

# Akane: Perplexity-Guided Time Series Data Cleaning

2024.06.11

Xiaoyu Han,  
Haoran Xiong,  
Zhenying He,  
Peng Wang,  
Chen Wang,  
X. Sean Wang

- Time series analysis

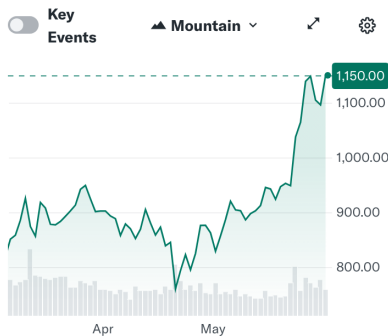
- Finance, Meteorology, Industry, Medicine, etc.

NasdaqGS - Nasdaq Real Time Price - USD

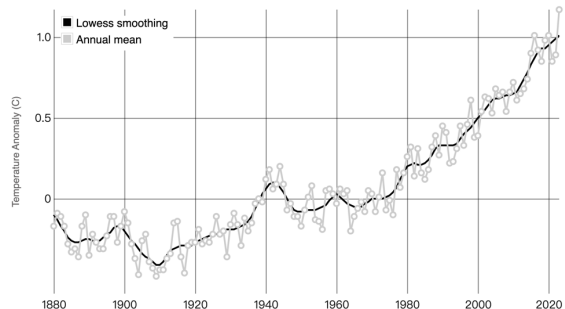
**NVIDIA Corporation (NVDA)**

**1,150.00** +53.67 (+4.90%)

At close: 4:00 PM EDT



Yahoo Finance

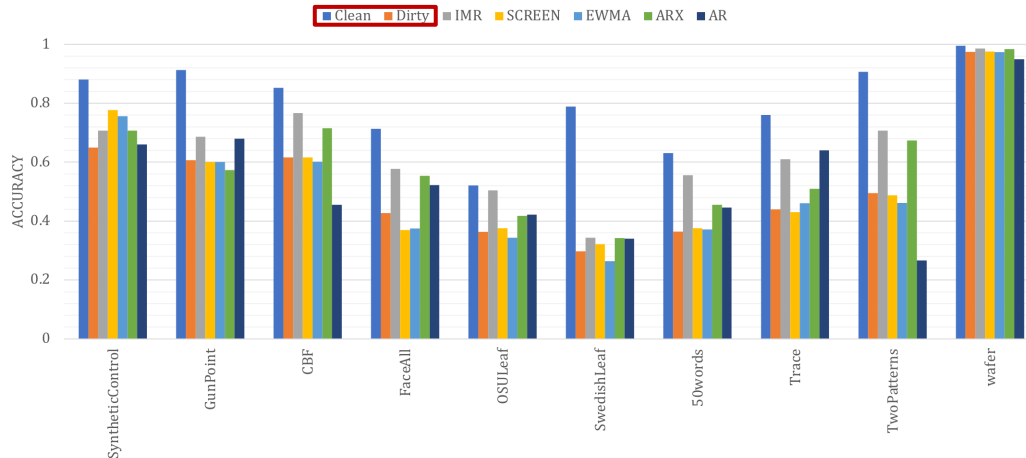


NASA Global Temperature

- Premise of reliable analysis

  - Good data quality

-- Data quality issues have been widely recognized in IoT data, and prevent the downstream applications. [CIKM'20]



- Time series reliability issues
  - Unreliable devices
  - Malfunctioning metering systems
  - Confuse of stock symbols

- Simply discarding might not work
  - Sometimes data are precious
  - Make data even more incomplete [CIKM'20]
  - Mislead downstream applications [CIKM'20]

- Simply discarding might not work
  - Sometimes data are precious
  - Make data even more incomplete [CIKM'20]
  - Mislead downstream applications [CIKM'20]



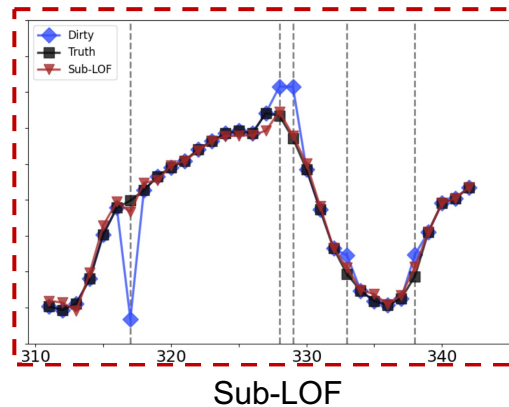
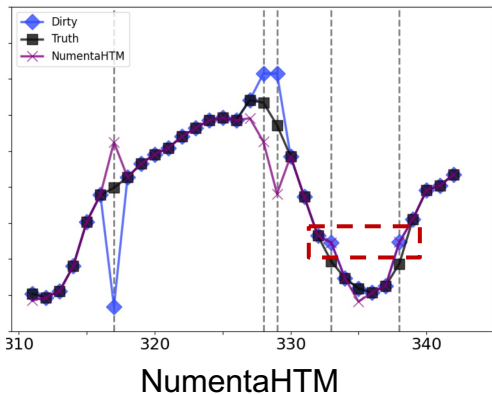
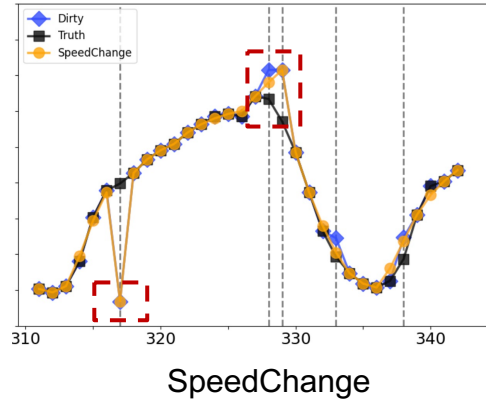
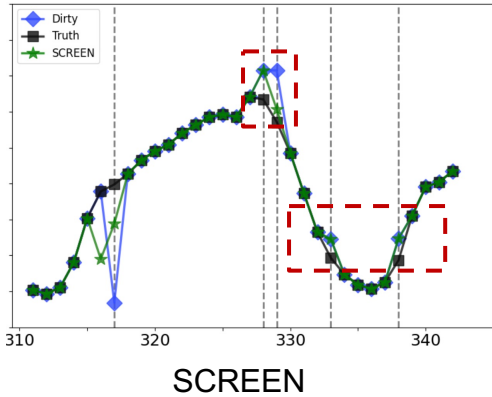
Time series data cleaning

- Classical algorithms
  - Moving Average (MA)
    - Change too many points
  - Auto Regression (AR)
    - Fit imprecisely

- Prevalent algorithms
  - SCREEN [SIGMOD'15]
  - SpeedChange [SIGMOD'16]
  - ~~IMR [VLDB'17] (supervised cleaning)~~
  - Numerous anomaly-detection based methods
  - ...

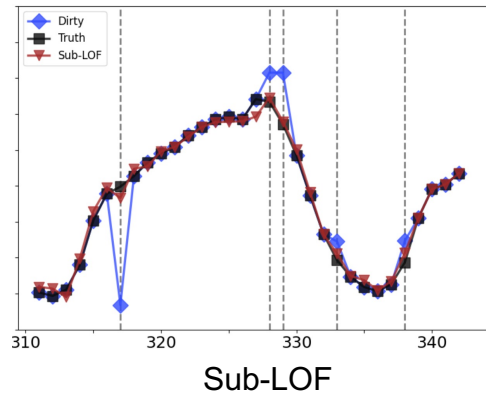
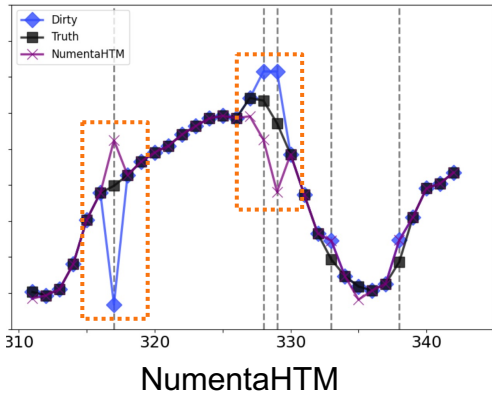
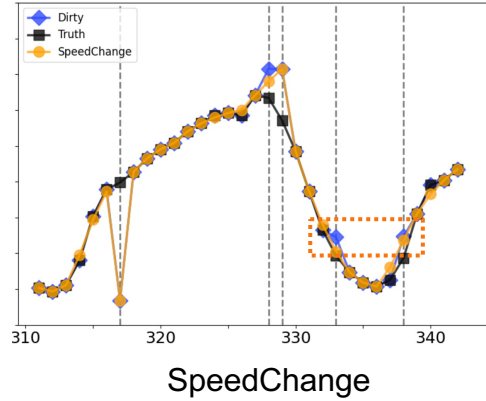
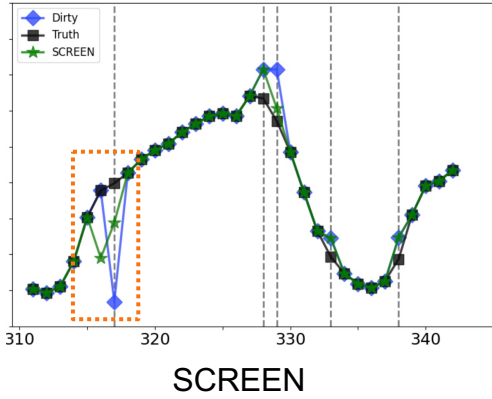


# Motivation



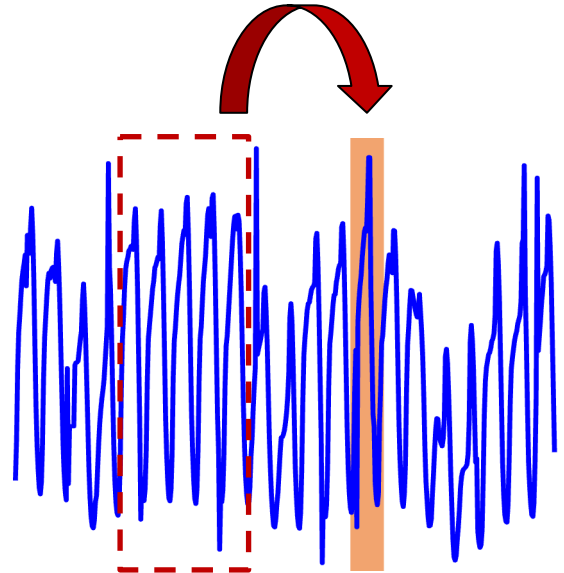
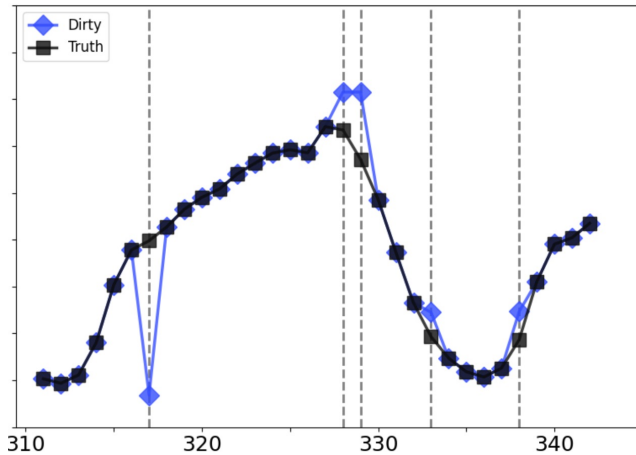
Mistakes and omissions in dirty data identification

# Motivation

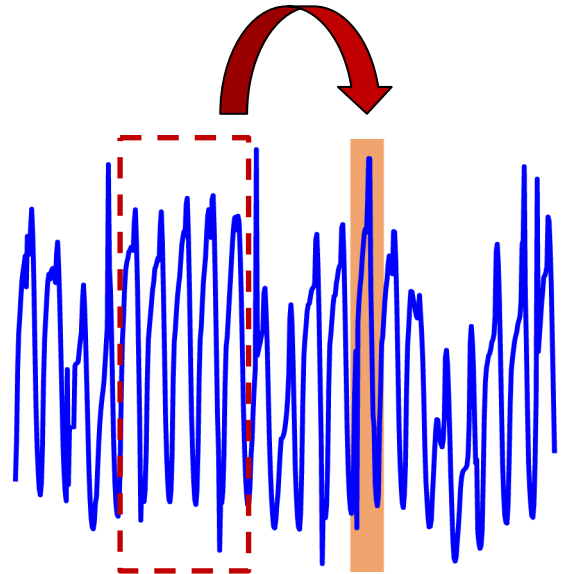
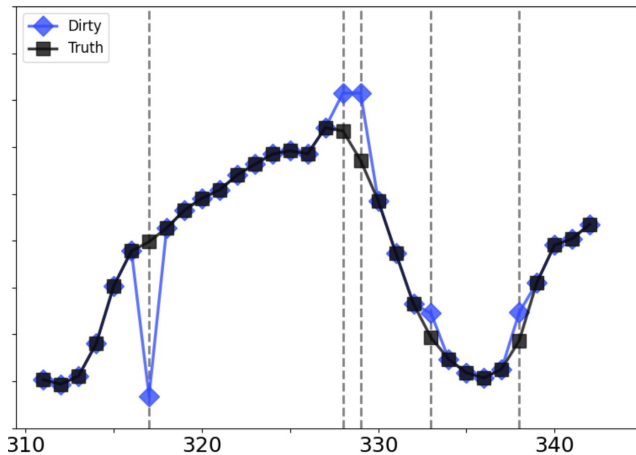


Deviates a lot from the truth in cleaning decisions

- Why fail?

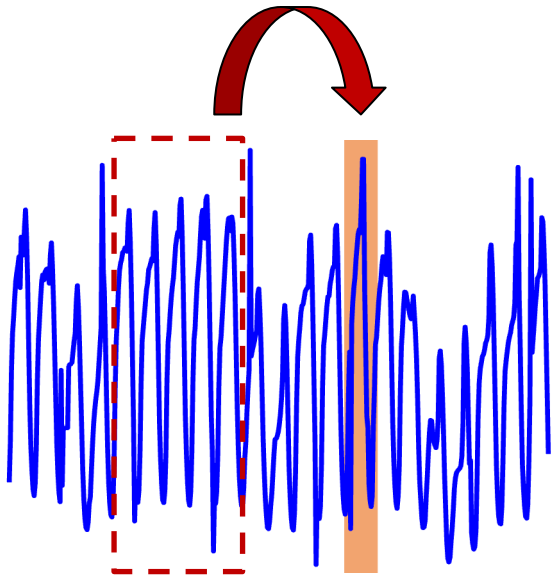


- Why fail?



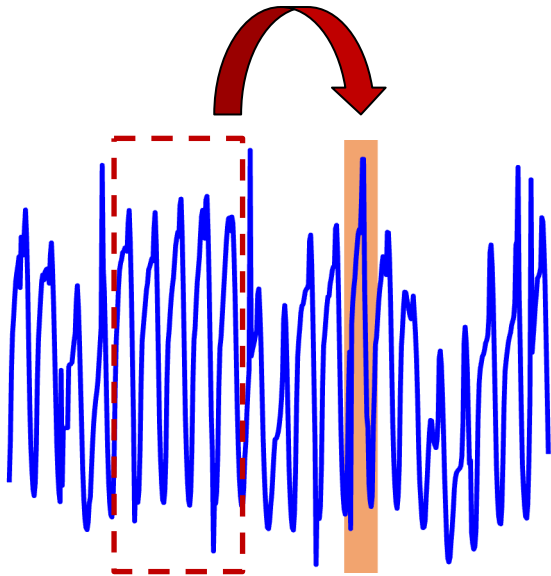
Plight: **Limited use** of inherent time series information

- What information can we use?



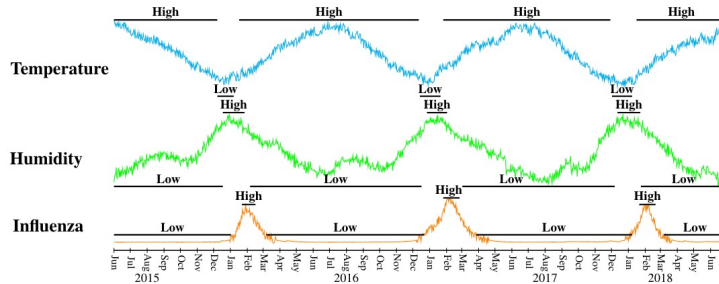
What are these?

- What information can we use?

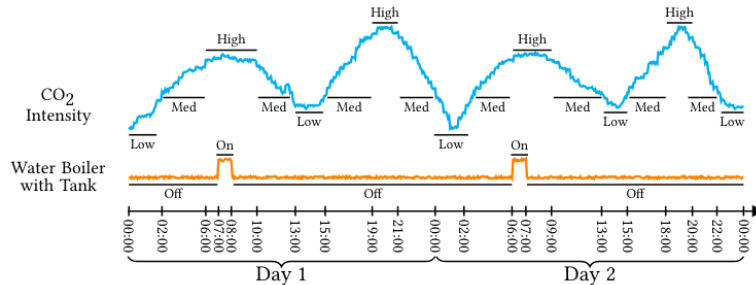


Recurrent patterns!

- Recurrent patterns in the time series

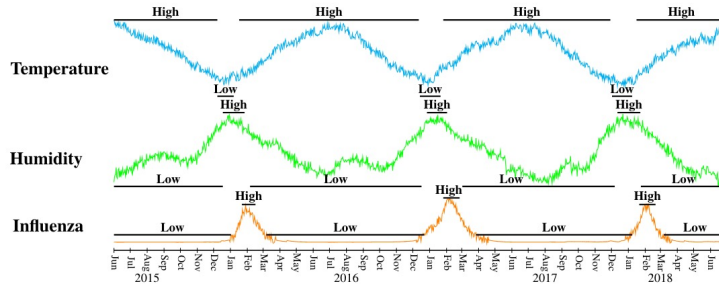


Weather and Influenza time series [ICDE'23]



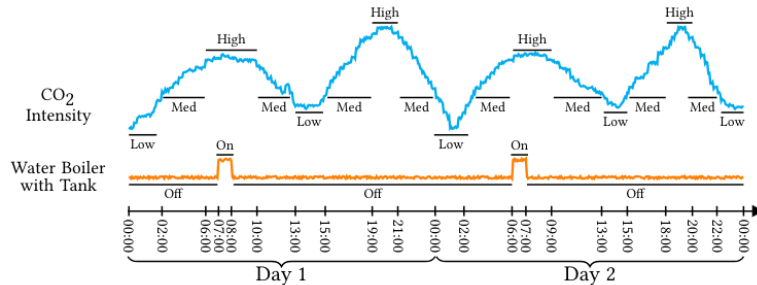
CO2 intensity and water boiler electricity usage [VLDB'21]

- Recurrent patterns in the time series



Recurrent  
≠  
Seasonal

Weather and Influenza time series [ICDE'23]



CO2 intensity and water boiler electricity usage [VLDB'21]



- How to use recurrent patterns?
- How to use similar patterns elsewhere in the time series to clean a point?

- Analogize recurrent patterns in a time series to fixed word combinations in a sentence

- Analogize **recurrent patterns** in a **time series** to **fixed word combinations** in a **sentence**
  - Repeated occurrence
  - Sequential dependency
  - A small part of whole data

- **Perplexity** in NLP

- An evaluation metric that measures the quality of language models, especially for text generation

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

- **Perplexity** in NLP

- An evaluation metric that measures the quality of language models, especially for text generation

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$



Introduce it to evaluate time series data quality

**Lower perplexity, better quality**

- **Perplexity-Guided** time series data cleaning
  - Given a dirty time series, cleaning the time series equals to **reducing** its perplexity **within the budget**

- **Perplexity-Guided** time series data cleaning
  - Given a dirty time series, cleaning the time series equals to **reducing** its perplexity **within the budget**

$$\text{Cost: } \Delta(X, X') = \|\{x_i \mid x_i \neq x'_i, 1 \leq i \leq n\}\|$$

- **Perplexity-Guided** time series data cleaning
  - Given a dirty time series, cleaning the time series equals to **reducing** its **perplexity** **within the budget**

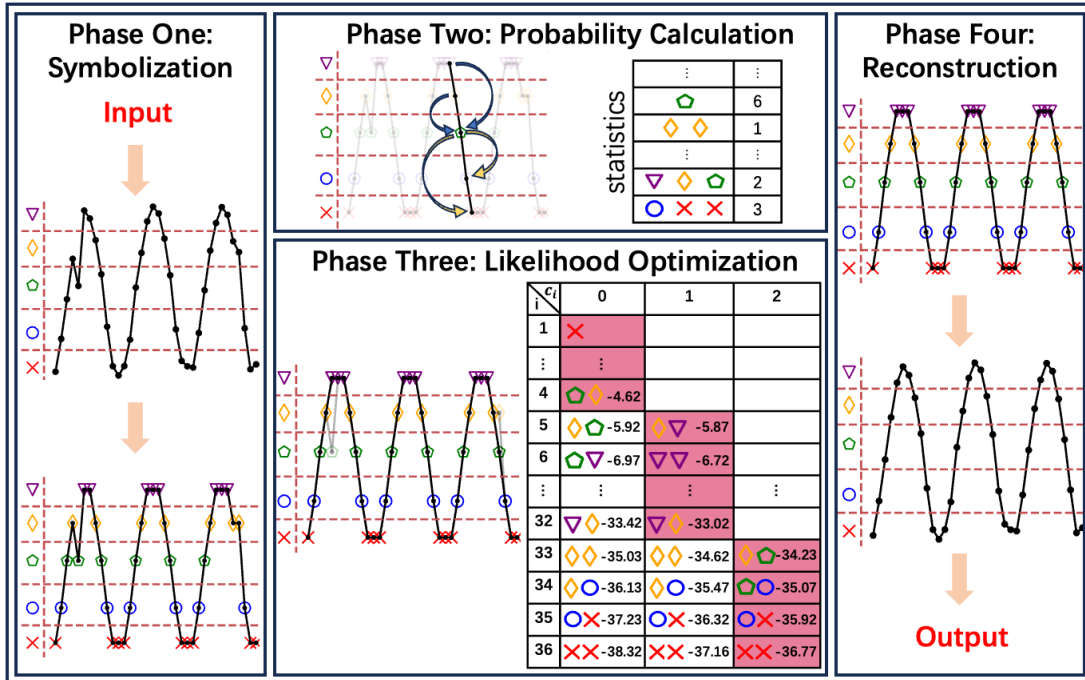
?

continuous value

$$\text{Cost: } \Delta(X, X') = \|\{x_i \mid x_i \neq x'_i, 1 \leq i \leq n\}\|$$



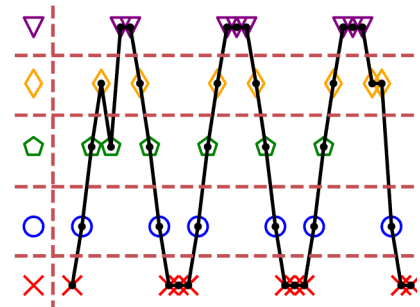
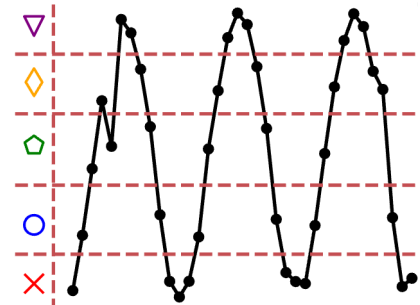
## • Overview



- Symbolization

- Change continuous value to discrete symbols to calculate the perplexity like textual data

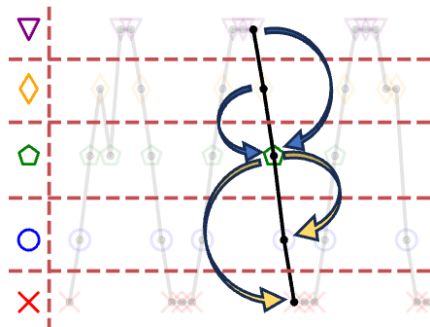
- K-means / uniform symbolization



- Probability calculation

- Use statistics (combinations) to calculate occurrence probability
- Introduce  $k$ -order Markov chain

$$P(x_i^w | x_{i-k:i-1}^w) = (C(x_{i-k:i}^w)) / (C(x_{i-k:i-1}^w))$$

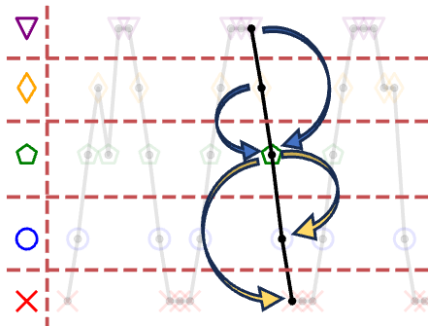


statistics	⋮	⋮
	⬠	6
	⬠ ⬠	1
	⋮	⋮
	⬡ ⬠ ⬠	2
	⬢ ⬠ ⬠	3

- Probability calculation

- Use statistics (combinations) to calculate occurrence probability
- Introduce  $k$ -order Markov chain

$$P(x_i^w | x_{i-k:i-1}^w) = (C(x_{i-k:i}^w) + 1) / (C(x_{i-k:i-1}^w) + r)$$



statistics	⋮	⋮
	⬠	6
	⬠ ⬠	1
	⋮	⋮
	⬡ ⬠ ⬠	2
	⬢ ⬠ ⬠	3

- Likelihood optimization

- Use likelihood rather than perplexity

$$PP(X) = PP(X^w) = \sqrt[n]{1/P(X^w)}$$

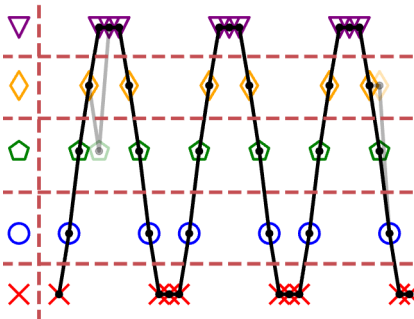


$$L(X) = L(X^w) = \log P(X^w)$$

- Likelihood optimization
  - Aim to maximize likelihood within budget

- Likelihood optimization

- A knapsack-like problem
- Use pseudo-polynomial time algorithm based on Dynamic Programming



$i \backslash c_i$	0	1	2
1	×		
⋮	⋮		
4	◇ ◇ -4.62		
5	◇ ◇ -5.92	▽ ▽ -5.87	
6	◇ ◇ -6.97	▽ ▽ -6.72	
⋮	⋮	⋮	⋮
32	▽ ▽ -33.42	▽ ▽ -33.02	
33	◇ ◇ -35.03	◇ ◇ -34.62	◇ ◇ -34.23
34	◇ ◇ -36.13	◇ ◇ -35.47	◇ ◇ -35.07
35	○ ○ -37.23	○ ○ -36.32	○ ○ -35.92
36	× × -38.32	× × -37.16	× × -36.77

- Likelihood optimization

- Recurrence equation

$$\Lambda(i, c_i, x'_{i-k+1:i}) = \max_{x'_{i-k} \in W_s} [\Lambda(i-1, c_{i-1}, x'_{i-k:i-1}) + \log P(x'_i | x'_{i-k:i-1})]$$

- $\Lambda(i, c_i, x'_{i-k+1:i})$  is the maximum likelihood of cleaned  $x'_{1:i}$  whose last  $k$  points are  $x'_{i-k+1}, \dots, x'_i$
- $c_i = \Delta(x'_{1:i}, x^W_{1:i})$  is the corresponding cleaning cost



- Likelihood optimization

- Recurrence equation

$$\Lambda(i, c_i, x'_{i-k+1:i}) = \max_{x'_{i-k} \in W_s} [\Lambda(i-1, c_{i-1}, x'_{i-k:i-1}) + \log P(x'_i | x'_{i-k:i-1})]$$

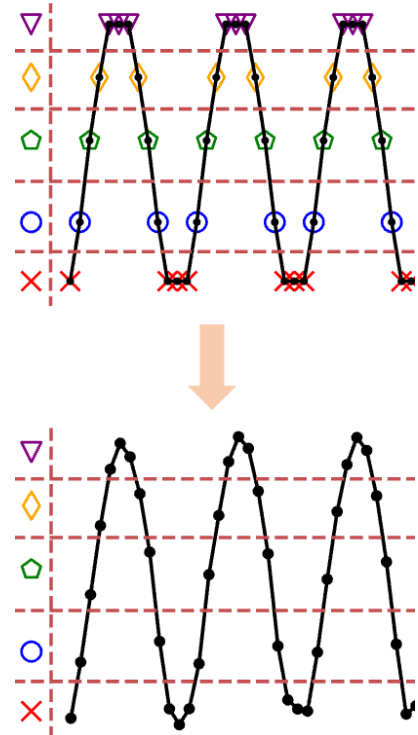
- Get cleaned result through retracement

- Time complexity  $O(r^{k+1}kn^2)$ ,  $r = |W_s|$

- Space complexity  $O(r^k n^2)$

- Reconstruction

- Change discrete symbols back to continuous values
- For each cleaned point, fit Linear Regression (LR) to obtain the most possible value



- Parameter advice

- K-means symbolization number  $r$

- argmin Davies-Bouldin Index

- Uniform symbolization number  $r$

- estimated range  $\left[ \frac{X_{max} - X_{min}}{2 * |\bar{v}|}, \frac{2 * (X_{max} - X_{min})}{|\bar{v}|} \right]$

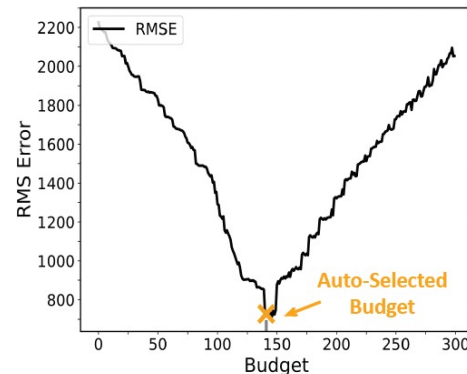
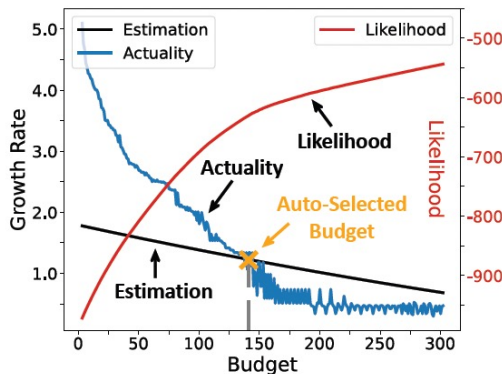
- Markov chain order  $k$

- directly choose 2 or 3

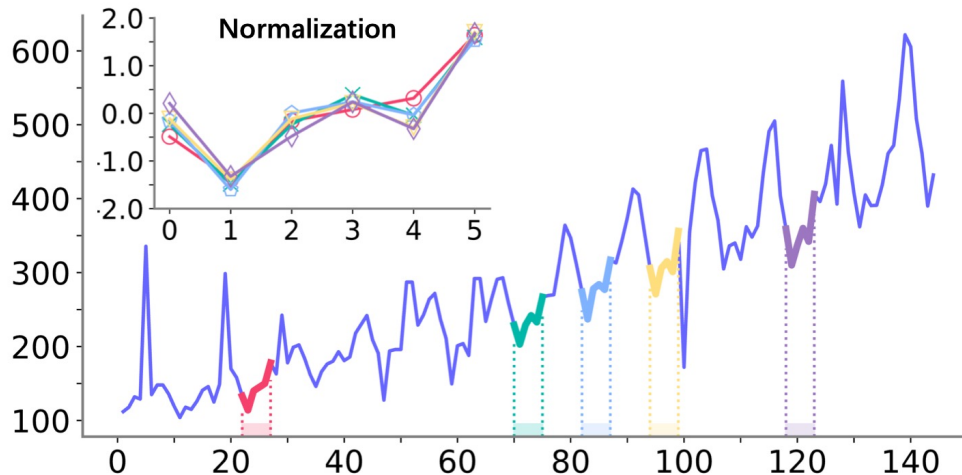
- Parameter advice

- Cleaning budget  $\delta$

- chase actual likelihood increase till  $\log \frac{m * (1-\tau)^{k+1} + 1}{m * (1-\tau)^k * \frac{\tau}{r-1} + 1}$
    - $\tau$  is current estimated dirty rate (we clean gradually)



- Homomorphic pattern
  - Different patterns but behaving similarly



- Homomorphic pattern aggregation

- Aggregate similar patterns to a homomorphic set

$$\text{Dist}(\hat{x}_{p-k:p}^w, \hat{x}_{q-k:q}^w) \leq \epsilon \quad (\text{distance after normalization})$$

- Homomorphic pattern aggregation

- Aggregate similar patterns to a homomorphic set

$$\text{Dist}(\hat{x}_{p-k:p}^w, \hat{x}_{q-k:q}^w) \leq \epsilon \quad (\text{distance after normalization})$$

An empirical value within 0.08 is good

- Homomorphic pattern aggregation

- Aggregate similar patterns to a homomorphic set

$$\text{Dist}(\hat{x}_{p-k:p}^w, \hat{x}_{q-k:q}^w) \leq \epsilon \quad (\text{distance after normalization})$$

- NP-Hard to find minimum sets from a time series
- We propose a heuristic method to solve it



- Homomorphic pattern aggregation
  - Use set frequency to calculate probability
  - For a pattern  $x_{p_i-k:p_i}'^w$  in a homomorphic set  $\{x_{p_1-k:p_1}'^w, \dots, x_{p_s-k:p_s}'^w\}$  of cardinality  $s$

$$P(x_{p_i}'^w | x_{p_i-k:p_i-1}'^w) = \frac{\sum_{j=1}^s C(x_{p_j-k:p_j}'^w) + 1}{\sum_{j=1}^s C(x_{p_j-k:p_j-1}'^w) + r}$$

- Homomorphic pattern aggregation
  - Use set frequency to calculate probability
  - For a pattern  $x_{p_i-k:p_i}'^w$  in a homomorphic set  $\{x_{p_1-k:p_1}'^w, \dots, x_{p_s-k:p_s}'^w\}$  of cardinality  $s$

$$P(x_{p_i}'^w | x_{p_i-k:p_i-1}'^w) = \frac{\sum_{j=1}^s C(x_{p_j-k:p_j}'^w) + 1}{\sum_{j=1}^s C(x_{p_j-k:p_j-1}'^w) + r}$$

- Homomorphic pattern aggregation
  - Use set frequency to calculate probability
  - For a pattern  $x_{p_i-k:p_i}'^w$  in a homomorphic set  $\{x_{p_1-k:p_1}'^w, \dots, x_{p_s-k:p_s}'^w\}$  of cardinality  $s$

$$P(x_{p_i}'^w | x_{p_i-k:p_i-1}'^w) = \frac{\sum_{j=1}^s C(x_{p_j-k:p_j}'^w) + 1}{\sum_{j=1}^s C(x_{p_j-k:p_j-1}'^w) + r}$$

Time complexity  $O(r^{k+1}kn^2)$   $\Rightarrow$   $O(r^{k+1}kn^2 \log n)$

- Greedy-based heuristic algorithm
  - Global optimum is too consuming
  - Each time we choose to clean one point to obtain maximum likelihood increase

$$\arg \max_{i, x_i''^w} \sum_{j=i}^{i+k} \left[ \log P(x_j''^w | x_{j-k:j-1}''^w) - \log P(x_j'^w | x_{j-k:j-1}'^w) \right]$$

- Greedy-based heuristic algorithm

- Global optimum is too consuming
- Each time we choose to clean one point to obtain maximum likelihood increase

$$\arg \max_{i, x_i''^w} \sum_{j=i}^{i+k} \left[ \log P(x_j''^w | x_{j-k:j-1}''^w) - \log P(x_j'^w | x_{j-k:j-1}'^w) \right]$$

No aggregation  $O(r^{k+1}kn^2)$   $\Rightarrow$   $O(rk^3n + kn \log n)$

With aggregation  $O(r^{k+1}kn^2 \log n)$   $\Rightarrow$   $O(rk^3n \log n)$

- Datasets

- 12 datasets with synthetic errors and real errors

Name	Dirty Type	Length	Rate	Source	Notes
CA	Synth	1K-10K	5-30%	California ISO	Energy
Romania	Synth	1K-10K	5-30%	Kaggle	Energy
Product	Synth	397	20%	FRED	Utilities
Retail	Synth	377	20%	FRED	Trade
Passenger	Synth	144	20%	Kaggle	Airline
Traffic	Real	200	2%	Dodgers [3, 26]	Flow
ID_7c18	Real	1500	0.73%	IOPS [36]	KPI
ID_7698	Real	1K	0.2%	IOPS [36]	KPI
ID_a40b	Real	10K	0.32%	IOPS [36]	KPI
MSL-T-5	Real	2218	0.23%	NASA-MSL [25]	Telemetry
SMAP-G-2	Real	7361	0.014%	NASA-SMAP [25]	Telemetry
Stock	Synth	12824	10%	SCREEN [41]	Finance

- Algorithms

- Akane, AkaneH, AkaneH+
- EWMA, AR-Linear, AR-LSTM, SCREEN, SpeedChange, Torsk, NumentaHTM, TripleES, FFT, GrammarViz, Sub-LOF

- Algorithms

 K-means

 Akane, AkaneH, AkaneH+

- EWMA, AR-Linear, AR-LSTM, SCREEN, SpeedChange, Torsk, NumentaHTM, TripleES, FFT, GrammarViz, Sub-LOF



- Algorithms

- Akane, AkaneH, AkaneH+



uniform + aggregation

- EWMA, AR-Linear, AR-LSTM, SCREEN, SpeedChange, Torsk, NumentaHTM, TripleES, FFT, GrammarViz, Sub-LOF

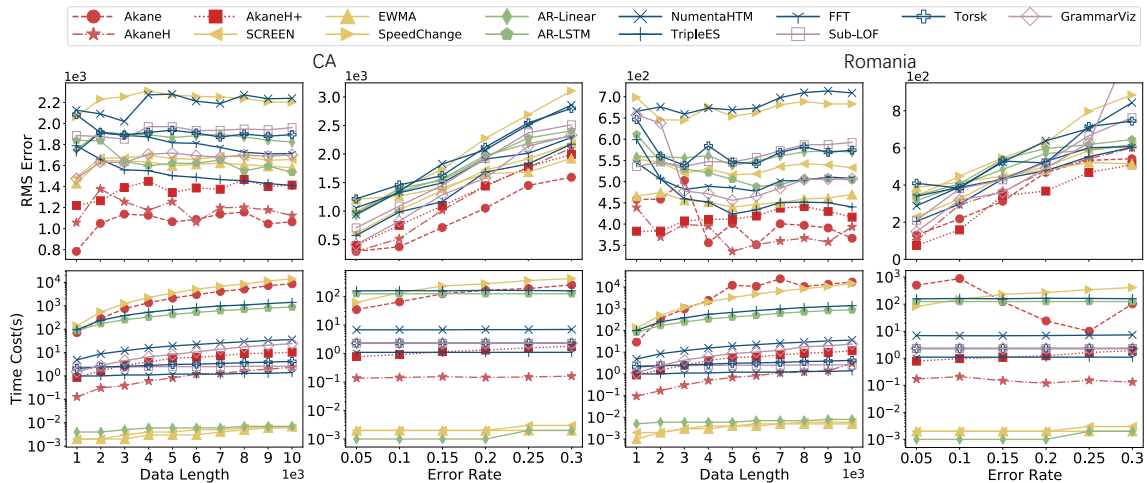
- Measurements

- Root Mean Square Error (RMSE)
- Accuracy

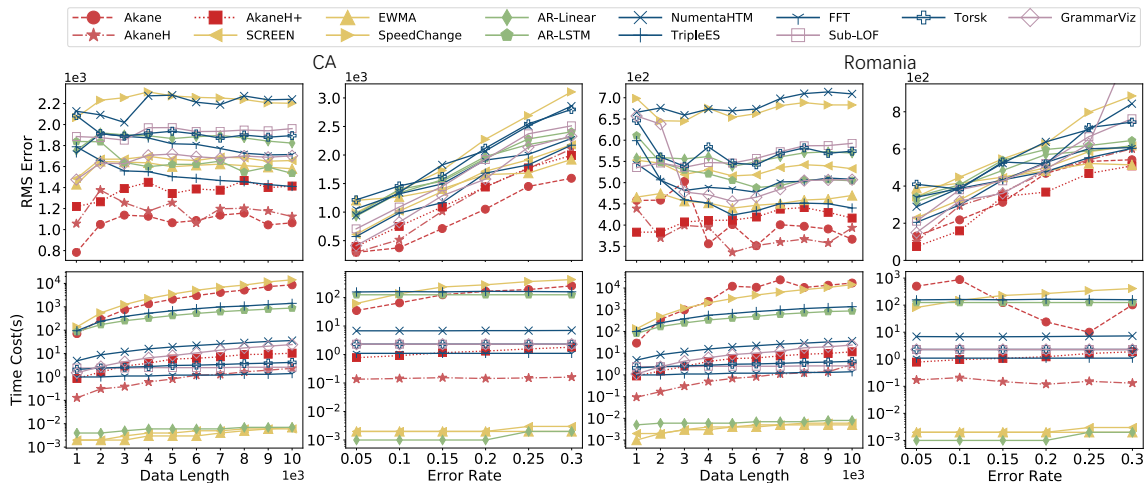
$$accuracy = \left( RMSE^{dirty} - RMSE^{cleaned} \right) / \left( RMSE^{dirty} - RMSE_{min}^{cleaned} \right)$$

- Execution time cost

- Vary length & error rate on CA and Romania

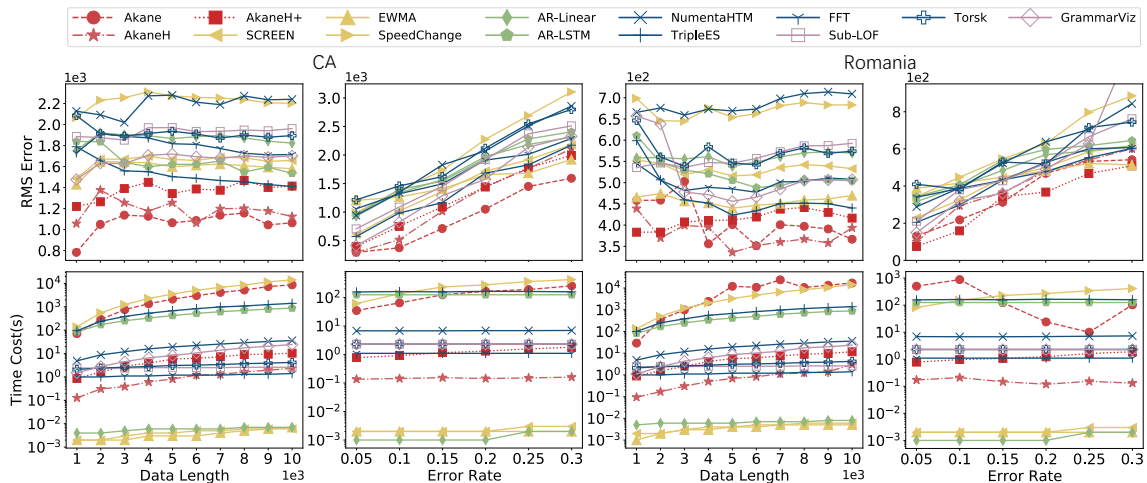


- Vary length & error rate on CA and Romania



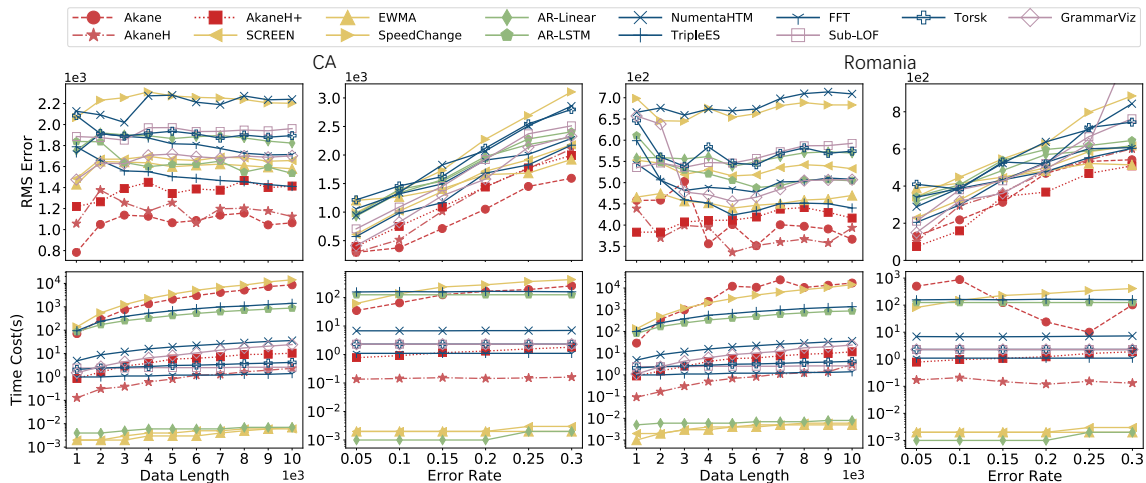
1. More powerful than existing algorithms

- Vary length & error rate on CA and Romania



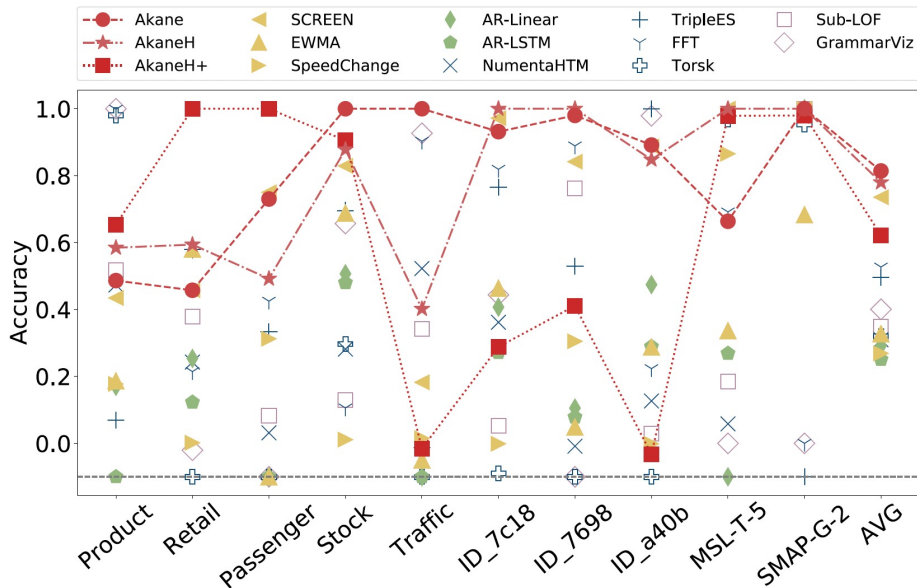
1. More powerful than existing algorithms
2. Relatively stable at all length; Better at low rate

- Vary length & error rate on CA and Romania



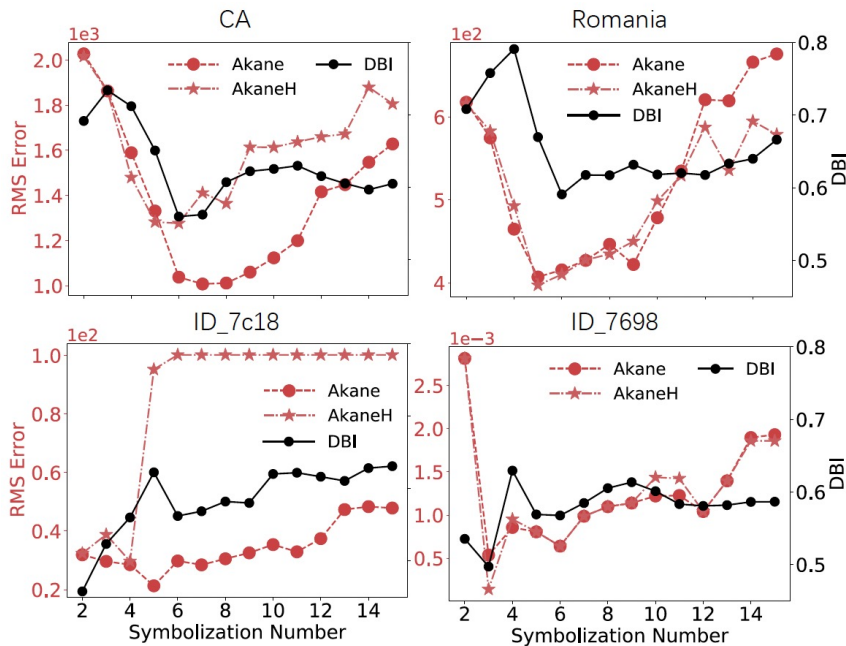
1. More powerful than existing algorithms
2. Relatively stable at all length; Better at low rate
3. Faster heuristic algorithm without many sacrifices

- Performance on various datasets



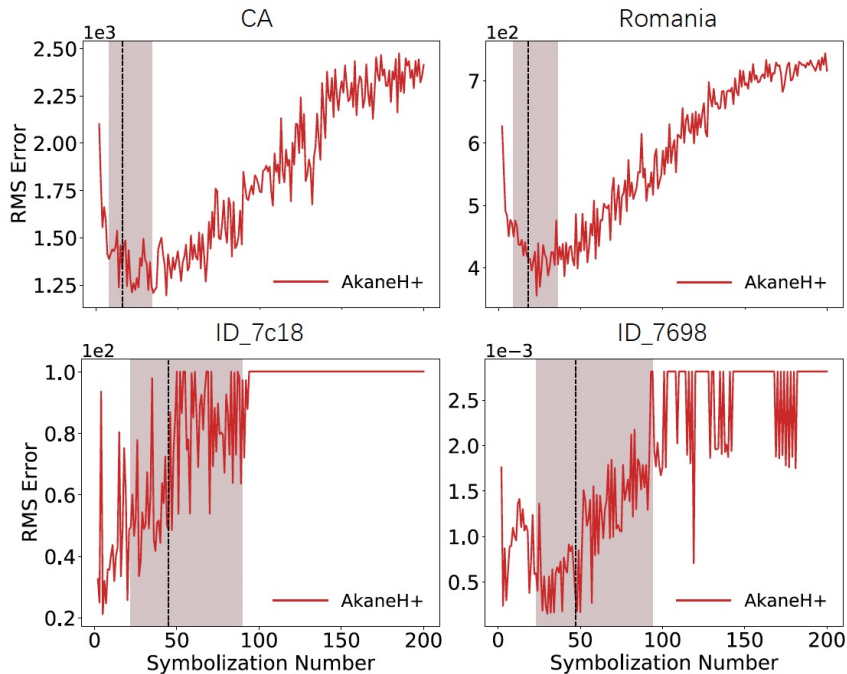
Overall great performance

- Vary K-means symbolization number  $r$

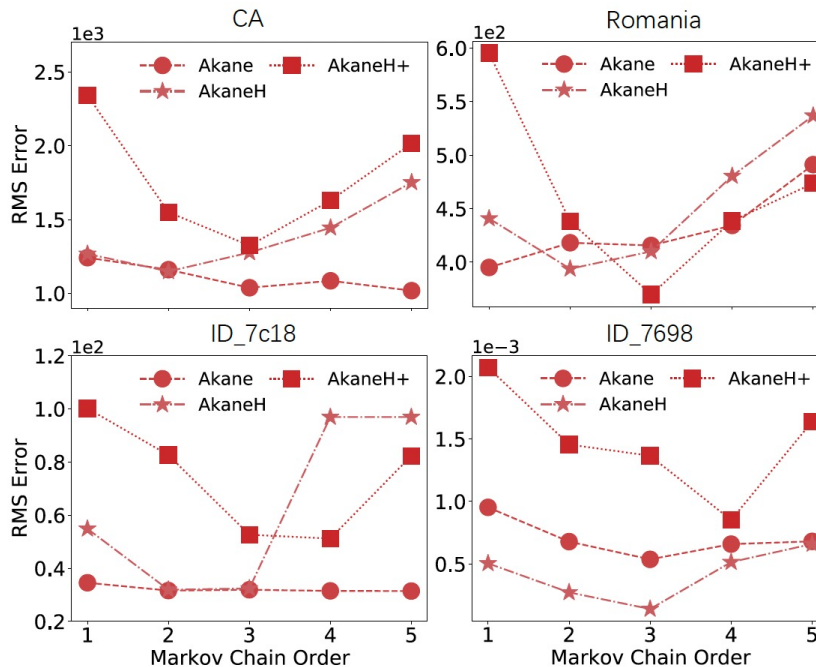




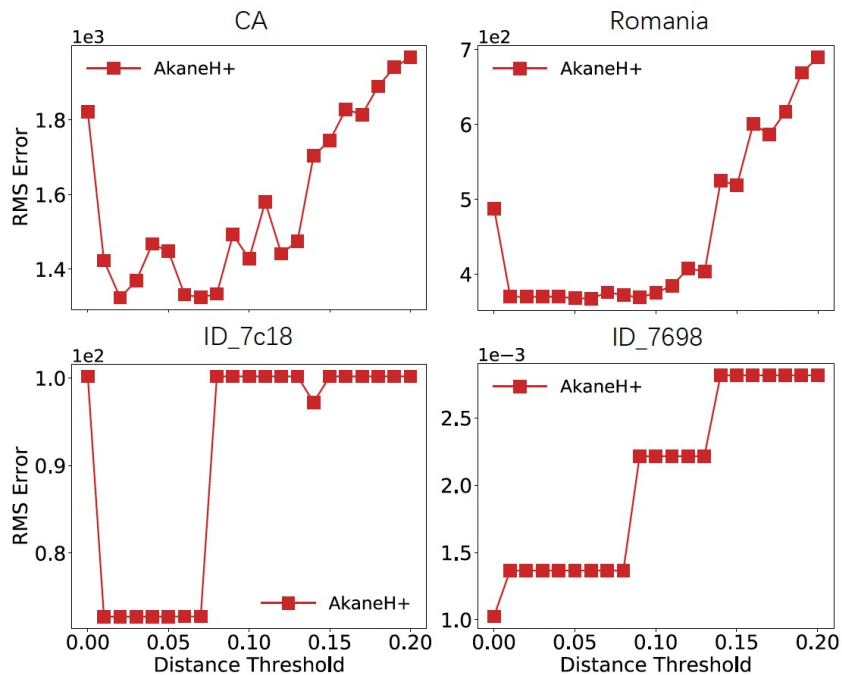
- Vary uniform symbolization number  $r$



- Vary Markov chain order  $k$



- Vary distance threshold  $\epsilon$



- Auto budget selection strategy evaluation

Evaluation	Akane	AkaneH	AkaneH+
Vary CA Length	0.953	0.944	0.966
Vary Romania Length	0.922	0.947	0.973
Vary CA Rate	0.974	0.894	0.939
Vary Romania Rate	0.892	0.853	0.991
10 fixed datasets	0.829	0.833	0.811

- Parameter experiments show
  - Performances have close relations to parameters
  - Our parameter advice is effective

- We use the insight of recurrent patterns to **formalize perplexity-guided time series data cleaning problem**
- We propose a **four-phase framework** to solve the cleaning problem, with parameter analysis and advice
- We further introduce **advanced solutions**, involving homomorphic pattern aggregation and greedy-based heuristic algorithm, to enhance generalization

**THANKS!**