

SYSTEM FAILURE DETECTION AND DIAGNOSIS BY ANALYZING SYSLOG AND SNS DATA: APPLYING BIG DATA ANALYSIS TO NETWORK OPERATIONS

D.Vimala¹, I. Mary Linda²

^{1,2}Assistant Professor, Dept of CSE, BIST, BIHER, Bharath University, Chennai – 73

¹vimala.cse@bharathuniv.ac.in, ²marylinda.cse@bharathuniv.ac.in

Abstract: We present two major information examination strategies for diagnosing the reasons for system disappointments and for identifying system disappointments early. Syslogs contain log information created by the framework. We dissected syslogs what's more, prevailing with regards to distinguishing the reason for a system disappointment via consequently learning more than 100 million logs without requiring any past learning of log information. Investigation of the information of an interpersonal interaction benefit (in particular, Twitter) empowered us to recognize conceivable system disappointments by extricating system disappointment related tweets, which represent under 1% of all tweets, continuously and with high exactness.

Keywords: big data, syslog, network failure detection

1. Introduction

Web convention (IP) systems comprise of numerous sorts of hardware from various merchants. These systems are winding up noticeably significantly more unpredictable in light of the fact that of the expanding interest for new and diverse applications. Also, large portions of these applications are given by various system administrators and gadgets, also, this makes it extremely hard to analyze[1,2] organize disappointments when they happen. Thus, it is exceptionally critical to create techniques to productively recognize organize disappointments and analyze their causes. In this article, we present two strategies for breaking down information from syslogs and from a person to person communication benefit (SNS) to accomplish early system disappointment recognition and to analyze the reason for the system disappointment that current working techniques can't address.

Log data analysis

Arrange administrators screen different sorts of data for example, trap data from system components, organize activity, CPU (focal preparing unit)/memory utility information, and syslog information. Specifically, the syslog information of system components, for example, switches, switches, and RADIUS (Remote Access Dial

In Client Service) servers incorporate nitty gritty and exact data for investigating and checking the wellbeing of systems when arrangements change. Be that as it may, breaking down log information has turned out to be extremely troublesome for the accompanying reasons[3]: (i) There are different sorts of logs, which list messages with low or high seriousness. What's more, the increment in the quantity of system components implies there is an enormous volume of complex log information, and it is along these lines important to extricate data precisely and productively all together to complete investigating and preventive maintenance[4-6].

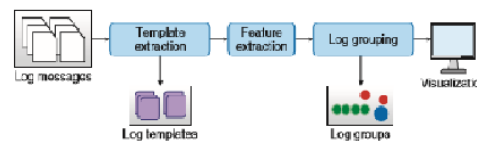


Figure 1. Flowchart for visualization of logs.

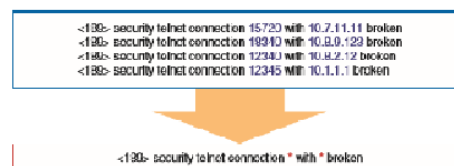


Figure 2. Conceptual image of log template extraction.

(ii) The log organize relies on upon every seller or administration. Subsequently, understanding the significance of each log message requires profound area information of each configuration. To defeat these issues, we have built up a method to dissect syslogs that includes consequently extricating the connections or variations from the norm from log messages utilizing machine-learning strategies without depending on any area information about the design or the seller of log information (Fig. 1). This examination system comprises of four stages: log layout extraction, log include extraction, log gathering, and perception of irregular occasions[9].

Log template extraction

log messages contain various parameters such as IP address, host name, and PID (process identification). Because parameter words are very rare, log messages with unique parameters may never appear twice even though the events the messages signify are the same. Therefore, we automatically extract a primary template from all log messages based on the observation that parameter words appear infrequently in comparison with template words in the other positions (Fig. 2). The log template enables us to easily correlate log messages[7-8].

Feature extraction

As mentioned before, the vendor’s severity of a log message is not necessarily reliable because it is not directly related to the actual network abnormality. Therefore, we need to quantify the abnormality and normality of logs without considering the severity of the message and without requiring any domain knowledge. For example, firewall logs and link down/up logs related to users’ connect/disconnect events contain very common messages and can be considered. Also, the logs generated by cron* jobs or in regular monitoring are not as frequent but are generated periodically on a daily basis. Therefore we define the frequency and periodicity features for log messages[10]

Log grouping

Ordinarily, arrange administrators don’t utilize a one-line log message, but instead, a gathering of logs. For instance,a switch reboot occasion prompts numerous logs, which shows that different procedures begin in the meantime. Subsequently, we have to gathering them as far as their co occurrence. Gathering logs lessens the volume of logs and helps administrators comprehend the logs. For log gathering, we utilize the machine learning procedure known as non-negative grid factorization (NMF) by changing over information log information into a lattice (Fig. 3)



Figure 3.Imageofloggrouping

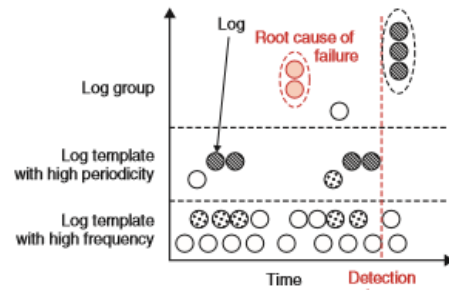


Figure 4. Conceptual image of log visualization

Visualisation

Ordinarily, organize administrators don't utilize a one-line log message, yet rather, a gathering of logs. For instance, a switch reboot occasion incites different logs, which shows that different procedures begin in the meantime. Along these lines, we have to gathering them as far as their cooccurrence[11] . Gathering logs lessens the volume of logs and helps administrators comprehend the logs. For log gathering, we utilize the machine learning system known as non-negative framework factorization (NMF) by changing over info log information into a network (Fig. 3).

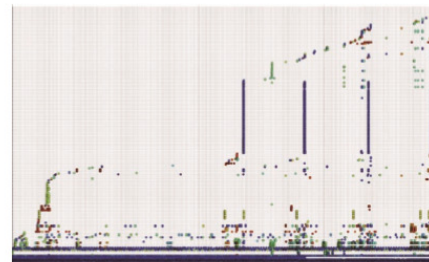


Figure 5. Example of log graph (one week’s syslog data).

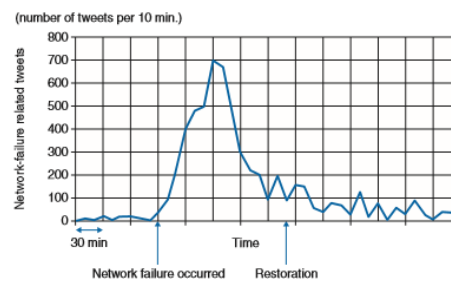


Figure 6. Time series of tweet counts related to an actual network failure

Twitter analysis

In this progression of the investigation, log information are communicated as a diagram. A calculated picture of log perception is appeared in Fig. 4, and a case of a log diagram is appeared in Fig. 5. In both figures, the even hub speaks to time, and the vertical pivot speaks to the format or log bunches specified before. Each point in the diagram speaks to the event of each log format or log assemble at each time[12]. Hosts are recognized by their diverse hues and examples in this illustration. The request of log layouts or log bunches on Arrange administrators can screen organize hardware by utilizing observing innovation, for example, SNMP (basic system administration convention). Despite the fact that they can distinguish equipment disappointments, it is troublesome for organize administrators to recognize disappointments brought about by programming bugs or to distinguish quality weakening due to clog. Thusly, a few cases end up plainly quiet disappointments, which can't be distinguished by system administrators. We have concentrated an approach to screen an informal communication benefit (SNS), to be specific, Twitter [13], to find issues influencing endorsers. For instance, we can see a surge in tweets about system disappointments when a system disappointment happens, as appeared in Fig. 6. We created a framework to screen Twitter progressively by checking for surges in these sorts of tweets[14].

2. System requirements

I.

Twitter is a prevalent stage for talking about incalculable discussion themes, and the quantity of tweets presently surpasses 400 million every day [3]. Japanese tweets alone record for 80–100 million tweets for every day. Since the quantity of tweets that identify with system issues is little in the aggregate number of tweets, we require an approach to extricate just applicable tweets (first necessity). Also, to identify the region where a organize disappointment happens, we require an approach to decide the area of the tweeters (second requirement).

Method to extract only network -failure related tweets

We found in our examination that catchphrase coordinating, a conventional approach to pursuit tweets, was not adequate for mechanized observing on the grounds that it brought about numerous false positives. This happens when the tweets contained the watchwords, however the tweets were most certainly not identified with issues with the system. For instance, if we look utilizing the watchwords call and drop, we may get tweets, for example, "I dropped my telephone in the can so I can't call or content". Since watchwords, for example, call what's more, drop are not organize particular words, catchphrase coordinating may prompt

a ton of false positive tweets that contain the watchwords however not the point of the arrange issue[15]. The system disappointment location engineering is appeared in Fig. 7.

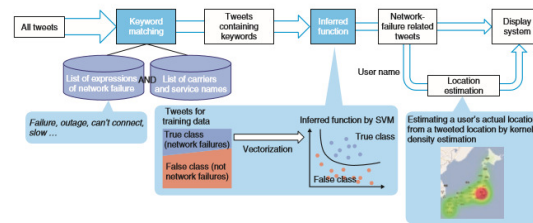


Fig. 7. Framework of network-failure detection system using Twitter.

We utilize regulated adapting, specifically, SVM (bolster vector machine), to stifle the false positives. Managed learning utilizes an informational collection of preparing cases. Each preparation illustration comprises of a couple of the content of a tweet and a mark demonstrating whether the tweet is identified with a system disappointment. An administered learning calculation dissects the preparation information and produces an induced capacity to gap tweets into those that are identified with system disappointments and those that are most certainly not. In our approach, each tweet is interpreted into a vector by utilizing the pack of-words strategy, which is a customary strategy in archive arrangement. This technique can be relied upon to stifle the false positives by factually considering all words showing up in one tweet. We assessed the viability of our strategy by applying it to a whole year of Twitter information. Six system disappointments were accounted for by a system transporter in that period. We assessed the system disappointment discovery framework by tallying the quantity of tweets that were grouped by our technique. At the point when the tally surpassed a specific edge, we viewed it as an alarm of a system disappointment. We additionally utilized the watchword coordinating strategy for correlation[16]. Both techniques distinguished the 6 real system disappointments. Nonetheless, the catchphrase just strategy likewise erroneously recognized 94 occasions, though the machine-learning strategy smothered all of those and had just 6 false detections.

Method to determine the location of tweeters

Twitter has a capacity to join the client's area by GPS (Global Positioning System) information, however most clients pick not to select into this capacity. Consequently, we have to gauge the area of Twitter clients who composed the system disappointment related tweets. A few reviews have utilized the predisposition of a circulation of words, which fundamentally includes tongue qualities, to assess a

client's area. Be that as it may, these reviews evaluate an unpleasant granularity of zones, for example, the Kanto area with a mistake of around 150 km and don't meet our necessity, which is to accomplish at any rate prefecture-level area (a mistake of under 50 km). Along these lines, we concentrated a high-precision area estimation strategy that utilizes gazetteer data, which incorporates the sets of a geographic name and its organizes. While most tweets don't contain GPS data, many tweets contain a geographic name. In spite of the fact that clients may tweet the geographic names of puts other than where they are really found, the covered areas of a significant number of their tweets will make it conceivable to assess their area on the grounds that Twitter is an administration for clients to post what they are doing. We utilized the part thickness estimation technique to cover the tweets of individual tweeters. We assessed the estimation mistake of clients whose areas were known and found that the estimation mistake was under 50 km for 66% of those clients. Moreover, the estimation mistake was under 25 km for half of all clients, which showed that our technique was powerful.

3. Conclusion

We introduced a big-data approach consisting of syslog and SNS analysis to predict or detect network failures. In cooperation with group companies, we are now evaluating the efficiency of syslog analysis using actual syslog data. We are also preparing a proposal for group companies for the use of SNS analysis as a tool for detecting silent failures.

References

- [1]. Udayakumar R., Kaliyamurthie K.P., Khanaa, Thooyamani K.P., Data mining a boon: Predictive system for university topper women in academia, *World Applied Sciences Journal*, v-29, i-14, pp-86-90, 2014.
- [2]. Kaliyamurthie K.P., Parameswari D., Udayakumar R., QOS aware privacy preserving location monitoring in wireless sensor network, *Indian Journal of Science and Technology*, v-6, i-SUPPL5, pp-4648-4652, 2013.
- [3]. BrinthaRajakumari S., Nalini C., An efficient cost model for data storage with horizontal layout in the cloud, *Indian Journal of Science and Technology*, v-7, i-, pp-45-46, 2014.
- [4]. BrinthaRajakumari S., Nalini C., An efficient data mining dataset preparation using aggregation in relational database, *Indian Journal of Science and Technology*, v-7, i-, pp-44-46, 2014.
- [5]. Khanna V., Mohanta K., Saravanan T., Recovery of link quality degradation in wireless mesh networks, *Indian Journal of Science and Technology*, v-6, i-SUPPL.6, pp-4837-4843, 2013.
- [6]. Khanaa V., Thooyamani K.P., Udayakumar R., A secure and efficient authentication system for distributed wireless sensor network, *World Applied Sciences Journal*, v-29, i-14, pp-304-308, 2014.
- [7]. Udayakumar R., Khanaa V., Saravanan T., Saritha G., Retinal image analysis using curvelet transform and multistructure elements morphology by reconstruction, *Middle - East Journal of Scientific Research*, v-16, i-12, pp-1781-1785, 2013.
- [8]. Khanaa V., Mohanta K., Saravanan. T., Performance analysis of FTTH using GEAPON in direct and external modulation, *Indian Journal of Science and Technology*, v-6, i-SUPPL.6, pp-4848-4852, 2013.
- [9]. Kaliyamurthie K.P., Udayakumar R., Parameswari D., Mugunthan S.N., Highly secured online voting system over network, *Indian Journal of Science and Technology*, v-6, i-SUPPL.6, pp-4831-4836, 2013.
- [10]. Thooyamani K.P., Khanaa V., Udayakumar R., Efficiently measuring denial of service attacks using appropriate metrics, *Middle - East Journal of Scientific Research*, v-20, i-12, pp-2464-2470, 2014.
- [11]. R.Kalaiprasath, R.Elankavi, Dr.R.Udayakumar, Cloud Information Accountability (Cia) Framework Ensuring Accountability Of Data In Cloud And Security In End To End Process In Cloud Terminology, *International Journal Of Civil Engineering And Technology (Ijciet)* Volume 8, Issue 4, Pp. 376–385, April 2017.
- [12]. R.Elankavi, R.Kalaiprasath, Dr.R.Udayakumar, A fast clustering algorithm for high-dimensional data, *International Journal Of Civil Engineering And Technology (Ijciet)*, Volume 8, Issue 5, Pp. 1220–1227, May 2017.
- [13]. R. Kalaiprasath, R. Elankavi and Dr. R. Udayakumar. Cloud. Security and Compliance - A Semantic Approach in End to End Security, *International Journal Of Mechanical Engineering And Technology (Ijmet)*, Volume 8, Issue 5, pp-987-994, May 2017.
- [14]. Thooyamani K.P., Khanaa V., Udayakumar R., Virtual instrumentation based process of agriculture by automation, *Middle - East Journal of Scientific Research*, v-20, i-12, pp-2604-2612, 2014.
- [15]. Udayakumar R., Thooyamani K.P., Khanaa, Random projection based data perturbation using geometric transformation, *World Applied Sciences Journal*, v-29, i-14, pp-19-24, 2014.
- [16]. Udayakumar R., Thooyamani K.P., Khanaa, Deploying site-to-site VPN connectivity: MPLS Vs IPsec, *World Applied Sciences Journal*, v-29, i-14, pp-6-10, 2014.

