Automatic Classification of System Logs

Siavash Ghiasvand[§] and Florina M. Ciorba*

§Technische Universität Dresden, Germany ★University of Basel, Switzerland

starting acron starting monitor Acron started on 2018-03-01 Jobs will be executed sequentially Normal exit (3 jobs run) finished acron (root) CMD (run-parts)

Sample system log entries

Proposed classification approach

Marking common terms in system log entries (marking terms as invariant)

Functions

Parentheses and quotations

File/application addresses

Current state of practice

- Log entry messages are unstructured[1] \bullet
- Most of log entry messages are repetitive (high frequency)[2] •
- Manual analysis is not efficient for large scale computing systems[4] \bullet
- Automatic classification requires pre-classified/labeled sample logs \bullet
- Available methods are system-specific^[5] lacksquare

Goal

Automatic classification of system logs

Contributions

- The approach does not require data pre-classification and labeling
- General approach

Sample system log entries

(sshd:session): session closed for root Job `cron.daily' terminated 96734

Live demo and sample script: ghiasvand.net/u/isc18

Dates, Hours, IP addresses

Names, Numbers

Classifying log entries into classes, based on Levenshtein[3] similarity metric

Harmonizing differences across entries of the same class

Refining and adjusting the resulting regular expressions

Evaluation and preliminary results

Datasets 1 month of system logs 108 nodes (1665 020 ontries)

Classes 700

Job `cron.weakly' ter (sshd:session): sessi Job `cron.hourly' ter (sshd:session): sessi (sshd:session): sessi Job `cron.daily' term	Eminated 537352 on closed for siavash on closed for florina on closed for s125342 minated 14038
CLASS 1	CLASS 2
#PARA#: session closed for root	Job `cron.hourly' terminated #DGIT#
#PARA#: session closed for siavash	Job `cron.daily' terminated #DGIT#
#PARA#: session closed for florina	Job `cron.weakly' terminated #DGIT#
#PARA#: session closed for s125342	Job `cron.daily' terminated #DGIT#
CLASS 1	CLASS 2
#PARA#: session closed for #VARI#	Job #VARI# terminated #DGIT#
<pre>RegEx 1 (\(.+?\))\: #: session closed for (.+)</pre>	RegEx 2 Job (.+) terminated ([0-9]+)

Conclusion

- Automatic classification of system logs is possible ${\color{black}\bullet}$
- Low frequency log entries reduce the overall classification \bullet accuracy
- No relation between number of log entries and number of automatically generated classes
- Strong relation between log frequency and classifiers accuracy

Future work

Addressing entries with two terms and syslog entries

		 100 node 100 node 100 node 100 node 99 node 270 node 677 node 	$e^{1,000,0237}$ $e^{1,699,687}$ $e^{1,699,697}$ $e^{1,699,697}$ $e^{$	7 entries) 8 entries) 9 entries) 7 entries) 1 entries)	Automatic Classification	546 568 637 622 860
		7. 99 node	1 year of systems (136,609,978	tem logs 3 entries)		1090
mber of system log entries (Thousands)	5000 4500 4000 3500 3000 2500 2000 1500 1000 500					
Nu	0	2001 to 2108	1100 to 1199	5100 to 51	6001 to 6099	1001 to 1270
🗕 # Rav	v entries	1665029	1699237	3555088	3 3795660	4549687
# Class	sses	722	546	568	637	622
				Nada IDa (

Node IDS (range

- Number of generated classes remains almost identical Heterogeneity of nodes has no significance on the number of classes High frequency log entries improve classification
- 41% of classes are common in all datasets

qO

63% of classes are common in several sets



that form blocks

References

[1] "The syslog protocol," http://tools.ietf.org/html/rfc5424, [Online; accessed March 2018] [2] S. Ghiasvand and F. M. Ciorba, "Anonymization of system logs for preserving privacy and reducing storage," in Proc. of the Future of Information and Communications Conference, 2018 [3] Levenshtein, VI. "Binary Codes Capable of Correcting Deletions, Insertions and Reversals." Soviet Physics Doklady 10 (1966): 707.

[4] T. Li et. al. "FLAP: An End-to-End Event Log Analysis Platform for System Management". In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '17). ACM, New York, NY, USA, 2017 [5] K. Berkay "Unsupervised Anomaly Detection in Unstructured Log-Data for Root-Cause-Analysis", master thesis, Tampere University of Technology, Finland, 2015

Acknowledgment

This work is in part supported by the German Research Foundation (DFG) within the Cluster of Excellence 'Center for Advancing Electronics Dresden (cfaed)', and by Eucor - The European Campus, within 'Data Analysis for Improving High Performance Computing Operations and Research'.

