When RL and Control meet: lessons learned

Ann Nowé

AI lab Vrije universiteit Brussel





2009-2013







2009-2013











 $r_k = -\Delta p$







- Basic PID to bootstrap RL-algorithm, i.e. to provide initial behavior to the Fuzzy Q-learning algorithm
- RL gradually improves with PID as backup, i.e. when RL proposes a potentially unsafe action, it can be overruled by the PID controller









2009-2013











Safe reinforcement learning for multienergy management systems with known constraint functions

Glenn Ceusters ^{a b c} 久 凶, Luis Ramirez Camargo ^{b d} 凶, Rüdiger Franke ^a凶, Ann Nowé ^c 凶 Maarten Messagie ^b 🖾

VLAIO Baekeland PhD with





ENERGY

1 Input: initialize RL algorithm, initialize constraint functions in sets X and U, initialize safe fallback policy π^{safe}

- 2 for k = 0, 1, 2, ..., do
 - Observe state s and select action a 3
 - if constraint check = True then keep selected action a as safe action a^{safe}
 - else

get safe action a^{safe} from safe fallback policy π^{safe}

end

Algorithm 1: SafeFallback

- Execute asafe in the environment
- Observe next state s', reward r and done signal d to indicate whether s' is terminal
- Give experience tuple (s, a^{safe}, r, s', d) and if $a^{safe} \neq a$: (s, a, r c, s', d) with cost c
- If s' is terminal, reset environment state 8
- end

5

Algorithm 2: GiveSafe		
1 Input: initialize RL algorithm, initialize constraint functions in sets X and U		
2 for $k = 0, 1, 2, \dots$ do		
3 Observe state <i>s</i> and select action <i>a</i>		
4 if constraint check = True then		
keep selected action a as safe action a^{safe}		
else		
s while constraint check = False do		
give experience tuple (s, a, c, s, d) with cost c		
agent selects new action a		
check constraints		
end		
return safe action <i>a^{safe}</i>		
end		
6 Execute <i>a^{safe}</i> in the environment		
7 Observe next state s', reward r and done signal d to indicate whether s' is terminal		
8 Give experience tuple (s, a^{safe}, r, s', d)		
9 If <i>s'</i> is terminal, reset environment state		
end		

Highlights

- RL learns policies that outperform MPC
- A (near-to) optimal multi-energy management policy can be learned safely.
- Constraints can be formulated independently from the (optimal) control technique
- Better policies can be found starting with an initial safe fallback policy.

Ann Nowé

SAFE Reinforcement Learning





03-04-2025 | 7

imec-int.com/safebot

Ann Nowé

CTRL x AI SBO with UGent

Do you have control structures available in your

2. machines/robots/processes? If so, what type of controllers do you typically use right now?

13 respondenten

Classical PID (single-input- single-output, cascaded)	77%	10 stemmen
PID Plus (Feedforward, anti- windup etc.)	46%	6 stemmen
Model Predictive Control (MPC)	23%	3 stemmen
Rule-based control	46%	6 stemmen
Machine Learning (ML) type	23%	3 stemmen
Combination of the above	15%	2 stemmen
Other	0%	0 stemmen



CTRL x Al

SBO with UGent









confidence verification robustness legal compliance performance trust safety education debugability acceptance autonomy



CTRL x AI SBO with UGent



ARTIFICIAL INTELLIGENCE LAB

CTRL x AI

SBO with UGent





Transparency / Interpretability / Explainability

Rule Distillation

• IF COND1 AND COND2 AND THEN ACTION1



IF (1, 5)=BRICK AND (1, 0)=COIN_RED THEN Action=RIGHT



Action Influence Graphs Structural Causal Models with actions $U = \begin{bmatrix} A_a \\ A_a \end{bmatrix}$ $W = \begin{bmatrix} A_a \\ A_a \end{bmatrix}$ $W = \begin{bmatrix} A_a \\ A_a \end{bmatrix}$ $A_b = \begin{bmatrix} A_a \\ A_a \end{bmatrix}$

QuestionWhy not build_barracks (A_b) ?ExplanationBecause it is more desirable to do action
build_supply_depot (A_s) to have more
Supply Depots (S) as the goal is to
have more Destroyed Units (D_u) and De-
stroyed buildings (D_b) .

Madumal, P., Miller, T., Sonenberg, L., Vetere, F.: Explainable reinforcement learning through a causal lens. AAAI'20

I'm building up momentum

IF COND1 AND COND2 AND THEN ACTION1



Visualisation





Local Advantage Networks for Multi-Agent Reinforcement Learning in Dec-POMDPs Raphaël Avalos et al, JMLR (2023) Starcraft Multi-agent Challenge (SMAC)

Post-HOC Feature Importance

Soft Decision Trees (SDT)

- Hybridization of Neural Networks and Decision Trees
 - Originally for image classification
- Branching nodes are perceptrons
 ⇒ soft decisions for tree traversal

$$p_i(x) = \sigma(\beta(xw_i + b_i))$$

• Leaves represent softmax distributions:

$$Q_k^l = \frac{\exp(\phi_k^l)}{\sum_{k'} \exp(\phi_{k'}^l)}$$





Frosst, N., & Hinton, G. (2017). Distilling a Neural Network Into a Soft Decision Tree. In T. R. Besold & O. Kutz (Eds.), *Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML* 2017 (Vol. 2071). Tarek R. Besold & Oliver Kutz.

Post-HOC Feature Importance



Soft Decision Trees (SDT)

- Hybridization of Neural Networks and Decision Trees
 - Originally for image classification
- Branching nodes are perceptrons
 ⇒ soft decisions for tree traversal
 - $p_i(x) = \sigma(\beta(xw_i + b_i))$
- Leaves represent softmax distributions:

$$Q_k^l = \frac{\exp(\phi_k^l)}{\sum_{k'} \exp(\phi_{k'}^l)}$$

Examining branching nodes' filters shows important regions of the state





Coppens, Y., Efthymiadis, K., Lenaerts T., and Nowé A., (2019) Distilling Deep Reinforcement Learning Policies in Soft Decision Trees, XAI workshop at IJCAI 2019.

POST-HOC Policy distillation

- Summarize policy behavior into a surrogate white box model
- Driven by performance rather than model accuracy.



Exploiting meta information



Greedy Policy

Distilled Policy



Inductive Rule Learning



From Interpretable Machine Learning. A Guide for Making Black Box Models Explainable, by Christoph Molnar, 2019. https://christophm.github.io/interpretable-ml-book/. (CC BY-NC-SA 4.0)



Mining rules from multi-labelled data



IF X<=1.0 THEN Class=A
 IF X>=1.01 THEN Class=B



Rule mining on greedy policy

Example Muddy world

greedy policy





Phase 1

2. IF X=19 THEN Class=UP





Example Muddy world

Phase 1

IF X<=18 THEN Class=RIGHT
 IF X=19 THEN Class=UP



Phase 2 1.1 IF X<=18 AND Y=10 AND X>=8 AND X<=9 THEN Class=UP

- 1.2 IF X<=18 AND X>=10 THEN Class=RIGHT
- 1.3 IF X<=18 AND X<=8 THEN Class=RIGHT
- 1.4 IF X<=18 AND Y=11 THEN Class=UP
- 1.5 IF X<=18 AND Y=9 THEN Class=DOWN
- 1.6 IF X<=18 THEN Class=RIGHT
- 2 IF X=19 THEN Class=UP



t. t. t.

Intrinsic transparency : Relational RL



Džeroski, S., De Raedt, L. & Driessens, K. Relational Reinforcement Learning. Machine Learning 43, 7-52 (2001).



Intrinsic transparency



Environment



BNAIC/BeNeLearn 2024

Critic-Moderated Genetic Programming







Selection



Mutation

Crossover

Intrinsic transparency











Reinforcement Learning with Formal Guarantees



Framework for learning discrete latent models of unknown continuous-spaces environment with bisimulation guarantees

- Can be learned by executing an RL policy in the environment
- Yields a distilled version of the RL policy
- New local losses bounds for (i) bisimulation guarantees (ii) discrete setting (iii) action embedding function
- **PAC schemes** to formally retrieve confidence metrics to asses the quality of the learned model



Distillation of RL Policies with Formal Guarantees via Variational Abstraction of Markov Decision Processes Florent Delgrange, Ann Nowé, Guillermo A. Pérez



Wasserstein Auto-encoded MDPs: Formal Verification of Efficiently Distilled RL Policies with Many-sided Guarantees Florent Delgrange, Ann Nowé, Guillermo A. Pérez

