

Handling Duplicated Tasks in Process Discovery by Refining Event Labels

<u>Xixi Lu</u> (x.lu@tue.nl)¹

Dirk Fahland¹

Frank van den Biggelaar²

Wil van der Aalst¹

¹Eindhoven University of Technology ²Maastricht University Medical Center

Overview

- (5 min) Research Problem
- (10 min) Approach
- (5 min) Evaluation and Conclusion



Ę



Research Problem

Input: Imprecise log



$$V \rightarrow A \rightarrow B \rightarrow C \rightarrow A \rightarrow S$$
$$V \rightarrow C \rightarrow B \rightarrow A \rightarrow S$$
$$V \rightarrow C \rightarrow A \rightarrow B \rightarrow S$$

Research Problem



So that

- 1. all traces & events preserved (no filtering)
- 2. model is more precise (than without refining)
- we can explore different refinement of labels interactively (because we don't know correct label refinement from given log)

Refining Event Labels imprecise **Precise** ►<mark>A1</mark>*B1*C1*A3* →<mark>A→B→C→A</mark>→ +C2+B2+A2+ ►<mark>A1</mark> B1 C1 A3 → C → B → A → +C2+A2+B2+ Refining Label +C2+<mark>B2</mark>+A2+ → C→ A→ B→ <mark>≁C3≁A4</mark>≁B3≁ Discover Discover Discover A1+++B1+++C1+++A3 A1+++B1+++C1+++A3

Ę

Refining Event Labels ●→ System unknown! → A → ● **Precise** imprecise *<mark>A1*B1*C1*A3</mark>* →<mark>A→B→C→A</mark>→ +C2+B2+A2+ *<mark>A1*B1*C1*A3</mark>* → C → B → A → *C2*A2*B2* **Refining Label** +C2+B2+A2+ → C→ A→ B→ <mark>→C3→A4</mark>→B3→ Discover Discover Discover A₁+•+B₁+•+C₁+•+A₃ A1+++B1+++C1+++A3 We don't know what is optimal!

Ę



Approach



Approach



Mapping Between Events



... based on "structural context of events"



... based on "structural context of events"(1) Differences in neighbors



- ... based on "structural context of events"
- (1) Differences in neighbors
- (2) Differences in structure



Distance = 3

- ... based on "structural context of events"
- (1) Differences in neighbors
- (2) Differences in structure



- ... based on "structural context of events"
- (1) Differences in neighbors
- (2) Differences in structure



Cost so far = 6+10+8+5+3

- ... based on "structural context of events"
- (1) Differences in neighbors
- (2) Differences in structure
- (3) #Dissimilar events



Total Cost = 6+10+8+5+3 +1

There is an algorithm to compute an optimal mapping :

Xixi Lu, Dirk Fahland, Frank J.H.M. van den Biggelaar, and Wil M.P. van der Aalst. **Detecting Deviating Behavior without Models**. In *BPM Workshops* 2015, volume 256, pp.126–139, Springer International Publishing, 2015.

Approach







a) Normalize costs w.r.t. maximal cost seen in the log



a) Normalize costs w.r.t. maximal cost seen in the log

b) Set variant threshold, say, 0.8

c) Remove edges if cost > variant threshold







Approach



Implemented in ProM

• There will be a demo in the demo track.



Evaluation

• 3 experiments, each 600 models, k transitions w/ same label

E1) Default parameters, k = 4
E2) Adaptive parameters, k = 4
E3) 1 in loop, k = 2



Experiment Result - Example



Experiment Result -F1-score of 🔀 and 🔀 w.r.t.





Improved 89% Models

E3) 1 in loop



Improved 60% Models

Experiment Results and Limitations



Real-life Example : Hospital Data



Conclusion and Future Work

- Improved logs (and process discovery) by refining labels for up to 87%
- Interactive and explorative
- Future work
 - Different ways to compute similar events using context of events
 - Log preprocessing framework
- Integrated in "Log to Model Explorer"
 - Supporting clustering, filtering and label refining
- Come to the Demo!

Questions?