

数据挖掘漫游

Mathew Shen(datahonor@gmail.com)

2023 年 8 月 3 日

1 Security

2 AIOps

- TSAD: Algorithm
- TSAD: System

3 QA

Security

帐号安全

Situation

- 自动机模拟登录验密
 - 暴力破解，拖/洗/撞库

帐号安全

Situation

- 自动机模拟登录验密
 - 暴力破解，拖/洗/撞库

Task

- 保证正常验密请求的前提下拦截黑产验密请求
 - 降低盗号率及其他衍生指标

帐号安全

Situation

- 自动机模拟登录验密
 - 暴力破解，拖/洗/撞库

Action

- FP-Growth(PFP[11]) + Spark
 - 基于请求特征挖掘关联规则，用于实时封禁

Task

- 保证正常验密请求的前提下拦截黑产验密请求
 - 降低盗号率及其他衍生指标

帐号安全

Situation

- 自动机模拟登录验密
 - 暴力破解，拖/洗/撞库

Action

- FP-Growth(PFP[11]) + Spark
 - 基于请求特征挖掘关联规则，用于实时封禁

Task

- 保证正常验密请求的前提下拦截黑产验密请求
 - 降低盗号率及其他衍生指标

Result

- 挖掘出来的规则符合专家经验
 - 拦截相当部分黑产验密请求

风控漫谈

查杀分离：信息差

对能够很好定位黑产的特征应给予保护，不直接用于账号封禁等

风控漫谈

查杀分离: 信息差

对能够很好定位黑产的特征应给予保护, 不直接用于账号封禁等

对抗的关键: ROI

提高攻击成本, 降低攻击 ROI 是关键

AIOps

AIOps 概览

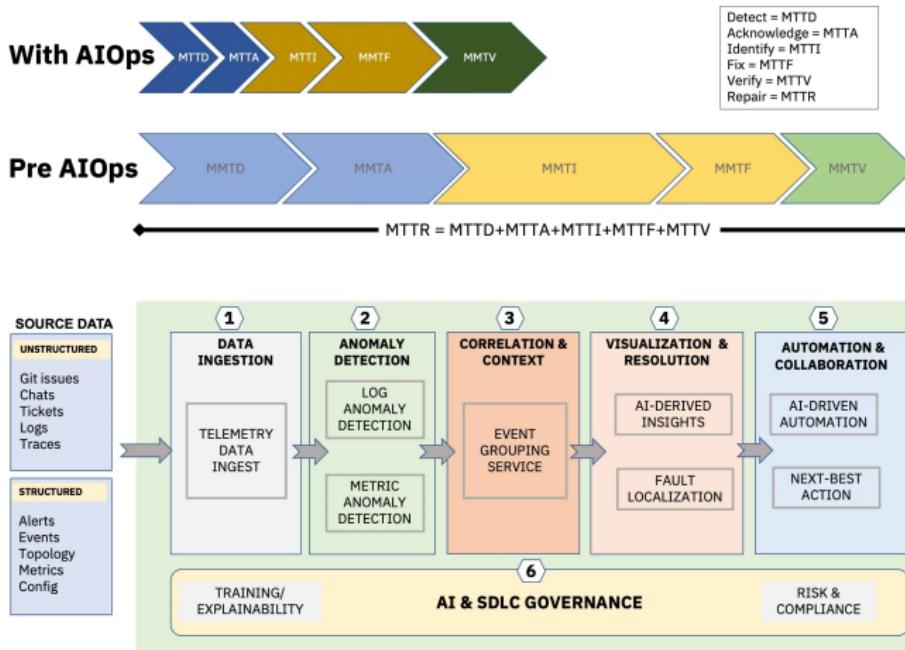


图 1: AIOps: A Path to Reliability at Cloud Scale¹

¹AIOps: A Path to Reliability at Cloud Scale

TSAD

Situation

- 大量指标需要监控以及时发现系统问题
 - 人工阈值的方法难以自适应，告警噪声大

TSAD

Situation

- 大量指标需要监控以及时发现系统问题
 - 人工阈值的方法难以自适应，告警噪声大

Task

- 检测各类时间序列存在的异常，及时发出预警信息
 - Precision, Recall, MTTD

TSAD

Situation

- 大量指标需要监控以及时发现系统问题
 - 人工阈值的方法难以自适应，告警噪声大

Action

- 构建算法库系统地检测各类时序异常
 - 基于算法库搭建 TSAD 系统来解决业务问题

Task

- 检测各类时间序列存在的异常，及时发出预警信息
 - Precision, Recall, MTTD

TSAD

Situation

- 大量指标需要监控以及时发现系统问题
 - 人工阈值的方法难以自适应，告警噪声大

Action

- 构建算法库系统地检测各类时序异常
 - 基于算法库搭建 TSAD 系统来解决业务问题

Task

- 检测各类时间序列存在的异常，及时发出预警信息
 - Precision, Recall, MTTD

Result

- Precision > 80%
 - Recall ≈ 100%
 - MTTD (< 3 Min)

TSAD: 算法/模型

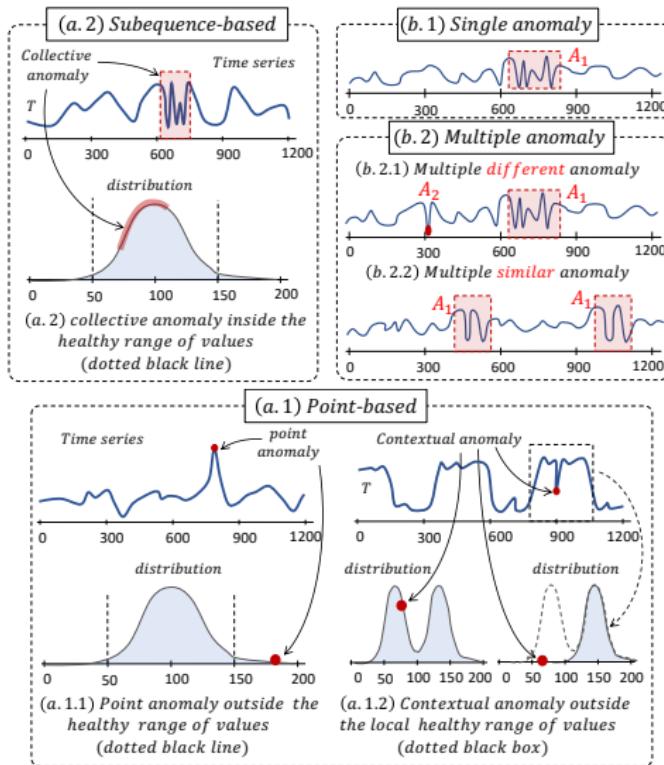


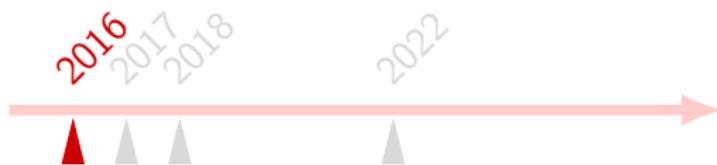
图 2: New Trends in Time-Series Anomaly Detection[4]

TSAD: 算法/模型

A Comprehensive Evaluation on TSAD[15]

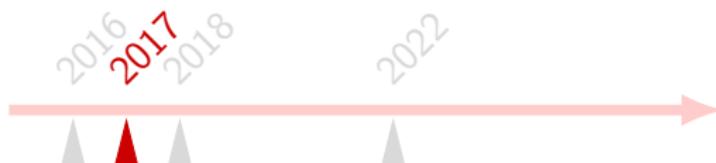
- ① Deep learning approaches are not (yet) competitive despite their higher processing effort on training data
 - ② There is no one-size-fits-all solution in the set of currently available algorithms...there is no clear winner
 - ③ Simple methods yield performance almost as good as more sophisticated methods
 - ④ Every practical algorithm deployment needs careful testing
 - ⑤ Anomalies on periodic time series are easier to detect than on non-periodic time series

TSAD: 算法/模型



- 2016: Amazon: RRCF[2]

TSAD: 算法/模型



- 2016: Amazon: RRCF[2]
 - 2017: Twitter: H-S-ESD[9]
 - 2017: IRISA: SPOT[16]

TSAD: 算法/模型



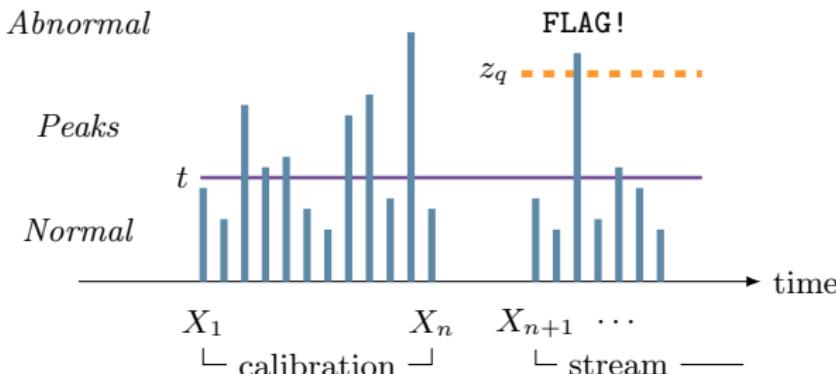
- 2016: Amazon: RRCF[2]
 - 2017: Twitter: H-S-ESD[9]
 - 2017: IRISA: SPOT[16]
 - 2018: Facebook(Meta): Prophet[17]

TSAD: 算法/模型



- 2016: Amazon: RRCF[2]
 - 2017: Twitter: H-S-ESD[9]
 - 2017: IRISA: SPOT[16]
 - 2018: Facebook(Meta): Prophet[17]
 - 2021: Huawei: FluxEV[12]

TSAD: 算法/模型



Algorithm 1 POT (Peaks-over-Threshold)

```

1: procedure POT( $X_1, \dots, X_n, q$ )
2:    $t \leftarrow \text{SETINITIALTHRESHOLD}(X_1, \dots, X_n)$ 
3:    $\mathbf{Y}_t \leftarrow \{X_i - t \mid X_i > t\}$ 
4:    $\hat{\gamma}, \hat{\sigma} \leftarrow \text{GRIMSHAW}(\mathbf{Y}_t)$ 
5:    $z_q \leftarrow \text{CALCTHRESHOLD}(q, \hat{\gamma}, \hat{\sigma}, n, N_t, t)$ 
6:   return  $z_q, t$ 
7: end procedure

```

TSAD: 算法 / 模型

Algorithm 2 SPOT (Streaming POT)

```

1: procedure SPOT( $(X_i)_{i>0}, n, q$ )
2:    $\mathbf{A} \leftarrow \emptyset$                                  $\triangleright$  set of the anomalies
3:    $z_q, t \leftarrow \text{POT}(X_1, \dots X_n, q)$ 
4:    $k \leftarrow n$ 
5:   for  $i > n$  do
6:     if  $X_i > z_q$  then                             $\triangleright$  anomaly case
7:       Add  $(i, X_i)$  in  $\mathbf{A}$ 
8:     else if  $X_i > t$  then                       $\triangleright$  real peak case
9:        $Y_i \leftarrow X_i - t$ 
10:      Add  $Y_i$  in  $\mathbf{Y}_t$ 
11:       $N_t \leftarrow N_t + 1$ 
12:       $k \leftarrow k + 1$ 
13:       $\hat{\gamma}, \hat{\sigma} \leftarrow \text{GRIMSHAW}(\mathbf{Y}_t)$ 
14:       $z_q \leftarrow \text{CALCTHRESHOLD}(q, \hat{\gamma}, \hat{\sigma}, k, N_t, t)$ 
15:    else                                          $\triangleright$  normal case
16:       $k \leftarrow k + 1$ 
17:    end if
18:  end for
19: end procedure

```

TSAD: 算法/模型

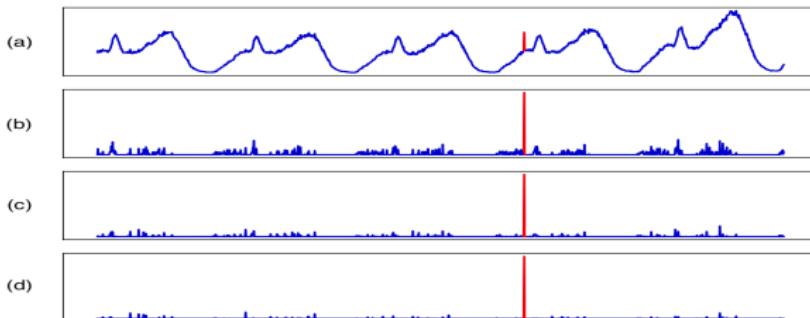


Figure 6: Smoothing. The first row (a) shows the raw time-series curves and anomalies are highlighted in red; (b) stands for originally extracted fluctuations; (c)-(d) represent the fluctuation features after the first and second smoothing respectively.

图 5: Huawei: FluxEV[12]

TSAD: 算法 / 模型

Algorithm 1: ExtAndSmooth

Input: input data $X = [X_1, \dots, X_n]$, window sizes s, p, d ,
period l

Output: $E = [E_1, \dots, E_n]$, $F = [F_1, \dots, F_n]$, $S = [S_1, \dots, S_n]$
 $M = [M_1, \dots, M_{n-d}]$

```

1 for  $i = 1$  to  $n$  do
2    $E_i = F_i = M_i = S_i = \text{None};$ 
3   if  $i > s$  then
4      $E_i = X_i - \text{EWMA}(X_{i-s, i-1});$ 
5   if  $i > 2s$  then
6      $\Delta\sigma = \sigma(E_{i-s, i}) - \sigma(E_{i-s, i-1});$ 
7      $F_i = \max(\Delta\sigma, 0);$ 
8   if  $i > 2s + 2d$  then
9      $M_{i-d} = \max(F_{i-2d, i});$ 
10  if  $i > 2s + d + l(p - 1)$  then
11     $\Delta F_i = F_i - \max(M_{i-l(p-1)}, \dots, M_{i-2l}, M_{i-l})$ 
12     $S_i = \max(\Delta F_i, 0);$ 

```

TSAD-System: 工业界开源系统



- 2017: Twitter: H-S-ESD[9]

TSAD-System: 工业界开源系统



- 2017: Twitter: H-S-ESD[9]
- 2020: Alibaba(DAMO): RobustX[19][20][7][21]
- 2020: Amazon: GluonTS[1]
- 2020: Zillow: Luminaire[5]
- 2020: MicroSoft: Auto-Selector[22]

TSAD-System: 工业界开源系统



- 2017: Twitter: H-S-ESD[9]
- 2020: Alibaba(DAMO): RobustX[19][20][7][21]
- 2020: Amazon: GluonTS[1]
- 2020: Zillow: Luminaire[5]
- 2020: MicroSoft: Auto-Selector[22]
- 2021: Salesforce: Merlion[3]
- 2021: LinkedIn: Silverkite[10]

TSAD-System: 工业界开源系统



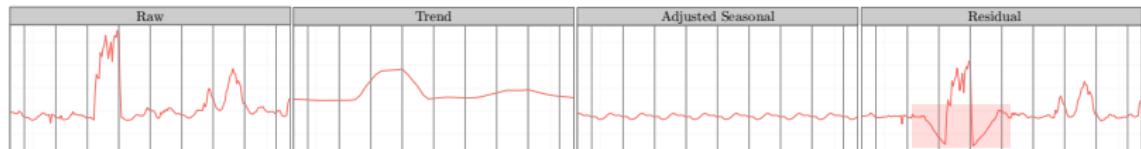
- 2017: Twitter: H-S-ESD[9]
- 2020: Alibaba(DAMO): RobustX[19][20][7][21]
- 2020: Amazon: GluonTS[1]
- 2020: Zillow: Luminaire[5]
- 2020: MicroSoft: Auto-Selector[22]
- 2021: Salesforce: Merlion[3]
- 2021: LinkedIn: Silverkite[10]
- 2022: IBM: AnomalyKiTS[14]
- 2022: MicroSoft: HEAT-RL[18]

TSAD-System: 工业界开源系统

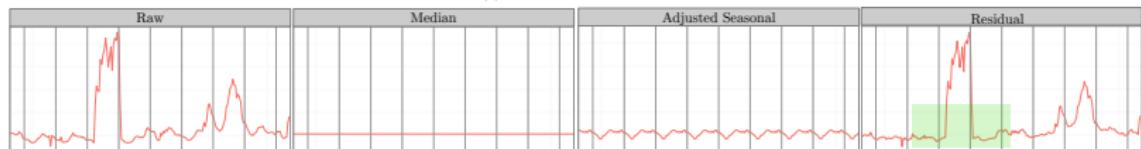


- 2017: Twitter: H-S-ESD[9]
- 2020: Alibaba(DAMO): RobustX[19][20][7][21]
- 2020: Amazon: GluonTS[1]
- 2020: Zillow: Luminaire[5]
- 2020: MicroSoft: Auto-Selector[22]
- 2021: Salesforce: Merlion[3]
- 2021: LinkedIn: Silverkite[10]
- 2022: IBM: AnomalyKiTS[14]
- 2022: MicroSoft: HEAT-RL[18]
- 2023: Amazon: Model Selection[8]

TSAD-System: H-S-ESD(Twitter, 2017)



(a) STL with Trend Removal



(b) STL with Median Removal

图 7: Twitter: ESD[9]

TSAD-System: H-S-ESD(Twitter, 2017)

Algorithm 1 S-ESD Algorithm

Input:

X = A time series

n = number of observations in X

k = max anomalies (iterations in ESD)

Output:

X_A = An anomaly vector wherein each element is a tuple
(*timestamp, observed value*)

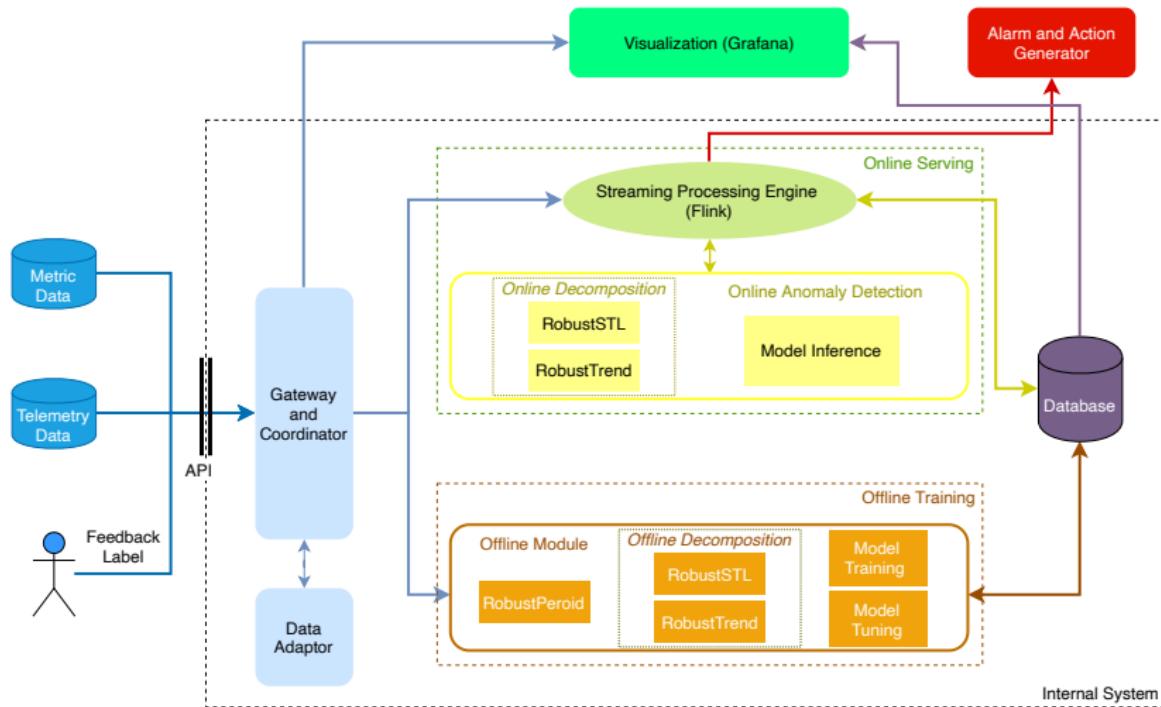
Require:

$$k \leq (n \times .49)$$

1. Extract seasonal component S_X using STL Variant
2. Compute median \tilde{X}
- /* Compute residual */
3. $R_X = X - S_X - \tilde{X}$
- /* Detect anomalies vector X_A using ESD */
4. $X_A = \text{ESD}(R, k)$

return X_A

TSAD-System: RobustX(Alibaba, 2020)



TSAD-System: Merlion(Salesforce, 2021)

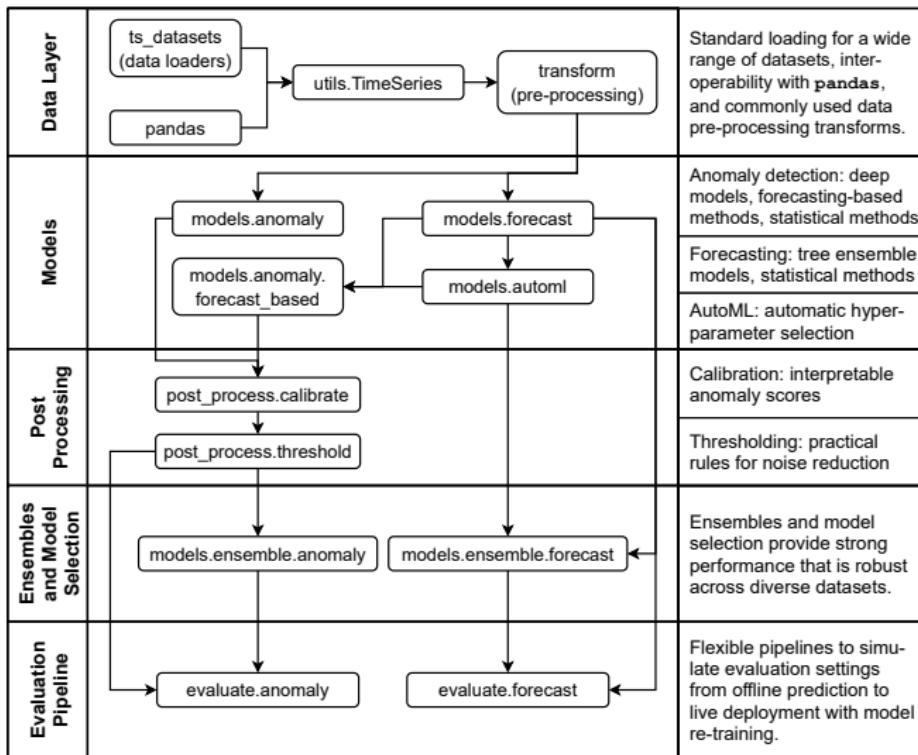
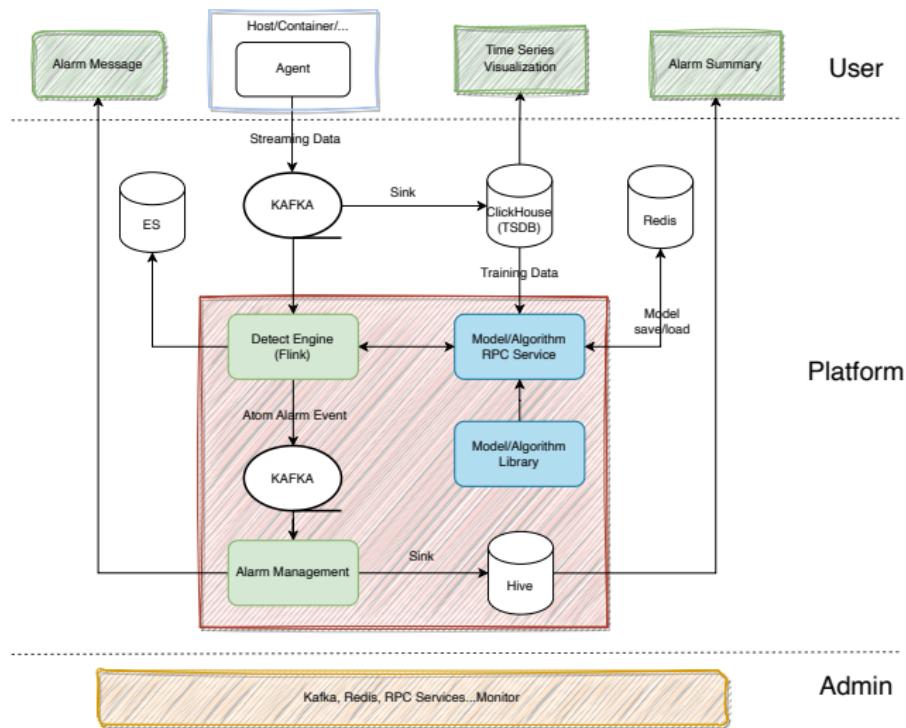


图 10: Salesforce: Merlion[3]

TSAD: 系统架构



Admin

TSAD: 系统性能问题与解决

Redis: 模型大 Key, 无法保证 SLA

优化算法实现, 降低模型内用到 10KB 以下

TSAD: 系统性能问题与解决

Redis: 模型大 Key, 无法保证 SLA

优化算法实现, 降低模型内用到 10KB 以下

Flink: 数据乱序

SlidingEventTimeWindows, 修改算法使能接受一段时间序列的输入, 并记录时间状态

TSAD: 系统性能问题与解决

Redis: 模型大 Key, 无法保证 SLA

优化算法实现, 降低模型内用到 10KB 以下

Flink: 数据乱序

SlidingEventTimeWindows, 修改算法使能接受一段时间序列的输入, 并记录时间状态

RPC: 序列化与反序列化开销过高

通过火焰图发现在进行 RPC 通讯时, Python 对象在序列化和反序列化时耗时较多

RCA: 根因分析

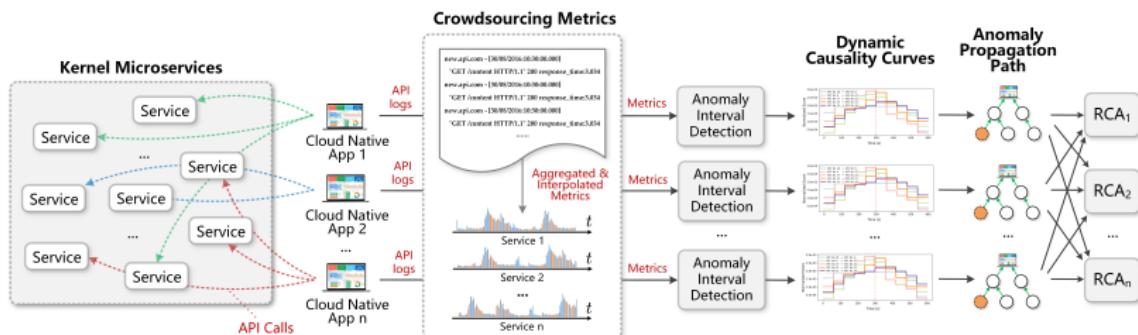


图 12: PKU: DyCause[13][24]

AI4X: 概览

内存故障预测

通过内存条中 CE 的模式提前预测 UE 的发生, 进而降低宕机率

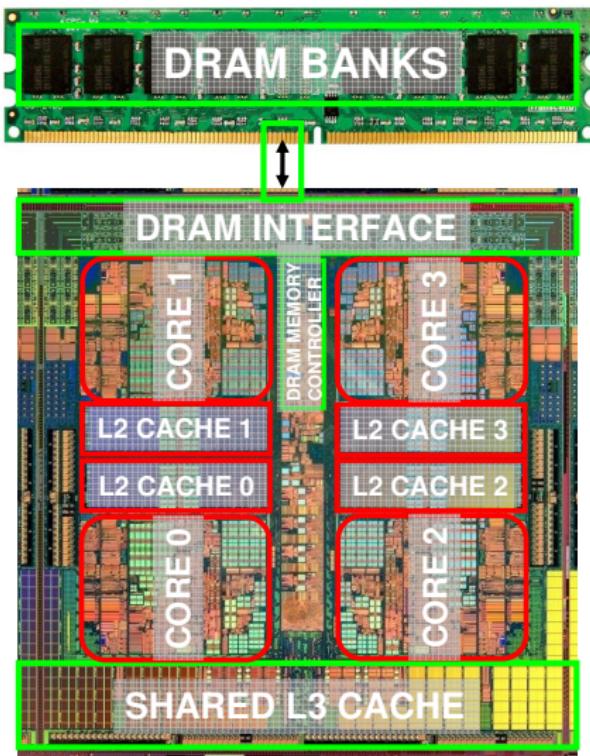
硬盘故障预测

通过硬盘 SMART Log 数据提前预测硬盘故障, 避免数据损坏

数据中心智能控制

结合更多业务侧的数据, 升级数据中心控制系统

AI4X: 内存故障预测



AI4X: 内存故障预测

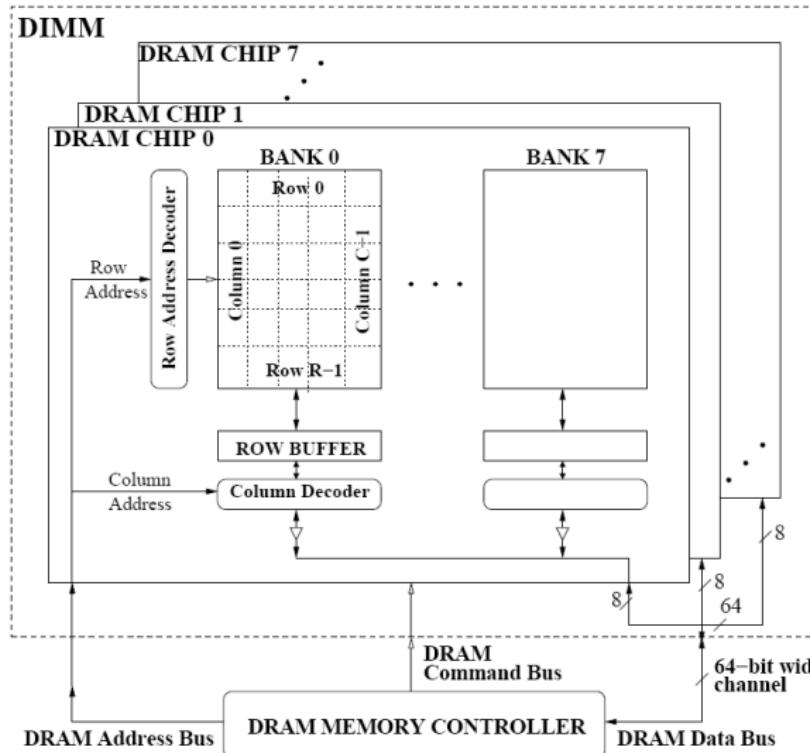


图 14: From CMU ECE740

AI4X: 内存故障预测

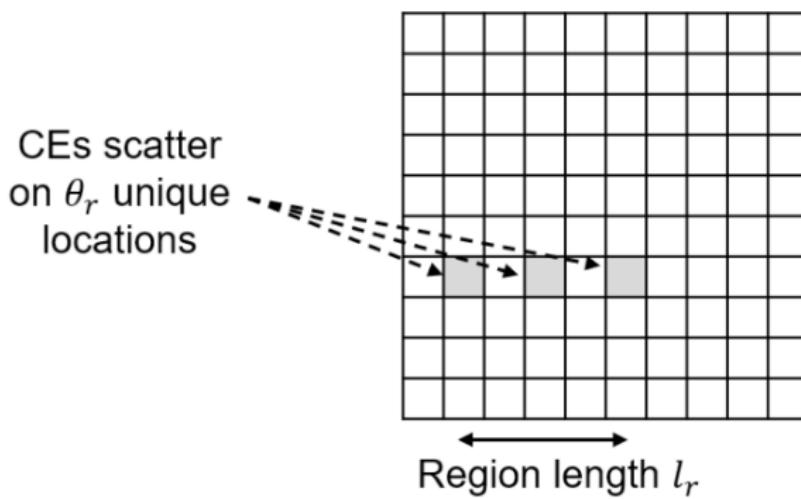


图 15: Intel, Bytedance: MFP[6]

AI4X: 硬盘故障预测

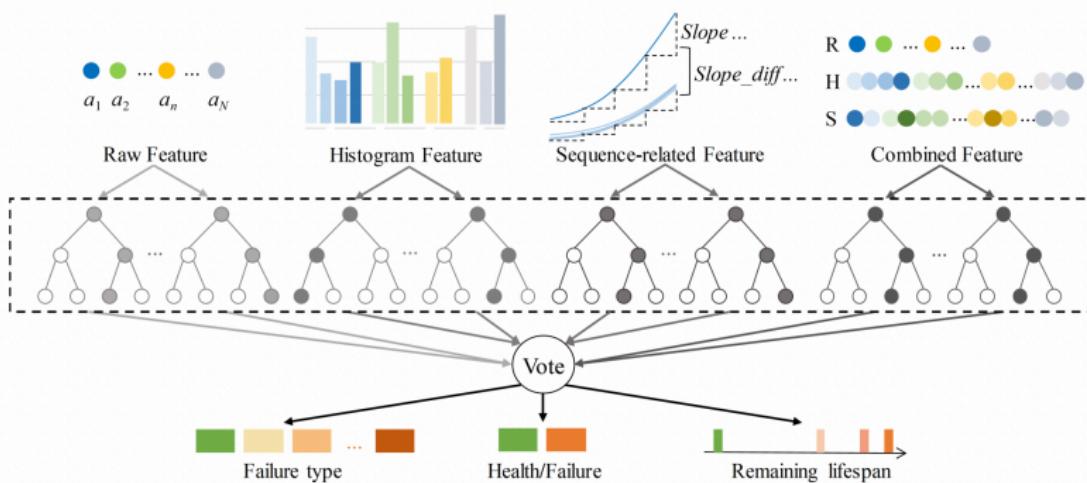


图 16: Samsung, Tencent: DFP(FAST'23)[23]

AI4X: 数据中心智能控制

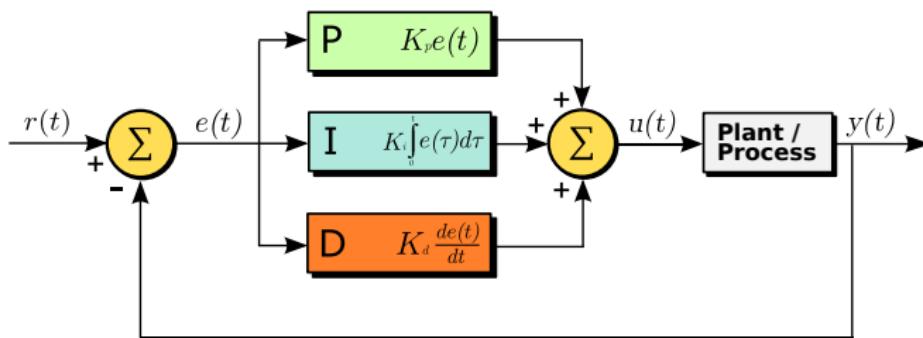


图 17: Wikipedia: PID controller²

²PID controller

AI4X: 数据中心智能控制

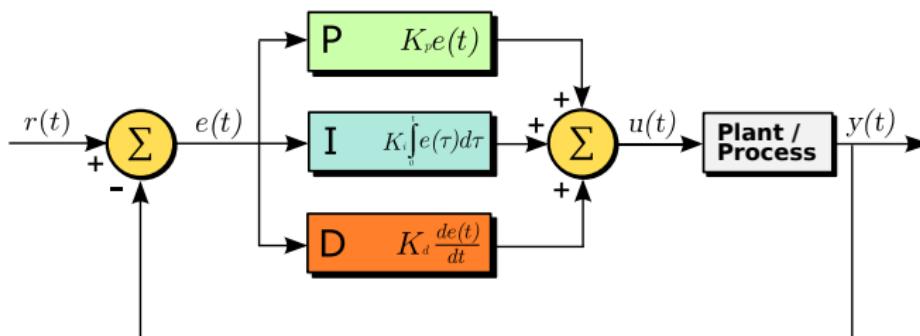


图 17: Wikipedia: PID controller²

PID -> MPC

机房内的机器的负荷预测可以提供很有利的信息, 借助 MPC 可以充分利用这一信息进行更加精准的温度调控 (字节液冷温控专利 [25])

²PID controller

参考文献 |

- [1] Alexander Alexandrov et al. "GluonTS: Probabilistic and Neural Time Series Modeling in Python.". In: *J. Mach. Learn. Res.* 21.116 (2020), pp. 1–6.
- [2] Matthew D Bartos, Abhiram Mullapudi, and Sara C Troutman. "rrcf: Implementation of the robust random cut forest algorithm for anomaly detection on streams". In: *Journal of Open Source Software* 4.35 (2019), p. 1336.
- [3] Aadyot Bhatnagar et al. "Merlion: A machine learning library for time series". In: *arXiv preprint arXiv:2109.09265* (2021).
- [4] Paul Boniol, John Paparrizos, and Themis Palpanas. "New Trends in Time-Series Anomaly Detection". In: (2023).
- [5] Sayan Chakraborty et al. "Building an automated and self-aware anomaly detection system". In: *2020 IEEE International Conference on Big Data (Big Data)*. IEEE. 2020, pp. 1465–1475.

参考文献 II

- [6] Xiaoming Du et al. "Fault-aware prediction-guided page offlining for uncorrectable memory error prevention". In: *2021 IEEE 39th International Conference on Computer Design (ICCD)*. IEEE. 2021, pp. 456–463.
- [7] Jingkun Gao et al. "Robusttad: Robust time series anomaly detection via decomposition and convolutional neural networks". In: *arXiv preprint arXiv:2002.09545* (2020).
- [8] Mononito Goswami et al. "Unsupervised Model Selection for Time-series Anomaly Detection". In: *arXiv preprint arXiv:2210.01078* (2022).
- [9] Jordan Hochenbaum, Owen S Vallis, and Arun Kejariwal. "Automatic anomaly detection in the cloud via statistical learning". In: *arXiv preprint arXiv:1704.07706* (2017).
- [10] Reza Hosseini et al. "A flexible forecasting model for production systems". In: *arXiv preprint arXiv:2105.01098* (2021).
- [11] Haoyuan Li et al. "Pfp: parallel fp-growth for query recommendation". In: *Proceedings of the 2008 ACM conference on Recommender systems*. 2008, pp. 107–114.

参考文献 III

- [12] Jia Li et al. "FluxEV: a fast and effective unsupervised framework for time-series anomaly detection". In: *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 2021, pp. 824–832.
- [13] Yicheng Pan et al. "Faster, deeper, easier: crowdsourcing diagnosis of microservice kernel failure from user space". In: *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 2021, pp. 646–657.
- [14] Dhaval Patel et al. "Anomalykits: Anomaly detection toolkit for time series". In: *AAAI*. AAAI Press. 2022.
- [15] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. "Anomaly detection in time series: a comprehensive evaluation". In: *Proceedings of the VLDB Endowment* 15.9 (2022), pp. 1779–1797.
- [16] Alban Siffer et al. "Anomaly detection in streams with extreme value theory". In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2017, pp. 1067–1075.

参考文献 IV

- [17] Sean J Taylor and Benjamin Letham. "Forecasting at scale". In: *The American Statistician* 72.1 (2018), pp. 37–45.
- [18] Yujing Wang et al. "Heat-RL: Online Model Selection for Streaming Time-Series Anomaly Detection". In: *Conference on Lifelong Learning Agents*. PMLR. 2022, pp. 767–777.
- [19] Qingsong Wen et al. "RobustPeriod: Time-frequency mining for robust multiple periodicities detection". In: *arXiv preprint arXiv:2002.09535* (2020).
- [20] Qingsong Wen et al. "RobustSTL: A robust seasonal-trend decomposition algorithm for long time series". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 01. 2019, pp. 5409–5416.
- [21] Qingsong Wen et al. "RobustTrend: A Huber loss with a combined first and second order difference regularization for time series trend filtering". In: *arXiv preprint arXiv:1906.03751* (2019).
- [22] Yuanxiang Ying et al. "Automated Model Selection for Time-Series Anomaly Detection". In: *arXiv preprint arXiv:2009.04395* (2020).

参考文献 V

- [23] Yuqi Zhang et al. "Multi-view Feature-based {SSD} Failure Prediction: What, When, and Why". In: *21st USENIX Conference on File and Storage Technologies (FAST 23)*. 2023, pp. 409–424.
 - [24] 张宇 et al. "一种主机故障的检测方法、装置、电子设备及存储介质". Chinese. Pat. CN116414653A. July 11, 2023.
 - [25] 陈昱 et al. "用于温度控制的方法、装置、设备和存储介质". May 2023.

QA