enterprise cloud and prepping for AWS re:Invent 2013').

The main characteristics of Big Data (4 V's) are: Volume, Velocity, Variety and Veracity (The Four V's of Big Data, 2015).

- Volume: The word "big" in Big Data defines the volume. The various sources of Big data include sensors, social media, activity generated data, data warehouse appliances, archives, business apps etc. (The Big 9 big data sources, 2014; Top 10 categories for big data sources and mining technologies, 2012).
- Velocity: This refers to the speed at which the data flows in and out of the system. Some of the examples for data generation points include mobile devices, microphones, sensors, social media etc.
- Variety: Big Data includes structured, semi-structured and unstructured data, which is being produced from various sources.
- Veracity: It refers to the inconsistencies and incompleteness in data which is collected from various sources.

In order to derive value out of this massive data, it should be collected and processed efficiently. This itself brings in a lot of challenges which includes preserving the privacy of data that is collected from various data sources at very high rates, in a variety of data formats. For processing and managing Big data, various technologies are used in the Hadoop ecosystem. This includes HDFS for storage and replication, MapReduce for distributed processing, Mahout for machine learning, Pig for scripting and so on (Khan, N et al., 2014). Data publishing plays a major role in the case of Big data as the data which is collected can be publicized for use or reuse by researchers in order to obtain valuable research output. The data can then be used for performing various data mining tasks, which helps to gain better insights about the data which is collected.

The chapter is organized as follows: Section II gives an outline about the various challenges with Big Data. Section III explains the various models used for privacy preservation. Section IV focuses on privacy preserving data mining, whereas section V discusses about privacy preserving data publishing. An architecture for privacy preserving data publishing has been proposed in section VI and Section VII concludes the chapter.

CHALLENGES WITH BIG DATA

Wu et al. (2014) has considered the three tiers of a Big Data processing framework to explain the various challenges faced in the Big data domain. The three tiers include big data mining platform, big data semantics and application knowledge and big data mining algorithms. Big data mining platform mainly focuses on the difficulty in accessing, processing and computing the massive quantity of data that is clearly impossible with the existing computing facilities. MapReduce computation model plays a major role here in efficiently handling the pool of data which is being gathered from various sources at very high velocities. Big data semantics and application knowledge deals with the privacy concerns and distributing data to the public. Data privacy and location privacy is given at most importance in this case. The third 16 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/privacy-preserving-big-data-publishing/280229

Related Content

Teaching Case for Addressing Risks with Strategies in an International Airport Project

Daly Paulose (2013). International Journal of Risk and Contingency Management (pp. 18-35). www.irma-international.org/article/teaching-case-addressing-risks-strategies/76655

Risk and Models of Innovation Hubs: MIT and Fraunhofer Society

Mohammad Baydoun (2015). *International Journal of Risk and Contingency Management (pp. 17-26).* www.irma-international.org/article/risk-and-models-of-innovation-hubs/145363

Investigation of Credit Risk based on Indian Firm Performance

Manoj Kumar (2017). International Journal of Risk and Contingency Management (pp. 35-46). www.irma-international.org/article/investigation-of-credit-risk-based-on-indian-firm-performance/177839

IT Security Risk Management Model for Handling IT-Related Security Incidents: The Need for a New Escalation Approach

Gunnar Wahlgrenand Stewart James Kowalski (2018). Security and Privacy Management, Techniques, and Protocols (pp. 129-151).

www.irma-international.org/chapter/it-security-risk-management-model-for-handling-it-related-security-incidents/202042

Effective Recognition of Stereo Image Concealed Media of Interpolation Error with Difference Expansion

Hemalatha J.and Kavitha Devi M. K. (2016). *Combating Security Breaches and Criminal Activity in the Digital Sphere (pp. 157-165).*

www.irma-international.org/chapter/effective-recognition-of-stereo-image-concealed-media-of-interpolation-error-withdifference-expansion/156458