



**Presenter:** Mozhgan Vazifehdoostirani

Jan Vaznendoostnam

m.vazifehdoostirani@tue.nl

Authors:

Mohammadreza Fani Sani, Mozhgan Vazifehdoostirani, Gyunam Park, Marco Pegoraro, Sebastiaan J. van Zelst, and Wil M.P. van der Aalst What would happen in the future?





It is an extremely active field of research!







Can we improve the computational performance by using sampled event logs, while maintaining the accuracy?









Select the best Candidates

Reduce the training time



**Quality / Accuracy** 

# Hypothesis

Maintain accuracy We do not loose too much information

RNTHAAC









8











10



How many cases should be chosen?







**Evaluation** 

1) Datasets



RWTHAACHEN

TECHNOLOGY

**Evaluation** 

2) Prediction objective





OLOGY

### Evaluation

4) Metrics

 $R_s = \frac{\text{Size of the whole event log}}{\text{Size of the sampled event log}}$ 

 $R_t = \frac{\text{Training time using whole data}}{\text{Training time using the sampled data}}$ 

 $R_{Acc} = \frac{\text{Accuracy using the sampled training log}}{\text{Accuracy using the whole training log}}$ 

 $R_{FE} = \frac{\text{Feature extraction time using whole data}}{\text{Feature extraction time using the sampled data}}$ 





### The reduction in the size of training logs and the improvement in the performance of feature extraction

Sample Methods			Divisi	on			Logarithmic distribution							unique	
	K = 2		K = 3		K = 10		log2		logȝ		log10				
Event Log	$R_s$	R <sub>FE</sub>	$R_s$	R <sub>FE</sub>	$R_s$	R <sub>FE</sub>	$R_s$	$R_{FE}$	$R_s$	$R_{FE}$	$R_s$	R <sub>FE</sub>	$R_s$	R <sub>FE</sub>	
RTFM	1.99	4.8	3	11.1	9.8	106.9	153.5	12527.6	236.3	23699.2	572.3	74912.8	285.1	24841.8	
BPIC-2012-W	1.22	1.37	1.41	1.8	1.66	2.51	6.06	22.41	9.05	37.67	28.5	208.32	1.73	2.36	





### The reduction in the size of training logs and the improvement in the performance of feature extraction

Sample Methods			Divisi	on			Logarithmic distribution							unique	
	K = 2		K = 3		K = 10		log2		log3		log10				
Event Log	$R_s$	$R_{FE}$	$R_s$	R <sub>FE</sub>	$R_s$	$R_{FE}$	$R_s$	R <sub>FE</sub>	$R_s$	$R_{FE}$	R <sub>s</sub>	$R_{FE}$	$R_s$	R <sub>FE</sub>	
RTFM	1.99	4.8	3	11.1	9.8	106.9	153.5	12527.6	236.3	23699.2	572.3	74912.8	285.1	24841.8	
BPIC-2012-W	1.22	1.37	1.41	1.8	1.66	2.51	6.06	22.41	9.05	37.67	28.5	208.32	1.73	2.36	

Lowest

Highest

RNTH



DLOGY

### Results

Highest

RWTHAACHEN TU

19

eindhoven UNIVERSITY OF TECHNOLOGY

### • LSTM

Lowest

Sample			Divisior	ı				Lo		unique				
Methods	K = 2		K = 3		K = 10		log2		log3		log10			
Event Log	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	$R_t$	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>
RTFM	1.001	2	1.004	2.9	0.99	9	0.716	26.7	0.724	33	0.767	41.8	0.631	29.1
BPIC-2012-W	1	1.4	0.985	1.3	0.938	1.3	0.977	4.7	0.97	5.8	0.876	11.9	0.996	1.6

### • XG Boost

Sample			Divisio	n				L		unique				
Methods	K =	K = 2		K = 3		K = 10		log2		log3			omqoe	
Event Log	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	$R_t$	R <sub>Acc</sub>	R <sub>t</sub>								
RTFM	1	2.4	1	1.4	1	84.1	0.686	126.4	0.706	191.8	0.772	355	0.582	297.7
BPIC-2012-W	0.999	2.3	0.998	2.4	0.997	3.4	0.923	10.7	0.97	16.7	0.883	64.8	0.997	2.8



### Results

20

EINDHOVEN UNIVERSITY OF TECHNOLOGY

9

Accuracy improvement

### • LSTM

Sample			Divisior	ı				unique						
Methods	K =	2	K = 3		K = 1	.0	log	2	lo	g3	log:	10	oning	
Event Log	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	$R_t$	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>	R <sub>Acc</sub>	R <sub>t</sub>
RTFM	1.001	2	1.004	2.9	0.99	9	0.716	26.7	0.724	33	0.767	41.8	0.631	29.1
BPIC-2012-W	1	1.4	0.985	1.3	0.938	1.3	0.977	4.7	0.97	5.8	0.876	11.9	0.996	1.6

### • XG Boost

Sample			Divisio	n				L		unique				
Methods	K = 2		K = 3		K = 10		log2		log3		log10		onique	
Event Log	R <sub>Acc</sub>	$R_t$	R <sub>Acc</sub>	$R_t$	R <sub>Acc</sub>	R <sub>t</sub>								
RTFM	1	2.4	1	1.4	1	84.1	0.686	126.4	0.706	191.8	0.772	355	0.582	297.7
BPIC-2012-W	0.999	2.3	0.998	2.4	0.997	3.4	0.923	10.7	0.97	16.7	0.883	64.8	0.997	2.8



Discussion

• Sampling event logs could increase the performance and keep the accuracy in some cases



21



- Sampling event logs could increase the performance and keep the accuracy in some cases
- We observed that different event logs needs different sampling methods





- Sampling event logs could increase the performance and keep the accuracy in some cases
- We observed that different event logs needs different sampling methods
- Characteristics of the given event log and suitable sampling parameters has more effect than number of sampled cases or prediction models





Discussion

- Sampling event logs could increase the performance and keep the accuracy in some cases
- We observed that different event logs needs different sampling methods
- Characteristics of the given event log and suitable sampling parameters has more effect than number of sampled cases or prediction models
- Using the proposed sampling method, we could speed up hyperparameters tunning and adapting with changes due to concept drift



### Future work













26

