

Causal Inference-Based Root Cause Analysis for Online Service Systems with Intervention Recognition

Mingjie Li, Zeyan Li, Kanglin Yin, Xiaohui Nie, Wenchi Zhang, Kaixin Sui, Dan Pei





Background

Analysis Methodology Evaluation Conclusion

Online Service Systems (OSS)



Online Service Systems (OSS)





.



Database

average active session log file sync logic read per second



Storage

disk utilization IO per second IO wait time

.

Online Service Systems (OSS)



Time	09:23	09:24	09:25	09:26	09:27
Value	302	4095	22142	44936	34745

Root Cause Analysis (RCA)





.



Database

average active session log file sync logic read per second



Storage

disk utilization IO per second IO wait time

.

Root Cause Analysis (RCA)





Database

.....

average active session log file sync logic read per second



Storage

disk utilization **IO per second** IO wait time

.....

A fault may propagate in the system

Root Cause Analysis (RCA)



5

Invariant Network-based KDD'16, ICDM'17

$$y(t) = 2x(t-1) + x(t-2) + 1$$

 $x - y$

Invariant Network-based KDD'16, ICDM'17



Invariant Network-based KDD'16, ICDM'17

Random Walk-based CCGRID'18, IWQoS'20, WWW'20, ASE'21



Invariant Network-based KDD'16, ICDM'17

Random Walk-based CCGRID'18, IWQoS'20, WWW'20, ASE'21





Invariant Network-based KDD'16, ICDM'17





Invariant Network-based KDD'16, ICDM'17



Search-based INFOCOM'14, SOC'18, ICSE'21



Invariant Network-based KDD'16, ICDM'17



Search-based INFOCOM'14, SOC'18, ICSE'21



Invariant Network-based KDD'16, ICDM'17



Search-based INFOCOM'14, SOC'18, ICSE'21



Background Causal Inference

Sage (ASPLOS'21) based on Counterfactual The latency is high now. Would the latency recover if CPU utilization decreased?

Background Causal Inference

Sage (ASPLOS'21) based on Counterfactual

How will the latency distribute *if we make CPU utilization low?*

 \mathscr{L}_2 Interventional

How will the latency distribute when CPU utilization is low?

 \mathscr{L}_1 Associational

The latency is high now. Would the latency recover if CPU utilization decreased?

 \mathscr{L}_3 Counterfactual

Judea Pearl's Ladder of Causation

Background **Causal Inference**

Sage (ASPLOS'21) based on Counterfactual

> How will the latency distribute if we make CPU utilization low?

> > \mathscr{L}_{2} , Interventional

How will the latency distribute when CPU utilization is low?

 \mathscr{L}_1 Associational



Judea Pearl's Ladder of Causation

 \mathscr{L}_3 Counterfactual \mathscr{L}_2 Interventional \mathscr{L}_1 Associational



The Ladder of Causation

$$P(\mathbf{V}_{\mathbf{m}} \mid \mathbf{v}')$$

$$V \mid do(\mathbf{m})$$

 \mathscr{L}_3 Counterfactual \mathscr{L}_2 Interventional \mathscr{L}_1 Associational



The Ladder of Causation

$$P(\mathbf{V}_{\mathbf{m}} \mid \mathbf{v}')$$

$$V \mid do(\mathbf{m})$$

 \mathscr{L}_3 Counterfactual \mathscr{L}_2 Interventional \mathscr{L}_1 Associational



M (Manipulation): the intervened variables m: the value of M

The Ladder of Causation

Definition (Intervention Recognition, IR). IR is to recognize **m** from $P(\mathbf{V} \mid do(\mathbf{m}))$ based on $P(\mathbf{V})$.

al



Causal Inference

Causal Bayesian Networks (CBN)



Causal Inference

Causal Bayesian Networks (CBN)

Causal Discovery: Mine the CBN from data



Causal Inference

Causal Bayesian Networks (CBN) Causal Discovery: Mine the CBN from data

work with assumptions which may not fit OSS





consider only a few metrics

which may not fit OSS

Challenges **Observational knowledge is incomplete**



Challenges **Observational knowledge is incomplete**



9

Challenges **Observational knowledge is incomplete**





Causal Bayesian Networks (CBN) Incomplete observational knowledge



Causal Bayesian Networks (CBN) Incomplete observational knowledge




Overview Causal Inference-based Root Cause Analysis (CIRCA) for OSS



Analysis Methodology Evaluation Conclusion



CPU Utilization



CPU Utilization

Intervention Recognition CPU Utilization 98%









$$\mathsf{R} \Rightarrow \mathscr{L}_2$$

$\mathscr{L}_2 \Rightarrow \mathsf{IR?}$

12

* red means the intervention





World 1







World 1







World 1



World 1



Analysis IR is at the second layer of the causal ladder



Assumption (Faithfulness). Any intervention makes an observable change, i.e., $P(V_i | \mathbf{pa}(V_i), do(v_i)) \neq P(V_i | \mathbf{pa}(V_i))$

Analysis IR is at the second layer of the causal ladder



Causal Hierarchy Theorem [1]

[1] Elias Bareinboim, Juan D. Correa, Duligur Ibeling, Thomas Icard. On Pearl's Hierarchy and the Foundations of Causal Inference. Last Revision: Mar, 2021

If we want to answer the question at Layer i, we need knowledge at Layer i or higher.

Analysis IR is at the second layer of the causal ladder



[1] Elias Bareinboim, Juan D. Correa, Duligur Ibeling, Thomas Icard. On Pearl's Hierarchy and the Foundations of Causal Inference. Last Revision: Mar, 2021

IR needs the knowledge of \mathscr{L}_2 (like the CBN [1])

Counterfactual knowledge of \mathscr{L}_3 is unnecessary

Analysis Intervention Recognition Criterion

A given variable is a root cause indicator $V_i \in \mathbf{M}$

Faithfulness

Change in the distribution conditioned on the parents in the CBN $P(V_i | \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i | \mathbf{pa}(V_i))$

Analysis

Intervention Recognition Criterion



[1] Elias Bareinboim, Juan D. Correa, Duligur Ibeling, Thomas Icard. On Pearl's Hierarchy and the Foundations of Causal Inference. Last Revision: Mar, 2021

Change in the distribution conditioned on the parents in the CBN $P(V_i | \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i | \mathbf{pa}(V_i))$ CBN

Analysis

Intervention Recognition Criterion



$V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i \mid \mathbf{pa}(V_i))$

[1] Elias Bareinboim, Juan D. Correa, Duligur Ibeling, Thomas Icard. On Pearl's Hierarchy and the Foundations of Causal Inference. Last Revision: Mar, 2021

Change in the distribution conditioned on the parents in the CBN $P(V_i | \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i | \mathbf{pa}(V_i))$ CBN

Nethodology

Background Analysis Methodology Evaluation Conclusion

Structural Graph Construction

Meta Metrics

Structural Graph Construction

Meta Metrics

Input



Time

17

Structural Graph Construction

Meta Metrics

Traffic (T)

Latency (L)

[2] Betsy Beyer, Chris Jones, Jennifer Petoff, and Niall Richard Murphy. Site Reliability Engineering (first ed.). O'Reilly Media, Inc. 2016.



Structural Graph Construction Meta Metrics



Assign directions among meta metrics as causal assumptions

Structural Graph Construction Skeleton with Architecture Extension



Structural Graph Construction Skeleton with Architecture Extension



Structural Graph Construction Monitoring Metric Plugging-in



Structural Graph Construction Monitoring Metric Plugging-in



Core Idea: $V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i \mid \mathbf{pa}(V_i))$

A few faulty data are available

Core Idea: $V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid \mathbf{pa}(V_i), do(\mathbf{m})) \neq P(V_i \mid \mathbf{pa}(V_i))$

A few faulty data are available

Core Idea:
$$V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid$$

Hypothesis Testing

$\mathbf{H}_0(V_i \notin \mathbf{M})$:

The need for faulty data is reduced



Core Idea:
$$V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid$$

Hypothesis Testing

$\mathbf{H}_0(V_i \notin \mathbf{M})$:

The need for faulty data is reduced



Core Idea:
$$V_i \in \mathbf{M} \Leftrightarrow P(V_i \mid$$

Hypothesis Testing

$\mathbf{H}_0(V_i \notin \mathbf{M})$:

The need for faulty data is reduced



 \sum

$$\mathcal{N}(expectation, \sigma_{residuals})$$

Descendant Adjustment Alleviate the bias introduced in hypothesis testing



Descendant Adjustment Alleviate the bias introduced in hypothesis testing



Descendant Adjustment Alleviate the bias introduced in hypothesis testing



Intuition: A variable may point to an actionable mitigation method more likely than its descendants

Evaluation

Background Analysis Methodology Evaluation Conclusion
Hyperparameters

a fault is detected

 t_d









Experimental Setup Evaluation Metrics

• Recall with the top-k results

•
$$AC@k = \frac{1}{|\mathcal{F}|} \sum_{\mathbf{M} \in \mathcal{F}} \frac{|\mathbf{M} \cap \{R\}\|}{|\mathbf{M} \in \mathcal{F}|}$$

•
$$k \leq K = 5$$

$R_i(\mathbf{M}) \mid i = 1, 2, \dots, k\}$

Simulation Study Data Generation

•
$$\mathbf{x}^{(t)} = \mathbf{A}\mathbf{x}^{(t)} + \beta \mathbf{x}^{(t-1)} + \epsilon^{(t)}$$

node (service level indicator) having no children.

	#Node	#Edge	#Graph	#Case/Graph
\mathcal{D}_{Sim}^{50}	50	100		
\mathscr{D}^{100}_{Sim}	100	500	10	100
\mathscr{D}_{Sim}^{500}	500	5,000		

• A encodes the CBN, enforced to be a connected DAG with only the first

• RHT-PG: RHT with the perfect graph

•
$$\mathbf{Pa}(X_i^{(t)}) = \mathbf{Pa}^{(t)}(X_i) \cup \left\{ X_i^{(t-1)} \right\}$$

\mathscr{D}_{Sim}^{50}				
AC@1	AC@5	T (s)		
0.432	0.733	0.306		
0.508	0.761	6.601		
0.541	0.682	0.308		
0.515	0.682	0.502		
0.178	0.217	0.501		
0.188	0.433	0.714		
0.188	0.433	0.437		
0.116	0.278	0.624		
0.074	0.223	4.844		
0.598	0.880	0.338		
0.615	0.952	0.346		
0.617	0.999			
	AC@10.4320.5080.5410.5410.5150.1780.1880.1880.1880.16150.5980.617	𝔅𝔅AC@1AC@50.4320.7330.5080.7610.5150.6820.1780.2170.1880.4330.1880.4330.1160.2780.0740.2230.5980.8800.6150.999		



• RHT-PG: RHT with the perfect graph

•
$$\mathbf{Pa}(X_i^{(t)}) = \mathbf{Pa}^{(t)}(X_i) \cup \left\{ X_i^{(t-1)} \right\}$$

- Takeaways
 - RHT has theoretical reliability.

\mathscr{D}_{Sim}^{50}				
AC@1	AC@5	T (s)		
0.432	0.733	0.306		
0.508	0.761	6.601		
0.541	0.682	0.308		
0.515	0.682	0.502		
0.178	0.217	0.501		
0.188	0.433	0.714		
0.188	0.433	0.437		
0.116	0.278	0.624		
0.074	0.223	4.844		
0.598	0.880	0.338		
0.615	0.952	0.346		
0.617	0.999			
	AC@10.4320.5080.5410.5410.5150.1780.1880.1880.1880.16150.5980.617	𝔅𝔅AC@1AC@50.4320.7330.5080.7610.5150.6820.1780.2170.1880.4330.1880.4330.1160.2780.0740.2230.5980.8800.6150.999		



• RHT-PG: RHT with the perfect graph

•
$$\mathbf{Pa}(X_i^{(t)}) = \mathbf{Pa}^{(t)}(X_i) \cup \left\{ X_i^{(t-1)} \right\}$$

- Takeaways
 - RHT has theoretical reliability.
 - A broken CBN cannot guarantee a correct answer to RCA.

Scoring	\mathscr{D}_{Sim}^{50}				
Method	AC@1	AC@5	T (s)		
NSigma	0.432	0.733	0.306		
SPOT	0.508	0.761	6.601		
DFS	0.541	0.682	0.308		
DFS-MS	0.515	0.682	0.502		
DFS-MH	0.178	0.217	0.501		
RW-Par	0.188	0.433	0.714		
RW-2	0.188	0.433	0.437		
ENMF	0.116	0.278	0.624		
CRD	0.074	0.223	4.844		
RHT	0.598	0.880	0.338		
RHT-PG	0.615	0.952	0.346		
Ideal	0.617	0.999			

26



• RHT-PG: RHT with the perfect graph

•
$$\mathbf{Pa}(X_i^{(t)}) = \mathbf{Pa}^{(t)}(X_i) \cup \left\{ X_i^{(t-1)} \right\}$$

- Takeaways
 - RHT has theoretical reliability.
 - A broken CBN cannot guarantee correct answer to RCA.
 - There may be statistical errors due to limited faulty data.

	2	
	٨	

Scoring		\mathcal{D}_{Sim}^{50}	
Method	AC@1	AC@5	T (s)
NSigma	0.432	0.733	0.306
SPOT	0.508	0.761	6.601
DFS	0.541	0.682	0.308
DFS-MS	0.515	0.682	0.502
DFS-MH	0.178	0.217	0.501
RW-Par	0.188	0.433	0.714
RW-2	0.188	0.433	0.437
ENMF	0.116	0.278	0.624
CRD	0.074	0.223	4.844
RHT	0.598	0.880	0.338
RHT-PG	0.615	0.952	0.346
Ideal	0.617	0.999	



- Classify faults based on the change of the root cause metrics when the service level indicator is abnormal
 - Weak: dramatically
 - **Strong**: slight
 - Mixed: both



- Classify faults based on the change of the root cause metrics when the service level indicator is abnormal
 - Weak: dramatically
 - **Strong**: slight
 - Mixed: both



Weak fault



- Classify faults based on the change of the root cause metrics when the service level indicator is abnormal
 - Weak: dramatically
 - Strong: slight
 - Mixed: both



Strong fault

Service Level Indicator



- Classify faults based on the change of the root cause metrics when the service level indicator is abnormal
 - Weak: dramatically
 - **Strong**: slight
 - Mixed: both
- Takeaway
 - RHT is more robust.

Scoring	Weak	(916)	Mixe	d (64)	Stron	g
Method	AC@1	AC@5	AC@1	AC@5	AC@1	ļ
NSigma	0.454	0.753	0.249	0.498	0.000	(
SPOT	0.534	0.783	0.293	0.503	0.000	(
DFS	0.558	0.707	0.282	0.368	0.550	(
DFS-MS	0.531	0.707	0.277	0.368	0.550	(
DFS-MH	0.184	0.223	0.069	0.123	0.250	(
RW-Par	0.194	0.445	0.142	0.300	0.050	(
RW-2	0.194	0.445	0.142	0.300	0.050	(
ENMF	0.111	0.269	0.124	0.321	0.300	(
CRD	0.071	0.207	0.088	0.353	0.150	(
RHT	0.613	0.888	0.325	0.730	0.800	
RHT-PG	0.624	0.954	0.358	0.914	1.000	
Ideal	0.627	1.000	0.358	0.995	1.000	



- Dataset
 - a large banking system
- Implementation
 - - denoted as Structural
 - Equip RHT with descendant adjustment
 - denoted as CIRCA

99 faults with high Average Active Sessions (AAS) from Oracle databases in

• Our structural graph contains 197 monitoring metrics with 2,641 edges



Empirical Study on Oracle Database Data Performance Evaluation

- Takeaways
 - CIRCA outperforms baselines.

Scoring Method	Graph Method	AC@1	AC@5	T (s)
NSigma	Empty	0.323	0.662	0.472
SPOT	Empty	0.152	0.419	5.027
DFS	Structural	0.187	0.313	0.483
DFS-MS	Structural	0.207	0.308	0.839
DFS-MH	Structural	0.268	0.439	0.844
RW-Par	PCTS	0.086	0.449	24.695
RW-2	PCTS	0.086	0.449	24.559
ENMF	Empty	0.111	0.374	0.771
CRD	Empty	0.035	0.313	4.787
CIRCA	Structural	0.404	0.763	0.578
ld	eal	0.929	1.000	



- Takeaways
 - CIRCA outperforms baselines.
 - Each of the 3 proposed techniques has a positive effect.
 - Search-based methods also benefit from the proposed structural graph.









RHT confronts incomplete observational knowledge



RHT confronts incomplete observational knowledge

Descendant adjustment helps CIRCA rank LFS ahead



RHT confronts incomplete observational knowledge

• Further advancement should handle this challenge

Descendant adjustment helps **CIRCA** rank LFS ahead





RHT confronts incomplete observational knowledge

• Further advancement should handle this challenge

Descendant adjustment helps **CIRCA** rank LFS ahead

- CIRCA outperforms pure RHT in this case
- Descendant adjustment needs more verification





Conclusion

Background Analysis Methodology Evaluation Conclusion

Contributions Formulate RCA as a causal inference task



The Ladder of Causation

Definition (Intervention Recognition, IR). IR is to recognize **m** from $P(\mathbf{V} \mid do(\mathbf{m}))$ based on $P(\mathbf{V})$.



Contributions Causal Inference-based Root Cause Analysis (CIRCA) for OSS



Causal Inference

Causal Bayesian Networks (CBN) Incomplete observational knowledge

Contributions Causal Inference-based Root Cause Analysis (CIRCA) for OSS



ContributionsEvaluation with both simulation

Simulation Study

Empirical Study

Evaluation with both simulation and real-world datasets

Contributions Evaluation with both simulation and real-world datasets

Simulation Study

Empirical Study





Causal Inference-based Root Cause Analysis

Thanks for listening

https://github.com/NetManAlOps/CIRCA

