

Security Monitoring System of IOT Networks Using Big Data

A.Satheesh¹, R.Santhiya²

¹Assistant Professor, Dept of CSE

²Dept of CSE

^{1,2}Sembodai Rukmani Varatharajan Engineering College

Abstract- Networks of the Internet of Things (IoT) nowadays find greater widespread in many domains. Particularities of creation of IoT make the problem of their security monitoring rather actual; it is caused by necessity of processing of big amounts of heterogeneous data in real time. The problem may be solved by means of implementation of the parallel system for security data processing within IoT on the fly basing on complex event processing (CEP) technology. Considers basic solutions for creating such a system. The proposed system is oriented for usage of software environment Hadoop and includes data collection, data storage, data normalization and analysis, and data visualization components. Discusses the issues of architecture of this system, its implementation and experimental estimation.

I. INTRODUCTION

Graph mining has gained major attention for a broad range of applications, such as social networks, protein-protein interaction networks, road networks, etc.

Probabilistic Graphical Models

Uncertainty is inescapable in real-world applications: we are able to nearly never predict with certainty what's going to happen within the future, and even within the gift and therefore the past, several necessary aspects of the planet aren't determined with certainty. Applied mathematics offers United States the fundamental foundation to model our beliefs concerning the various potential states of the planet, and to update these beliefs as new proof is obtained. These beliefs are often combined with individual preferences to assist guide our actions, and even in choosing that observations to create. Whereas applied mathematics has existed since the seventeenth century, our ability to use it effectively on massive issues involving several inter-related variables is fairly recent, and is due mostly to the event of a framework called Probabilistic Graphical Models (PGMs)[1][2]. This framework, that spans strategies reminiscent of theorem networks and Mark off random fields, uses ideas from distinct information structures in computing to with efficiency write in code and manipulates likelihood distributions over high-

dimensional areas, typically involving lots of or maybe several thousands of variables. These strategies are employed in a vast vary of application domains, that include: internet search, medical and fault designation, image understanding, reconstruction of biological networks, speech recognition, language process, decryption of messages sent over a loud line, mechanism navigation, and lots of a lot of. The PGM framework provides a vital tool for anyone World Health Organization desires to find out a way to reason coherently from restricted and vociferous observations.

In a probabilistic graph, any two edges[3][4] e_i and e_j are called conditionally Independent if $p(e_i, e_j) = p(e_i)p(e_j)$, and conditionally dependent if $p(e_i, e_j) \neq p(e_i)p(e_j)$ [2]. For the standard probabilistic graph model, any two edges are conditionally autonomous of each other. Typically, beginners and mutual exclusion among adjacent edges area unit usually determined in numerous graph headed applications. As one of the fundamental data processing techniques, bunch is wide employed in numerous graph analysis [5][6] applications admire community exposure, index construction, etc. This paper focuses on clustering correlated probabilistic graphs that aims to partition the vertices into many disconnected clusters with high intra-cluster and low inter-cluster similarity, as illustrated. Next, we are going to inspire the matter of bunch correlative probabilistic graphs exploitation many applications.

II. RELATED WORK

In Protein-Protein Interaction (PPI) networks [1], the interaction between 2 proteins is mostly established with a likelihood property thanks to the limitation of observation strategies. additionally, it's been verified that the interaction between super molecules A and B can influence the interaction between protein A and another protein C, if A, B and C have some common options. It's been verified that the likelihood of pairwise interaction and correlation among edges [4] will be derived from applied mathematics models. Bunch applied to such related to probabilistic protein-protein interaction network knowledge is useful to find complexes to research the structure properties of the PPI Network.

To cluster a correlate probabilistic graph G , a possible world graph G_i of G can be sculptural as a settled internal representation sampled from the correlate probabilistic graph in step with the chance distribution. The edit distance[6] $D(G_i, Q)$ from G_i to the cluster graph letter is outlined because the variety of edges that require to be superimposed or removed to remodel G_i into Q . By evaluating all the potential world graph instances, the estimated edit distance, denote as $D(G, Q)$ is obtained and viewed as a activity to gauge the deviation from a correlate probabilistic graph to the cluster graph. Hence, a smaller deviation implies a a lot of specific result, and our objective turns to the goal of result a cluster graph Q that may minimize $D(G, Q)$. However, it's very long if we have a tendency to calculate the expected edit distance by considering all potential world graphs. To resolve this drawback, we have a tendency to propose a completely unique estimation model which needs the high-octane generation of a position access order once hard conditional chances. The estimation model has obvious error bounds.

III. PROPOSED METHODOLOGY

The analysis of relevant works on the systems for parallel big data processing and application of CEP for IoT, we should arrive to the following conclusions. First, there are works devoted to the implementation of the CEP in networks similar to IoT. However they do not target issues of parallel big data processing. Second, there are CEP systems with parallel processing of big data.

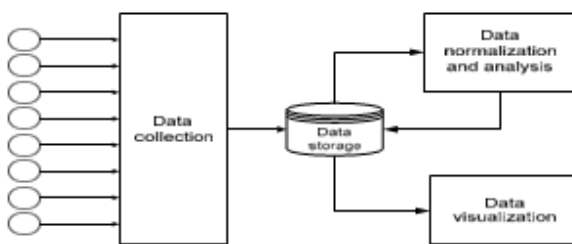


Fig:1 Flow diagram of the proposed methodology

However, these systems cannot be implemented in IoT because of existing constraints on communication channels' bandwidth and performance of computational nodes. Finally, monitoring of computer network security is relatively new and insufficiently investigated area for systems of parallel big data processing.

The parallel data processing system for IoT security monitoring is intended for collection and pre-processing of large amounts of information about security events that are received from the terminal devices ('things') and elements of

network infrastructure (routers, antiviruses, operation systems, database management systems, firewalls, etc.).

The security event sources generate large volumes of heterogeneous traffic that can be classified as Big Data. Therefore, this system can compete with traditional systems of network security if they are not able to cope with the growing volume of traffic in the IoT network.

The system is intended for implementation in the Hadoop environment. It consists of the following functional components: the data collector component that is responsible for the timely and accurate receipt of information about security events from sources of different types; component of the data storage, providing secure storage of data and efficient processing of requests; component for data normalization and analysis, converting the data collected to the unified format and performing on them the main preprocessing steps; the data visualization component, allowing real-time visual analysis using previously developed models of visualization.

3.1 Merits

Designing computer output should proceed in an organized, well thought out manner; the right output must be developed while ensuring that each output element is designed so that people will find the system can use easily and effectively.

IV. SYSTEM IMPLEMENTATION

4.1 Data Collection

The data collection component is designed to collect traffic on machines of the controlled IoT network and to send it to the data storage component. Machines involved in the data collection component, by their functional purpose are divided into the following types: collecting, used to collect the data; coordinating, designed to control data collection; receiving, intended to receive the data collected.

4.2 Data Storage

The data storage component supports operation of the distributed data storage, which is represented with the help of Hadoop Distributed File System (HDFS). HDFS divides massive amounts of file into chunk (64MB) units and stores them in nodes, called Data Nodes. Metadata, showing how input data are distributed across Data Nodes are stored in a node called Name Node. Name Node manages the metadata of HDFS, controlling input/output data between users and Data Nodes. Data Nodes save actual data, while directly performing

the input/output data with users and the block copying. Data Node processes are running on each receiving machine. The Name Node process is running on the coordinating machine. All the logs of the input traffic are loaded in the HDFS and are made available for further processing by the normalization and analysis component.

4.3 Data Normalization And Analysis

Data stored by the collection component are available from the HDFS for the normalization and analysis component. Data normalization is in reducing all the input data to the common internal format. CSV (Comma- Separated Values) is selected as such common format. The program that implements the normalization of the data has two modes: basic and advanced. In basic mode, the program scans the contents of the input file. For each encountered record it generates a separate line in the output file with a fixed set of attributes that describe the record structure.

4.4 Data Visualization

The data visualization component is intended for submission to the system operator of the result data obtained at the stage of analysis, in a form suitable for visual analysis. The component takes raw data from HDFS, where they were stored after processing by the normalization and analysis component. The visualization component uses a variety of standard and custom visualization models. Standard visualization models are: histograms; pie charts; linear graphics.

4.5 Indexing

Indexing collects, parses, and stores data to facilitate fast and accurate information retrieval. Index design incorporates interdisciplinary concepts from linguistics, cognitive psychology, mathematics, informatics, physics, and computer science. An alternate name for the process in the context of search engines designed to find web pages on the Internet is web indexing. Popular engines focus on the full-text indexing of online, natural language documents. Media types such as video and audio and graphics are also searchable. Meta search engines reuse the indices of other services and do not store a local index, whereas cache-based search engines permanently store the index along with the corpus. Unlike full-text indices, partial-text services restrict the depth indexed to reduce index size. Larger services typically perform indexing at a predetermined time interval due to the required time and processing costs, while agent-based search engines index in real time.

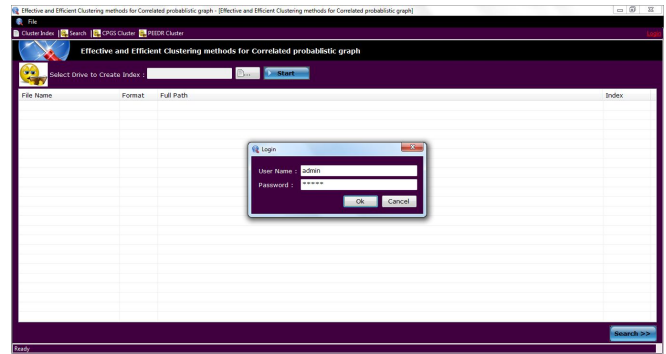


Fig A2.1 User Login

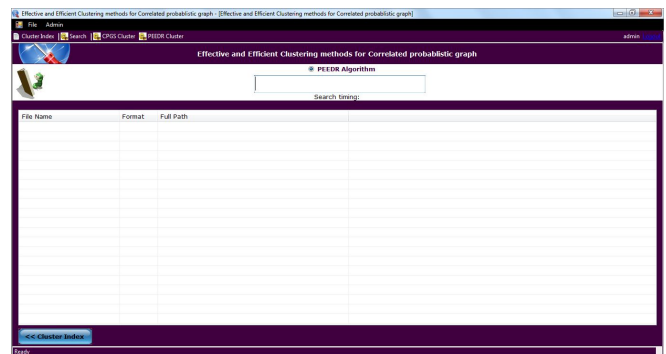


Figure A2.2 PEEDR Algorithm Search

V. CONCLUSION AND FUTUREWORK

In this paper we studied how to improve to the construction of security monitoring systems for IoT networks, based on the principles of parallel processing of data about security events. According to this approach, the network security monitoring system is implemented on the Hadoop platform with the ability to implement CEP technology. The proposed architecture of the developed parallel data processing system for IoT security monitoring consists of the components responsible for data collection, storage, normalization and analysis, and data visualization. Data normalization, analysis and visualization are performed on the fly. Data are stored in the distributed Hadoop storage that improves the reliability and efficiency of processing of the data requests. Implementation of the system for carrying out its experimental estimation was performed taking into account the computational constraints typical to IoT networks. Input streams were simulated by use of the special tested that generated security events in a fragment of the IoT network, and using an external database about traffic in the real computer network.

Experimental evaluation of the developed system showed that the developed system, despite the presence of limitations in computing resources, has high enough performance, comparable and in some cases considerably

exceeding implementations known from the literature. Further research is planned in direction of enhancement of the developed network security monitoring system by improving the speed of requests to data blocks of HDFS and the wider usage of CEP implementation tools.

VI. ACKNOWLEDGMENT

I would like to express my special thanks of gratitude to my college management. I would like to give my thanks to guide who gave me the golden opportunity to do this wonderful project on the. In addition, I would also like to thank my parents who helped me a lot in finalizing this project within the limited time frame.

REFERENCES

- [1] Bonchi.F, Gionis.A, Kollios.G and Potamias.M, PVLDB, vol. 3, no. 1, pp. 997–1008, Sept. 2010. “K-nearest neighbors in uncertain graphs,”
- [2] B, Jin.R, Liu.L and Wang.H, PVLDB, vol. 4, no. 9, pp. 551–562, Jun. 2011. “Distance-constraint reachability computation in uncertain graphs,”
- [3] Chen.L, Wang.G, Wang.H and Yuan.Y, PVLDB, vol. 5, no. 9, pp. 800–811, May 2012. “Efficient subgraph similarity search on large probabilistic graph databases,”
- [4] Hua.M and Pei.J, in Proc. 13th Int. EDBT, New York, NY, USA, 2010, pp. 347–358. “Probabilistic path queries in road networks: Traffic uncertainty aware path selection,”
- [5] Flynn.P.J, Jain.A.K and Murty.M.N, ACM Comput. Surv. vol. 31, no. 3, pp. 264–323, Sept. 1999 “Data clustering: A review,”
- [6] Shamir.R, Sharan.R, and Tsur.D, Discrete Appl. Math., vol. 144, no. 1–2, pp. 173–182, 2004. “Cluster graph modification problems,”
- [7] Cetindil.I, Esmaelnezhad.J, Li.C, and Newman.D, in WebDB, 2012, pp. 7–12. “Analysis of instant search query logs,”
- [8] Miller.R.B, in Proceedings of the December 9-11, 1968, fall joint computer conference, part I, ser. AFIPS '68 (Fall, part I). New York, NY, USA: ACM, 1968, pp. 267–277. “Response time in man-computer conversational transactions,”
- [9] Henzinger.M.R, Marais.H, Moricz.M and Silverstein.C, “Analysis of a very large web search engine query log,”
- [10] Ackermann.M.R, Blömer.J, Kuntze.D, and Sohler.C, Algorithmica, vol. 69, no. 1, pp. 184–215, May 2014. “Analysis of agglomerative clustering,”
- [11] Broschart.A, Schenkel.R, Theobald.M, won Hwang.S and Weikum.G, in SPIRE, 2007, pp. 287–299. “Efficient text proximity search,”
- [12] Shi.S, Suel.T, Wen.J.R, Yan.H and Zhang.F. in CIKM, 2010, pp. 1229–1238. “Efficient term proximity search with term-pair indexes,”
- [13] Shi.S, Wen.J.R, Yu.N and Zhu.M. in CIKM, 2008, pp. 679–688. “Can phrase indexing help to process non-phrase queries?”
- [14] Jain.A and Pennacchiotti.M, in COLING, 2010, pp. 510–518. “Open entity extraction from web search query logs,”
- [15] Grabski.K and Scheffer.T, in SIGIR, 2004, pp. 433–439. “Sentence completion,”