

ANALYZING RAW LOG FILES TO FIND EXECUTION ANOMALIES

Viktor Jovanoski, Jan Rupnik,
Mario Karlovčec, Blaž Fortuna

SiKDD - October 2017

OVERVIEW

- Anomaly detection
- Log file analysis
- Algorithm
- Evaluation
- Future work

ANOMALY DETECTION

- **Anomaly: any data point that significantly deviates from the remaining data**
- **Goals:**
 - Detect anomalous sudden changes (short term)
 - Detect slow degradation (long term)
 - Root-cause analysis
- Applicable to many fields, we monitor specific IT system

ANOMALY DETECTION - CHALLENGES

- Diverse data types – numeric, discrete, text, sequences, ...
- Running on stream/batch
- Actionability – users require understandable explanations
- Proper quantity – not too many, not too few
- “Concept drift” – reality changes over time

GENERAL STREAMING ALGORITHM

```
while input data available do  
    parse data and extract record  
    calculate anomaly score  
    if score above threshold then  
        report anomaly  
    end if  
    add record to model  
end while
```

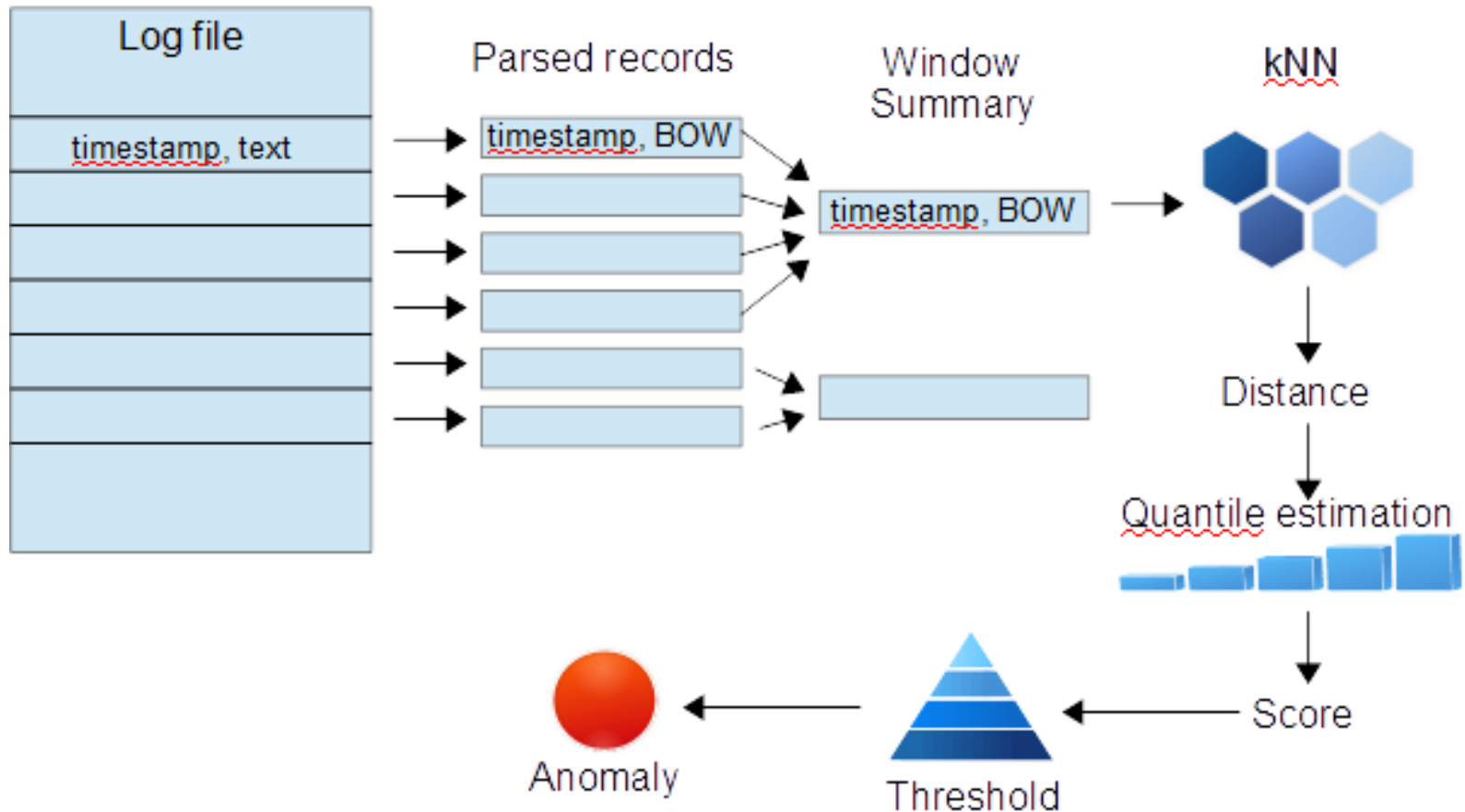
LOG FILE ANALYSIS

- Why log files? Sometimes this is all we have.
- Unstructured text, but contain timestamp
- Parse location identifiers, detect errors and warnings
- Identify “execution profiles” – descriptions of time intervals

ALGORITHM

```
while input data available do  
    read data from log file  
    parse data and extract BOW  
    insert into time-windows  
    calculate distance in kNN  
    if score quantile above threshold then  
        create explanation using distance  
        report anomaly  
    end if  
    add record to model  
end while
```

PIPELINE



ALGORITHM PARAMETERS

- k - # of nearest neighbors for distance comp.
- `time window` – the length of summary windows
- `history window` – how many historical examples do we compare against
- `nn_rate` – quantile treshold

- Must be hand-tuned to find balance between quantity, statistical quality and actionability

ANOMALY EXPLANATION

- Must be ACTIONABLE
- Show difference to nearest neighbor
 - Most often, unusual processes stick out
 - Must be interpreted by an expert

EVALUATION

- Not many labeled datasets
- Mostly unlabeled data
 - Previously undetected anomalies
- Precision and recall
- Test1: Synthetic data
- Test2: Real web-server logs

EVALUATION – SYNTHETIC DATA

- Simulating parallel execution of several processes, with deterministic pattern
 - 8 extra anomalous instances, 2 per week
- 1 year of data – 1 min and 1 h windows, history = 10 days, $k = 1$, $nn_rate = 0.003$
- 1 min window – recall=1.0, precision=0.06
- 1 h window – recall=0.66, precision=0.14

EVALUATION — WEB-SERVER LOGS

- Static web site – not many content changes
- 1 month of data – 5 min windows, history = 10 days, $k = 1$
- Manual inspection of generated anomalies
- Those with **one strong dimensional outlier** were “correct” – new content has been published

FUTURE WORK

- Full-spectrum analysis
 - Many different data sources
 - Correlation across data sources
 - Root-cause analysis
 - Prediction
 - Active learning
- Deep learning