Anomaly Detection in Data Center IT Infrastructure using Natural Language Processing and Time Series Solutions

Elisabetta Ronchieri^{1,2} L. Giommi¹, D. C. Duma¹, A. Costantini¹, D. Salomoni¹, F. Pacinelli²



¹INFN CNAF, ²Univ. Bologna

ISGC 2023, March 19-24, 2023, Taipei, March 22, 2023

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Table of Contents

1 INFN CNAF Data Center IT Infrastructure

- 2 Goal adding intelligence
- 3 Data Sources
- 4 Methodology
- 6 Results
- Discussions and Future Works

7 Q&A

E. Ronchieri (INFN CNAF, Univ. Bologna)

INFN CNAF Data Center in Bologna, Italy

 It focuses on technological development and knowledge transfer to heterogeneous experiments, national and European projects, and industry (whenever possible).
 Image provided by Pier Paolo Ricci, INFN CNAF.



Facts

- over 1700 active users
- over 50 experiments
- over 65,000 CPU cores
- 41 PB of disk

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• 98 PB of tapes

INFN CNAF Data Center Pillars

Pillars already supported

✓ Connectivity:

to cloud and infrastructure for data transmission

✓ Security:

for data and privacy protection

New Pillar

<u>المارم</u>

🛦 Intelligence

- through infrastructure and algorithms in order to <u>extract value from service</u> and machine data and convert them into useful information
- perform predictive maintenance

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ 三日= のへぐ

Table of Contents

1 INFN CNAF Data Center IT Infrastructure

② Goal - adding intelligence

- 3 Data Sources
- 4 Methodology

6 Results

Discussions and Future Works

7 Q&A

Anomaly Detection through Logs and Monitoring Metrics

Facts

- large amount of data
- over 1,000 different services running
- data flowing 24/7

Goal

perform predictive maintenance

How?

 define intelligence based on logs and monitoring metrics in order to <u>detect anomalies</u> and properly intervene



Table of Contents



2 Goal - adding intelligence



4 Methodology



Discussions and Future Works

7 Q&A

Data Sources

Time range for this study: 6th June, 2020 - 21st July, 2021

Log Files

- 40 kinds of log files
- 1143 distinct machines

Monitoring Metrics
• 18 metrics
• 3 categories of monitoring metrics: io-stat,
load average, and memory

Log Files and Monitoring Metrics

Log Files

- Logs are created by Linux system services and contain semi-structured texts.
- They are used to keep records of occurring events, analyze and debug system failures.
- The services may be highly verbose [Dec+20b].
- On a single machine, there are several different log files.
- Each log file corresponds to <u>a distinct service</u>.
- The log files have different formats.

Monitoring Metrics

- Each machine usually runs different services and we can collect information on:
 - memory statistics, representing memory usage;
 - Ithe load the machine has been under, averaged over multiple time frames (i.e., load_avg)ç
 - the central processing unit (CPU) statistics and input/output (I/O) statistics (i.g., io_stat).

An Example of the ALRT Service Log File

date	time	timestamp	hostname	ip	process_na
2020/08/11	11:57:17	1597139837.0	ddn-04-a.cnaf.infn.it	*	ALRT
2020/07/08	14:36:39	1594211799.0	ddn-04-a.cnaf.infn.it	*	ALRT
2020/07/08	14:36:39	1594211799.0	ddn-04-a.cnaf.infn.it	*	ALRT
2020/07/27	16:24:23	1595859863.0	ddn-04-a.cnaf.infn.it	*	ALRT
2020/07/25	12:24:08	1595672648.0	ddn-04-a.cnaf.infn.it	*	ALRT

Log message



ne	msg
	ALRT CLI_MAIN Failing Disk 51P S/N JK1101YAJXPDEZ Reason(User Requested)
	ALRT AVR_MON Left Power Supply Failure
	ALRT AVR_MON Left Side AC Line Low
	ALRT CLI_MAIN Failing Disk 19G S/N JK1101YAJSDYXV Reason(User Requested)
	ALRT DC_REC Failing Disk 42P S/N JK1101YAKAB69V Reason(IO Timeout)

- date, time and timestamp express when the instance has been registered in the local time zone.
- hostname and ip columns provide information on the machine and network on which the event occurred, respectively.
- process_name is the name of the service (log file).
- msg reports the textual message.

An Example of the memory.used Monitoring Metric

name	tags	time	domain	duration	metric	value
memory.used	host=api-int.cnsa.cr.cnaf.infn.it	1,62643E+18	cnsa.cr.cnaf.infn.it	0.22133	metrics-memory	1,37128E+09
memory.used	host=api-int.cnsa.cr.cnaf.infn.it	1,62644E+18	cnsa.cr.cnaf.infn.it	0.23058	metrics-memory	1,39689E+16
memory.used	host=api-int.cnsa.cr.cnaf.infn.it	1,62644E+18	cnsa.cr.cnaf.infn.it	0.22633	metrics-memory	1,46820E+16
memory.used	host=api-int.cnsa.cr.cnaf.infn.it	1,62644E+18	cnsa.cr.cnaf.infn.it	0.21558	metrics-memory	1,49160E+16
memory.used	host=api-int.cnsa.cr.cnaf.infn.it	1,62645E+18	cnsa.cr.cnaf.infn.it	0.21825	metrics-memory	1,49583E+16

0.00 0.00<		0		A	_	_	_		-		-		п.,
Tasks 94 total, 1 (rmining, 0) 3 topping, 0 stopping,	ton	97:37:49	un 5	dave		28	1.000		100	d av	verage: 0.0		
Const. A Bunc. B Bunc. A Bunc. <th< td=""><td>Tasks</td><td>94 tota</td><td>1</td><td>1</td><td>inning</td><td>0</td><td>1 5100</td><td></td><td>ina</td><td></td><td>stopped</td><td>a zombie</td><td>100</td></th<>	Tasks	94 tota	1	1	inning	0	1 5100		ina		stopped	a zombie	100
Her: 15658288.total. 137568.total. 137684.total. 137684.total. USER PR NZ VAR HZS SNR S SCOUND 217924.toched 1 PR NZ VAR HZS SNR S SCOUND 217924.toched 1 POT ZAL HZS S.R.S SCOUND 217924.toched 1 POT ZAL HZS S.R.S SCOUND 217924.toched 1 POT ZAL HZS S.R.S SCOUND 11766.toched 1 POT ZAL SCOUND S.R.S SCOUND 11766.toched 1 POT ZAL SCOUND S.R.S SCOUND 11767.toched 1 POT ZAL SCOUND S.R.S SCOUND 11767.toched 1 POT ZAL S.R.S S.R.S SCOUND 11877.toched 11777.toched 1 POT ZAL S.R.S S.R.S S.R.S SCOUND 11877.toched 1 POT ZAL S.R.S S.R.S S.R.S 118788.toched 117787.toched </td <td>Coule</td> <td>. 0 0bus</td> <td>· .</td> <td>0.00</td> <td>/ A</td> <td>ann i</td> <td>100 0</td> <td>1</td> <td>id a</td> <td>0.5</td> <td>wa 0 0khi</td> <td>0.0%ci 0.0%ct</td> <td></td>	Coule	. 0 0bus	· .	0.00	/ A	ann i	100 0	1	id a	0.5	wa 0 0khi	0.0%ci 0.0%ct	
Supp. 409784k test, 10164k used, 401548k free, 210792k cached 1 root 26 10356 1125 60 10355.01 10164.01 1 root 26 10356.01 1125 60 10356.01 1056.01 1 root 26 10356.01 1125 60 10556.01 10146.01 3 root RT 0 0 5 0.0 0.0 008.06 101770/0 4 root RT 0 0 5 0.0 0.0 008.06 101770/0 5 root RT 0 0 5 0.0 0.0 008.06 101770/0 7 root 20 0 0 5 0.0 0.0 0.0 0.0 9 root 20 0 0 5 0.0 0.0 0.0 0.0 9 root 20 0 </td <td>Mont</td> <td>16058284</td> <td>tota</td> <td>1</td> <td>16775</td> <td>LEAL I</td> <td>us od</td> <td></td> <td>1746</td> <td>01</td> <td>free 131</td> <td>336k buffors</td> <td></td>	Mont	16058284	tota	1	16775	LEAL I	us od		1746	01	free 131	336k buffors	
Total Close Close <th< td=""><td>Suan</td><td>4997844</td><td>tota</td><td></td><td>181</td><td>644</td><td>used,</td><td></td><td>48154</td><td>aL .</td><td>free, 218</td><td>792k cached</td><td></td></th<>	Suan	4997844	tota		181	644	used,		48154	aL .	free, 218	792k cached	
100 105 <td>Swup :</td> <td>4997046</td> <td></td> <td></td> <td>10.</td> <td>1046</td> <td>1960,</td> <td></td> <td>40134</td> <td>U.</td> <td>1100, 210</td> <td>Park cucifed</td> <td></td>	Swup :	4997046			10.	1046	1960,		40134	U.	1100, 210	Park cucifed	
1 root 20 0 0 1 0 <th0< th=""> 0 0 0</th0<>	PTD	IICED	DD	MT	VIDT	DEC	CHD	c	S-C DI I	AMER	M TIMEA	COMMAND	
2 root 2 root<	1	root	20	0.1	0356	1369	1132	c	0.0	0.1	1 0.56.52	init	- U.
3 root 8 0 <th0< th=""> <th0< th=""> 0 <th0< th=""></th0<></th0<></th0<>	5	root	20		0	1300	1152	č	0.0		0.00.00	kthroadd	
4 root 20 0 </td <td>2</td> <td>root</td> <td>DT</td> <td></td> <td></td> <td></td> <td></td> <td>c</td> <td>0.0</td> <td></td> <td>0 0:00.00</td> <td>minration/A</td> <td></td>	2	root	DT					c	0.0		0 0:00.00	minration/A	
S root RT 0 </td <td></td> <td>root</td> <td>20</td> <td></td> <td></td> <td></td> <td></td> <td>0</td> <td>0.0</td> <td></td> <td>0 0.00.00</td> <td>haoftired (8</td> <td></td>		root	20					0	0.0		0 0.00.00	haoftired (8	
6 root RT 0 <th0< th=""> 0 0 0</th0<>		root	DT					0	0.0		0 0:00.04	migration/8	
Troot 20 0 <td></td> <td>root</td> <td>DT</td> <td></td> <td></td> <td></td> <td></td> <td>6</td> <td>0.0</td> <td></td> <td>0 0.00.00</td> <td>inigration/e</td> <td></td>		root	DT					6	0.0		0 0.00.00	inigration/e	
8 root 20 0 </td <td>7</td> <td>root</td> <td>20</td> <td></td> <td></td> <td></td> <td></td> <td>0</td> <td>0.0</td> <td></td> <td>0 0:01.90</td> <td>watchuog/e</td> <td>10</td>	7	root	20					0	0.0		0 0:01.90	watchuog/e	10
0 0		root	20					5	0.0		0 0:33.09	evencs/o	- 11
9 FOOL 20 0 <td></td> <td>root</td> <td>20</td> <td></td> <td></td> <td>0</td> <td></td> <td>5</td> <td>0.0</td> <td>0.1</td> <td>0 0:00.00</td> <td>cgroup</td> <td>- 11</td>		root	20			0		5	0.0	0.1	0 0:00.00	cgroup	- 11
11 root 20 0 <td></td> <td>root</td> <td>20</td> <td></td> <td></td> <td></td> <td></td> <td>5</td> <td>0.0</td> <td></td> <td>0 0:00.00</td> <td>kne (per</td> <td>- 11</td>		root	20					5	0.0		0 0:00.00	kne (per	- 11
1 root 20 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	10	root	20			0		5	0.0		0 0:00.00	netns	- 11
12 root 20 0 0 0 0 5 0.0 0.0 0.0 0.00 0.00 pm 13 root 20 0 0 0 0 5 0.0 0.0 0.00 0.00 0.00 0.0	11	root	20					2	0.0		0 0:00.00	async/ngr	- 11
14 TODA 20 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	12	root	20					5	0.0	0.1	0 0:00.00	pm	- 11
14 root 20 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	13	root	20	0	0	0	0	5	0.0	0.0	0 0:00.00	xenwatch	- 11
15 FOOT 20 0 0 0 0 5 0.0 0.0 0:03.43 Sync_supers 16 Foot 20 0 0 0 0 5 0.0 0.0 0:02.94 bdi-default 17 Foot 20 0 0 0 0 5 0.0 0.0 0:00.00 kintegrityd/0	14	root	20			0		5	0.0	0.1	0 0:00.16	xenbus	- 11
15 root 20 0 0 0 0 0 0 0.0 0:02.94 bd1-default 17 root 20 0 0 0 0 5 0.0 0.0 0:00.00 kintegrityd/0	15	root	20		0	0	0	5	0.0	0.0	0 0:03.43	sync_supers	- 11
1/ root 20 0 0 0 0 5 0.0 0.0 0:00.00 kintegrityd/0	16	root	20	0	0	0		5	0.0	0.0	0 0:02.94	bdi-detault	- 11
	17	root	20	0	0	0	0	S	0.0	0.0	0 0:00.00	kintegrityd/0	

- **time** expresses when the metric has been registered in the local time zone.
- tags and domain provide information about the machine values are collected.
- name and metric provide the metric name and category.
- value contains the metric value.

Table of Contents

- 1 INFN CNAF Data Center IT Infrastructure
- 2 Goal adding intelligence
- 3 Data Sources
- 4 Methodology
- 6 Results
- Discussions and Future Works

7 Q&A

Main Challenges

Working with

- Different types of data (textual and numerical);
- Thousands of machines to be analysed;
- A completely unsupervised task.

Methodology



Anomaly Dictionary



Approach: The red highlights terms that may be likely anomalies [VRC22].

Correlation between Monitoring Metrics



- 18 monitoring metrics are correlated.
- Heatmap on the left shows how each pair of metrics is correlated.
- Metrics are both positively and negatively correlated.
- Selected one metric per class:

iostat.avg-cpu.pct_idle, load_avg.fifteen, memory.usage

Principal Component Analysis (PCA)

PCA

- It is used to perform dimensionality reduction.
- PCA
 - with reconstruction error RE allows us to recognize an anomalous observation.
- Approach: RE is larger for uncommon terms, so less observed observations.

Image shows a dimensional reduction to 2 components.



Clustering Algorithms

Density-Based Spatial Clustering of Applications with Noise

- <u>DBSCAN builds clusters from</u> the highly populated area of <u>observations</u>, <u>close</u> at most an *epsilon* value from each other.
- epsilon is computed performing parameter tuning
 - for log data, on overall distances between word-vectors
 - for monitoring metrics, by means of an elbow curve
- Approach: non-anomalous observations are expected to be concentrated closer to each other.

K-means

- K-means is a simple approach to partition data into K distinct, non-overlapping clusters.
- The only parameter is K that is computed performing parameter tuning.
- <u>The observations have been</u> partitioned among *K* clusters and then considered as potential anomalies provided belonging to the least populated cluster.
- Approach: expected large distances between anomalies and non-anomalies.

Anomaly Score (AS) - for log data

- For the entire message, <u>AS is</u> computed by doing the average value of the labels corresponding to the words that are part of the message.
- As this work is unsupervised, this process allows us to give a continuous measure of <u>anomalies</u> rather than a purely binary value, i.e. anomalous-non anomalous.
 - The binary value 0 or 1 for word represents respectively non-anomalous and anomalous terms.
- In the following, the words in red are the words identified as anomalous terms.



daily database available for update (local version: 26052, remote version: 26053)

Time Series Analysis - Meam and Variance Outlier Detection

For log data

- Messages have been vectorized and word-vectors from each message have been averaged to get the vector belonging to the message.
- The difference between consecutive vector-messages has been calculated (measured in terms of Euclidean distance between vectors), thus obtaining a time-series of the differences between messages.

For monitoring data

• metrics values have been standardized as the orders of magnitudes of the various metrics are really different from each other.

Approach: an interval of non-anomalous values is defined, including all the observations whose values are at most s tandard deviations σ away from the value of mean μ :

threshold = $\mu + s\sigma$, $s \in [1, 3]$.

E. Ronchieri (INFN CNAF, Univ. Bologna)

An example of Mean and Variance Outlier Detection

A <u>window size</u> of 20 observations and <u>tolerance</u> of 2 standard deviations on a specific machine, i.e. **cloud-ctrl01**.



Validation Criteria

For log data

- the anomaly score AS for each message has been considered, particularly the most recurring words in log messages with high/low AS;
- (a) a similarity percentage of more than 20%/25% has been considered as an alarming sign.

For monitoring data

• the Mann-Whitney test has been used to check for significant differences between anomalous and non-anomalous classified observations for every metric.

◆□ ▶ < @ ▶ < E ▶ < E ▶ E = のQ @ 23/35</p>

Table of Contents

1 INFN CNAF Data Center IT Infrastructure

- 2 Goal adding intelligence
- 3 Data Sources
- ④ Methodology

6 Results

Discussions and Future Works

7 Q&A

An Example of Log and Monitoring Occurrences for a Given Machine



◆□ ▶ < @ ▶ < E ▶ < E ▶ E = のQ @ 25/35</p>

Examples of Combining Results

Service	Technique	Parameter	% of correspondence
auditd	DBSCAN	eps = 0.1	50%
	K-Means	$n_clusters = 3$	97.5%
	PCA	thresh. $= 0.85\%$	20%
smartd	DBSCAN	eps = 0.1	55%
	K-Means	$n_clusters = 3$	53%
	PCA	thresh. $= 0.85\%$	55%

• Correspondence has been verified by checking logs and monitoring data occurring at less than 900 (15 minutes) seconds of difference.

• Anomalous and non-anomalous events are, at least to some extent, consistent in time.

Table of Contents

1 INFN CNAF Data Center IT Infrastructure

- 2 Goal adding intelligence
- 3 Data Sources
- 4 Methodology
- 6 Results

6 Discussions and Future Works

7 Q&A

E. Ronchieri (INFN CNAF, Univ. Bologna)

Discussions and Future Works

Discussions

- For log data, DBSCAN, K-means and Principal Component Analysis were very effective.
- For monitoring data, time series analysis provided interesting results, while treating messages as time-series observations did not show very meaningful results.
- We have got over 50% of correspondence between anomalous found in log files and monitoring metrics.

Future Works

- <u>To analyse in more detail the anomalies' nature</u> in order to perform a more precise identification of the root causes of anomalies.
- To develop an advanced model, e.g. a deep learning one, that uses both types of data with respect to timestamp, machine name and network info.

◆□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶

Table of Contents

1 INFN CNAF Data Center IT Infrastructure

- 2 Goal adding intelligence
- 3 Data Sources
- 4 Methodology
- 6 Results
- 6 Discussions and Future Works

🕖 Q&A

E. Ronchieri (INFN CNAF, Univ. Bologna)

< □ ▶ < @ ▶ < E ▶ < E ▶ E = の Q @ 29/35

Q&A

- Elisabetta Ronchieri, elisabetta.ronchieri@cnaf.infn.it
- Thanks to Diego Michelotto and Antonio Falabella, INFN CNAF, for their support in collecting monitoring metrics.

References I

- [Dec+20a] Leticia Decker et al. "Real-Time Anomaly Detection in Data Centers for Log-based Predictive Maintenance using an Evolving Fuzzy-Rule-Based Approach". In: (Apr. 2020). arXiv: 2004.13527 [cs.AI].
- [Dec+20b] Leticia Decker et al. "Real-time anomaly detection in data centers for log-based predictive maintenance using an evolving fuzzy-rule-based approach". In: 2020 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE). IEEE. 2020, pp. 1–8.
- [Gee21] GeeksforGeeks. *iostat command in Linux with examples*. www.geeksforgeeks.org/ iostat-command-in-linux-with-examples/". Online; accessed 08 June 2022. Dec. 2021.
- [K] Mathangi K. Load average: What is it, and what's the best load average for your Linux servers? https://www.site24x7.com/blog/load-average-what-is-it-andwhats-the-best-load-average-for-your-linux-servers. Online; accessed 12 July 2022.

◆□ → < @ → < E → < E → E = のへで 31/35</p>

References II

[Kap]	Vladimir Kaplarevic. <i>How to Check Memory Usage in Linux, 5 Simple Commands.</i> https://phoenixnap.com/kb/linux-commands-check-memory-usage. Online; accessed 12 July 2022.
[KBR16]	Tom Kenter, Alexey Borisov, and Maarten de Rijke. "Siamese CBOW: Optimizing Word Embeddings for Sentence Representations". In: (2016).
[VRC22]	L. Viola, E. Ronchieri, and C. Cavallaro. "Combining Log Files and Monitoring Data to Detect Anomaly Patterns in a Data Center". In: <i>Computers</i> 11.117 (2022). DOI: https://doi.org/10.3390/computers11080117.

Table of Contents

8 Appendix I

- A Subset of Monitoring Metrics
- Useful Parameters

A Subset of Monitoring Metrics

Category	Metric	Description
load_avg [K]	one	load average value of the machine every minute
	five	load average value of the machine every five minutes
	fifteen	load average value of the machine every fifteen minutes
memory [Kap]	available	estimation of how much memory is available for starting new applications, without
	buffers	swapping memory reserved by the operating system to allocate as buffers when the process needs them
	cached	recently used files stored in RAM
	dirty	memory waiting to be written back to disk
	free	unused memory
	total	total installed memory on the device
	used	memory currently in use by running processes
ostat [Gee21]	pct_idle	the percentage of time that the CPU or CPUs were idle and the system did not have an outstanding disk I/O request.
	pct_iowait	the percentage of the time that the CPU or CPUs were idle during which the system had an outstanding disk I/O request
	pct_nice	the percentage of CPU utilization that occurred while executing at the user level with a nice priority
	pct_user	the percentage of CPU being utilization that while executing at the user level
		 < ロ > < 豆 = < つへの

Useful Parameters

Parameters

- **Dim. reduction**: if applicable, whether or not PCA was applied on data to reduce their dimension
- **N. comps**: if applicable, (whether dimensionality reduction was applied), the number of principal components obtained.
- N. clusters: if applicable, the number of clusters defined for partitioning.
- Window size: if applicable, the length of the sliding windows (measured in a number of observations per partition).
- **Epsilon**: Estimating the Epsilon was very problematic as, potentially, every machine has its own optimal epsilon, whose computation is pretty tedious. The percentages refer to the quantile of the series obtained from the distance of all the observations (words) between each other.
- **Tolerance**: Tolerance has actually been different yet depending on the techniques in which it was adopted.

N.A.: Not Applicable, so that such parameter could not, or simply was not, specified in performing a certain algorithm. $(\Box \rightarrow \langle \Box \rangle + \langle \Xi = \langle \Xi \rangle + \langle \Xi = \langle \Xi \rangle + \langle \Xi = \langle$

E. Ronchieri (INFN CNAF, Univ. Bologna)

Hyperparameters Tuning with Grid Search Analysis Model

Algorithm	Dim. re- duction	N. comps	N. clus- ters	Window size	Epsilon	Tolerance
DBSCAN	TRUE,	2, 3,	variable	N.A.	0.05%, 0.1%,	N.A.
	FALSE	10, 20			0.25%, 0.5%	
K-Means	TRUE,	2, 3,	2, 3	N.A.	N.A.	N.A.
	FALSE	10, 20				
PCA De-	implicit	2, 3,	N.A.	N.A.	N.A.	0.75, 0.85, 0.95
composition		10, 20,				(measured as
& Recon-		30				the quantiles
struction						of the error
						vector)
Time-Series	N.A.	N.A.	N.A.	0, 10,	N.A.	2, 3 (the <i>k</i>
Outlier De-				30, 60 (0		number of σ
tection				means no		away from μ in
				window		the outlier de-
				partition)	< □ > < □ > < □	tection_alg.)