On the Role of Fitness, Precision, Generalization and Simplicity in Process Discovery

Joos Buijs Boudewijn van Dongen *Wil van der Aalst*



http://www.win.tue.nl/coselog/

Technische Universiteit **Eindhoven** University of Technology

Where innovation starts

TU

Advances in Process Mining

- Many process discovery and conformance checking algorithms and tools are available (cf. the various **ProM** packages).
- Also commercial software based on these ideas: Disco (Fluxicon), Reflect (Futura/Perceptive), BPMOne (Pallas Athena/Perceptive), ARIS Process Performance Manager (Software AG), Interstage Automated Process Discovery (Fujitsu), QPR ProcessAnalyzer/Analysis (QPR Software), flow (fourspark), Discovery Analyst (StereoLOGIC), etc.
- We applied process mining in over 100 organizations.



A monifestra is a "public dedication of principles and intentions" by a group of people. This monifesto is written by mandhers and supporters of the IEEE Tosk Force on Process Mining. The good of this task force is to promose the research, development, education, implementhate, reducation, and understanding of process mining. ng is a sharaky zame ang is a sharaky zame ang is a sharaky at yang is a More than 75 people involving more than 50 organizations created the Process Mining Manifesto in the context of the IEEE Task Force on Process Mining.

Available in 13 languages



Example Process Discovery (Vestia, Dutch housing agency, 208 cases, 5987 events)



PAGE 2



Challenge: Four Competing Quality Criteria



Example: one log four models



 N_4 : fitness = +, precision = +, generalization = -, simplicity = -

Model N₁



#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

Model N₂



#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	



Model N₄



#	trace
455	acdeh
191	abdeg
177	adceh
144	abdeh
111	acdeg
82	adceg
56	adbeh
47	acdefdbeh
38	adbeg
33	acdefbdeh
14	acdefbdeg
11	acdefdbeg
9	adcefcdeh
8	adcefdbeh
5	adcefbdeg
3	acdefbdefdbeg
2	adcefdbeg
2	adcefbdefbdeg
1	adcefdbefbdeh
1	adbefbdefdbeg
1	adcefdbefcdefdbeg
1391	

Another challenge: Huge search space



... with just a few interesting candidates

M M M Μ Μ Μ Μ M Μ Μ Μ M _ _M Μ M M _M ^м м М M M M M M M Μ Μ ΜM M _ _M MM M Μ Μ Μ Μ Μ Μ Μ MM Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ мМ Μ Μ м М_ Μ ΜΜ Μ MM M M Μ мМ ΜΜ ΜΜ M M_M M_M M M Μ ΜM Μ Μ M_N Μ ΜM Μ Μ Μ Μ Μ MM M M Μ Μ M M ΜM Μ Μ Μ Μ M M Μ Μ Μ Μ Μ Μ <u>_M</u>M Μ., M M Μ Μ ММ МM Μ Μ Μ Μ Μ Μ ^м м Μ Μ M Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ M Μ Μ ΜΜ Μ м м м Μ Μ Μ Μ мМ Μ Μ Μ ΜM M_M M_M ΜΜ M_ Μ Μ IVI IVI M M M Μ Μ Μ ΜΜ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ MMM Μ MM._MM Μ ммм Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ M M M Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ Μ мМ Μ Μ Μ Μ м ^{м м м} Μ Μ Μ Μ MM Μ Μ ΜΜ Μ Μ ΜΜ M M M Μ Μ Μ M M M Μ M Μ Μ

Two requirements

- 1. It should be possible to seamlessly balance the different quality criteria based on user-defined preferences.
- 2. The algorithm should always return a "correct" process model and not waste time on model having deadlocks and other anomalies.





Proposal: Evolutionary Tree Miner (ETM)

- Process trees as representation (= limit search space to "good" models).
- Genetic approach (= very flexible)
- Fitness function uses all four criteria (= seamlessly balance the different "forces")

Representational Bias: Process Trees



Petri Net Semantics

(used for comparison and conformance checking only)



Steps of the Genetic ETM Algorithm



Population Change



Four Metrics (see paper)



 $Q_{rf} = 1 - \frac{\text{cost for aligning model and event log}}{\text{Minimal cost to align arbitrary event log on model and vice versa}}$ $Q_s = 1 - \frac{\#\text{duplicate activities} + \#\text{missing activities}}{\#\text{nodes in process tree} + \#\text{event classes in event log}}$ $Q_p = 1 - \frac{\sum_{\text{visited markings}} \#\text{visits} * \frac{\#\text{outgoing edges} - \#\text{used edges}}{\#\text{outgoing edges}}}{\#\text{total marking visits over all markings}} \quad 1 = \text{optimal}$ $Q_g = 1 - \frac{\sum_{\text{nodes}} (\sqrt{\#\text{executions}})^{-1}}{\#\text{nodes in tree}} \quad 0 = \text{very bad}$





Conventional Algorithms (1/3) ("best effort" mapping to process trees to allow for comparison)

alpha miner





f: 0,992	p: 0,995
s: $1,000$	g: 0,889

low fitness

ILP miner



language-based region miner



f: 1,0	00 p:	0,784
s: 0,9	33 g:	0,830

low precision



low fitness





sound

Conventional Algorithms (2/3)



f: 1,000	p: 0,986
s: 0,875	g: 0,852

multi-phase miner



f: 1,000	p: 0,830
s: 1,000	g: 0,889

Conventional Algorithms (3/3)

genetic miner







state-based region miner





f: 1	,000	p:	0,893
s: 0	,933	g:	0,830

Often unsound result and no mechanism to seamlessly balance the four criteria



Genetic Mining (ETM) While Considering Only One Criterion



Considering Replay Fitness and One Other Criterion

Considering 3 of 4 Criteria

replay fitness needs to have a larger weight

Considering All Four Criteria with Emphasis on Fitness

fitness has weight 10

Initial Model Versus Discovered Model

Trace	#
ABCDEG	6
ABCDFG	38
A B D C E G	12
A B D C F G	26
A B C F G	8
A C B E G	1
A D B C F G	1
A D B C E G	1
A D C B F G	4
ACDBFG	2
ACBFG	1

G

 $\mathbf{F}\mathbf{F}$

f:	1,000	p: 0,893	
s:	0,933	g: 0,830	

simulated

А

Real-Life Event Logs

- Event log L0 is the event log used before. L0 contains 100 traces, 590 events and 7 activities.
- Event Log L1 contains 105 traces, 743 events in total, with 6 different activities.
- Event Log L2 contains 444 traces, 3.269 events in total, with 6 different activities.
- Event Log L3 contains 274 traces, 1:582 events in total, with 6 different activities.

Event logs L1, L2 and L3 are extracted from the information systems of municipalities participating in the CoSeLoG project (http://www.win.tue.nl/coselog/).

PAGE 29

Results

	LO	L1	L2	L3
α -algorithm	f: 0,992 p: 0,995	f: 1,000 p: 0.510	f: 1.000 p: 0.468	f: 0.976 p: 0.532
	s: 1,000 g: 0,889	s: 0.923 g: 0.842	s: 0.923 g: 0.885	s: 0.923 g: 0.866
	overall: 0,969	overall: 0,819	overall: 0,819	overall: 0,824
	f: 1,000 p: 0,748	f: 1.000 p: 0.551	f: 1.000 p: 0.752	f: 1.000 p: 0.479
ILP Miner	s: 0,933 g: 0,830	s: 0.857 g: 0.775	s: 0.923 g: 0.885	s: 0.857 g: 0.813
	overall: 0,887	overall: 0,796	overall: 0,890	overall: 0,787
	f: 1,000 p: 0,986	f: 0.966 p: 0.859	f: 0.917 p: 0.974	f: 0.995 p: 1.000
Heuristics	s: 0,875 g: 0,852	s: 0.750 g: 0.746	s: 0.706 g: 0.716	s: 1.000 g: 0.939
	overall: 0,928	overall 0,830	overall: 0,828	overall: 0,983
Genetic	f: 1,000 p: 0,922	f: 0,997 p: 0,808	f: 0.905 p: 0.808	f: 0.987 p: 0.875
	s: 0,737 g: 0,790	s: 0,750 g: 0.707	s: 0,706 g: 0.717	s: 0.750 g: 0.591
	overall: 0,862	overall: 0,815	overall: 0,784	overall: 0,801
ETM	f: 0,992 p: 0,995	f: 0,901 p: 0,989	f: 0,863 p: 0,982	f: 0,995 p: 1,000
	s: 1,000 g: 0,889	s: 0,923 g: 0,894	s: 0,923 g: 0,947	s: 1,000 g: 0,939
	overall: 0,969	overall: 0,927	overall: 0,929	overall: 0,983

If unsound, the sound behavior is approximated when creating the process tree.

Equal weights for all criteria.

PAGE 30

Conclusion

- First algorithm that allows for balancing all four perspectives.
- Genetic algorithm is very flexible, but also very slow.
- Process trees only used internally (choose your favorite representation)
- Future work:
 - Improve speed
 - Distribute PM tasks
 - Discover configurable process trees

process mining workbench

and data mining necessary to understand the remainder of the book. Part II focuses on process discovery as the most important process mining task. Part III moves beyond discovering the control flow of processes and highlights conformance checking, and organizational and time perspectives. Part IV guides the reader in successfully applying process mining in practice, including an introduction to the widely used open-source tool ProM. Finally, Part V takes a step back, reflecting on the material presented and the key open challenges.

Overall, this book provides a comprehensive overview of the state of the art in process mining. It is intended for business process analysts, business consultants, process managers, graduate students, and BPM researchers.

Features and Benefits:

- First book on process mining, bridging the gap between business process modeling and business intelligence.
- Written by one of the most influential and most-cited computer scientists and the best-known BPM researcher.
- Self-contained and comprehensive overview for a broad audience in academia and industry.
- The reader can put process mining into practice immediately due to the applicability of the techniques and the availability of the open-source process mining software ProM.

Process Mining

Process Mining

Wil M. P. van der Aalst

Discovery, Conformance and Enhancement of Business Processes

www.processmining.org

ISBN 978-3-642-19344-6

www.win.tue.nl/ieeetfpm/

springer.com

Computer Science