

Reinforcement Learning for Dynamic Strategies in TCT

Manuel Schneckeneither

(supervised by Prof. Dr. Georg Moser)

University of Innsbruck, Austria

CL Seminar Summer 2020



Motivation

The Tyrolean Complexity Tool, i.e. TCT ^[1]

- Automatic tool to infer runtime (and derivational) complexity
- Allows (first or higher-order) functional programs as generalised form of term rewrite systems as input
- In case of success: TCT outputs a certification proving the presented asymptotic complexity



[1] Manuel Schneckenreither, “Dynamic Strategies for TCT ”, 2020

[2] Landi, “Undecidability of static analysis”, 1992

Motivation

The Tyrolean Complexity Tool, i.e. TCT [1]

- Automatic tool to infer runtime (and derivational) complexity
- Allows (first or higher-order) functional programs as generalised form of term rewrite systems as input
- In case of success: TCT outputs a certification proving the presented asymptotic complexity

Static Analyses

- Static analyses are undecidable [2]
- Thus: Some sort of creativity is required (Strategies)
- Idea: Reinforcement learning for dynamic strategies

[1] Manuel Schneckenreither, “Dynamic Strategies for TCT ”, 2020

[2] Landi, “Undecidability of static analysis”, 1992

Reinforcement Learning for Dynamic Strategies in TCT

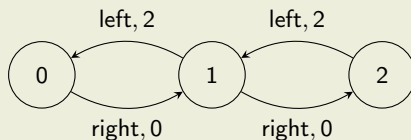
Table of Contents

1. Motivation
2. Discounted Reinforcement Learning (Discounted RL)
3. Average Reward Adjusted Discounted RL
4. Proposed Implementation in TCT



Reinforcement Learning (RL)

Reinforcement Learning Processes [3]

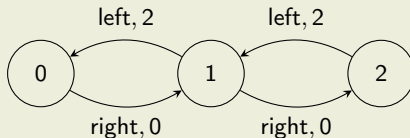


- Based on dynamic programming: states, actions, rewards

[3] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Reinforcement Learning (RL)

Reinforcement Learning Processes [3]



- Agent iteratively assesses state values $V_{\gamma}^{\pi}(s)$, where
- $0 < \gamma < 1$ is the discount factor, π the policy function, and
- for two consecutive observed states s_t, s_{t+1} with reward r_t

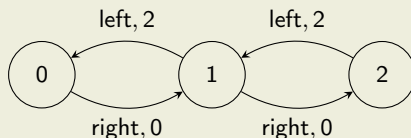
$$V_{\gamma}^{\pi}(s_t) \stackrel{\alpha}{\leftarrow} r_t + \gamma V_{\gamma}^{\pi}(s_{t+1})$$

($\stackrel{\alpha}{\leftarrow}$ is exp. smoothed updating with rate α)

[3] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Reinforcement Learning (RL)

Reinforcement Learning Processes [3]



- Thus, state values $V_{\gamma}^{\pi}(s)$ are

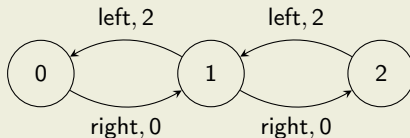
$$V_{\gamma}^{\pi}(s) = \lim_{N \rightarrow \infty} E\left[\sum_{t=0}^{N-1} \gamma^t R_t^{\pi}(s)\right]$$

where $R_t^{\pi}(s)$ the reward received at time t upon starting in state s by following policy π .

[3] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Reinforcement Learning (RL)

Reinforcement Learning Processes [3]



- RL in summary: Iteratively solving constraint problems composed of $|\mathcal{S}|$ constraints and $|\mathcal{S}|$ variables, where $s \in \mathcal{S}$:

$$V_{\gamma}^{\pi}(0) = 0 + \gamma V_{\gamma}^{\pi}(1) \qquad V_{\gamma}^{\pi}(2) = 2 + \gamma V_{\gamma}^{\pi}(1)$$

$$\text{left: } V_{\gamma}^{\pi}(1) = 2 + \gamma V_{\gamma}^{\pi}(0) \quad \text{right: } V_{\gamma}^{\pi}(1) = 0 + \gamma V_{\gamma}^{\pi}(2)$$

($|\cdot|$ is the size of a set)

[3] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Aim in Discounted Reinforcement Learning

Aim of Discount Reinforcement Learning

- Finding an optimal policy π^*
- π^* maximises the state value for all states s as compared to any other policy π :

$$V_{\gamma}^{\pi^*}(s) - V_{\gamma}^{\pi}(s) \geq 0$$



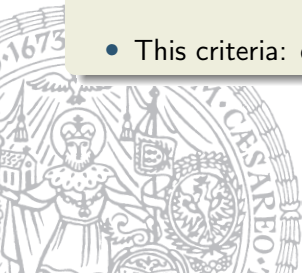
Aim in Discounted Reinforcement Learning

Aim of Discount Reinforcement Learning

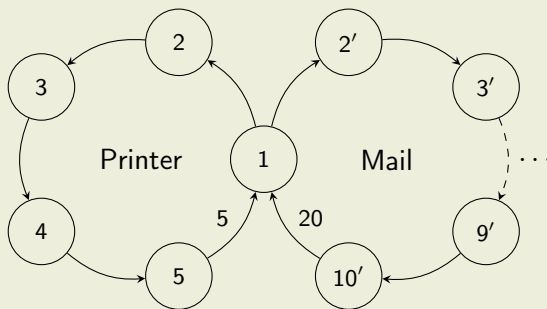
- Finding an optimal policy π^*
- π^* maximises the state value for all states s as compared to any other policy π :

$$V_{\gamma}^{\pi^*}(s) - V_{\gamma}^{\pi}(s) \geq 0$$

- This criteria: discounted-optimality as γ is **fixed**



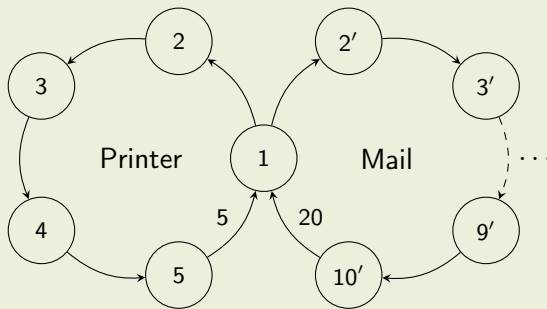
The issue of discounted-optimality illustrated [4]



- Average reward received is
 - 1 for the printer-loop
 - 2 for the mail-loop

[4] Adapted from Mahadevan, "Optimality criteria in reinforcement learning", 1996

The issue of discounted-optimality illustrated [4]



- Average reward received is
 - 1 for the printer-loop
 - 2 for the mail-loop
- BUT: if $\gamma < 3^{-\frac{1}{5}} \approx 0.8027$ the agent prefers the printer loop

[4] Adapted from Mahadevan, "Optimality criteria in reinforcement learning", 1996

Laurent Series Expansion of Discounted State Values

Laurent Series Expansion of Discounted State Values^[5]

The Laurent series expansion of $V_\gamma^\pi(s)$ of discounted state values:

$$V_\gamma^\pi(s) = \frac{\rho^\pi}{1-\gamma} + V^\pi(s) + e_\gamma^\pi(s)$$

Average Reward ρ^π (simplified for unichain Processes)

The **average reward** is defined as

$$\rho^\pi = \lim_{N \rightarrow \infty} \frac{\mathbb{E}[\sum_{t=0}^{N-1} R_t^\pi]}{N}$$

where R_t^π is the reward received at time t .

[5] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Laurent Series Expansion of Discounted State Values

Laurent Series Expansion of Discounted State Values^[5]

The Laurent series expansion of $V_\gamma^\pi(s)$ of discounted state values:

$$V_\gamma^\pi(s) = \frac{\rho^\pi}{1-\gamma} + \textcolor{red}{V}^\pi(s) + e_\gamma^\pi(s)$$

Bias value $V^\pi(s)$

The **bias value** is defined as

$$V^\pi(s) = \lim_{N \rightarrow \infty} \mathbb{E} \left[\sum_{t=0}^{N-1} (R_t^\pi(s) - \rho^\pi) \right]$$

It is the additional reward that sums up when starting in state s .

[5] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Laurent Series Expansion of Discounted State Values

Laurent Series Expansion of Discounted State Values^[5]

The Laurent series expansion of $V_\gamma^\pi(s)$ of discounted state values:

$$V_\gamma^\pi(s) = \frac{\rho^\pi}{1-\gamma} + V^\pi(s) + e_\gamma^\pi(s)$$

Error term of e_γ^π ^[6]

Error term $e_\gamma^\pi(s)$ consists of infinitely many terms, but

$$\lim_{\gamma \rightarrow 1} e_\gamma^\pi(s) = 0$$

If $\gamma < 1$ it takes number of steps and reward values into account.

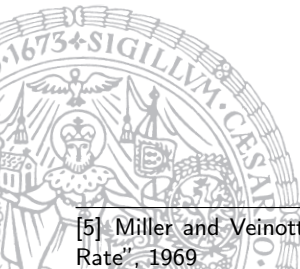
[6] Puterman, "Markov Decision Processes: Discrete Stochastic Dynamic Programming", 1994

Laurent Series Expansion of Discounted State Values

Laurent Series Expansion of Discounted State Values^[5]

The Laurent series expansion of $V_\gamma^\pi(s)$ of discounted state values:

$$V_\gamma^\pi(s) = \frac{\rho^\pi}{1-\gamma} + V^\pi(s) + e_\gamma^\pi(s)$$



[5] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

Laurent Series Expansion of Discounted State Values

Laurent Series Expansion of Discounted State Values^[5]

The Laurent series expansion of $V_\gamma^\pi(s)$ of discounted state values:

$$V_\gamma^\pi(s) = \frac{\rho^\pi}{1-\gamma} + V^\pi(s) + e_\gamma^\pi(s)$$

Note

- Recall: High $\gamma \approx 1$ values required for optimal policies
- First term converges to ∞ as $\gamma \rightarrow 1$
- Recall: RL is an iterative method

[5] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

So What?

$$V_{\gamma}^{\pi}(s) = \frac{\rho^{\pi}}{1 - \gamma} + V^{\pi}(s) + e_{\gamma}^{\pi}(s)$$



So What?

$$V_{\gamma}^{\pi}(s) = \frac{\rho^{\pi}}{1 - \gamma} + V^{\pi}(s) + e_{\gamma}^{\pi}(s)$$

Consider this simple gridworld problem [6]

		up	
(0,0)	left	(0,1)	right
random		down	
up		up	
left	(1,0)	right	(1,1)
	down		down

State/Action Space

		$\mathcal{U}(0,8)-1$	
(0,0)	$\mathcal{U}(0,8)$	(0,1)	$\mathcal{U}(0,8)-1$
10		$\mathcal{U}(0,8)$	
$\mathcal{U}(0,8)$		$\mathcal{U}(0,8)$	
$\mathcal{U}(0,8)-1$	(1,0)	$\mathcal{U}(0,8)$	(1,1)
	$\mathcal{U}(0,8)-1$		$\mathcal{U}(0,8)-1$

Reward Function

[6] Manuel Schneckeneither, "Average Reward Adjusted Discounted Reinforcement Learning: Near-Blackwell-Optimal Policies for Real-World Applications", 2020

So What?

$$V_{\gamma}^{\pi}(s) = \frac{\rho^{\pi}}{1 - \gamma} + V^{\pi}(s) + e_{\gamma}^{\pi}(s)$$

Same problem as 5x5 grid

	23.113	162.243	354.839	398.391
(0,0) 101.260	67.160 (0,1) 61.783 49.812	49.812 (0,2) 349.515 306.761	306.761 (0,3) 399.560 399.376	399.376 (0,4) 398.444
	37.331	153.653	351.907	399.539
...

- All states values completely assessed independently
- For $\gamma \approx 1$: $\rho^{\pi}/(1 - \gamma) \gg V^{\pi}(s) + e_{\gamma}^{\pi}(s)$
- Thus, all values need to increase to $\approx \rho^{\pi}/(1 - \gamma)$
- BUT: Greater $V_{\gamma}^{\pi}(s)$, then more likely to be picked

Average Reward Adjusted Reinforcement Learning

Basic Idea

Separately assessing

- average reward ρ^π
- bias value $V^\pi(s)$



Average Reward Adjusted Reinforcement Learning

Basic Idea

Separately assessing

- average reward ρ^π
- bias value $V^\pi(s)$

Algorithm

$$\rho^\pi \stackrel{\alpha}{\leftarrow} r_t + V^\pi(s_{t+1}) - V^\pi(s_t)$$

$$V^\pi(s_t) \stackrel{\beta}{\leftarrow} r_t + \gamma_1 V^\pi(s_{t+1}) - \rho^\pi$$

Note: We can set $\gamma_1 = 1$ here, as the subtraction of the average reward bounds the state values.

Discounted vs. Average Reward Adjusted RL

Reconsider the 5x5 Gridworld [6]

Algorithm	Sum Reward	Avg. Steps
Avg. Rew. Adjusted $\gamma_1 = 0.99$	51894.094	5.039
Avg. Rew. Adjusted $\gamma_1 = 0.999$	51878.069	5.063
Avg. Rew. Adjusted $\gamma_1 = 1.00$	51856.529	5.055
Standard Discounted $\gamma = 0.99$	34409.464	7661.833
Standard Discounted $\gamma = 0.999$	33931.917	7379.155
Standard Discounted $\gamma = 0.50$	30171.837	9999.000

(500k learning steps, $40 \times 10k$ greedy evaluation steps)

[6] Manuel Schneckeneither, "Average Reward Adjusted Discounted Reinforcement Learning: Near-Blackwell-Optimal Policies for Real-World Applications", 2020

Average Reward Adjusted Reinforcement Learning

More Insights

Average Reward Adjusted Reinforcement Learning ...

1. requires fewer number of learning steps
2. allows more natural specification of reward function (continuous feedback), which helps to
 - find the goal
 - distinguish between good and bad sequences of actions

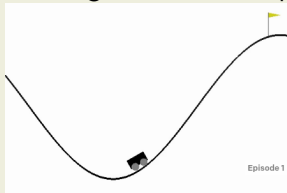


Average Reward Adjusted Reinforcement Learning

More Insights

Average Reward Adjusted Reinforcement Learning ...

1. requires fewer number of learning steps
2. allows more natural specification of reward function (continuous feedback), which helps to
 - find the goal
 - distinguish between good and bad sequences of actions



Source: <https://gym.openai.com/>

Average Reward Adjusted Reinforcement Learning

More Insights (cont'd)

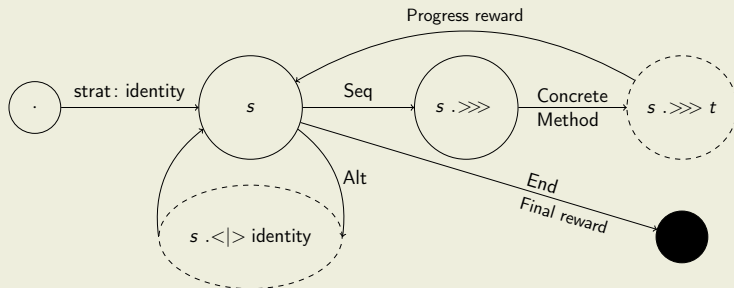
Average Reward Adjusted Reinforcement Learning ...

3. state values are easier to interpret
4. normed state values have a higher spread
 - important when approximating the state-value function (ANN)



Implementation in $T_C T$

Process Schema [7]

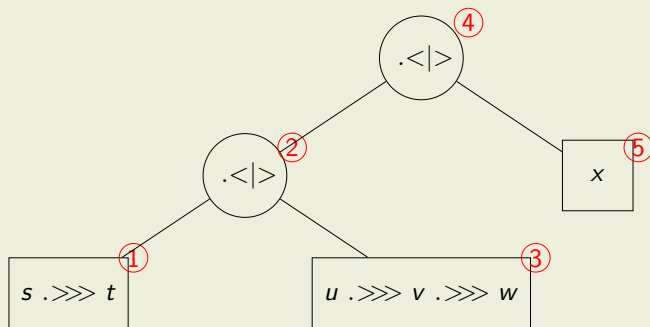


- Starting strategy: identity
- Strategy operators: Sequencing $.>>>$ and alternative $.<|>$
- Two agents: decomposition and complexity analyser

[7] Manuel Schneckeneither, "Dynamic Strategies for $T_C T$ ", 2020

Implementation in $T_C T$

Binary Strategy Tree [8]



- No ambiguity by enforcing right-associativity of $.<|>$

[8] Manuel Schneckeneither, "Dynamic Strategies for $T_C T$ ", 2020

Implementation in TCT

State Space

- Characteristics of input problem, e.g.
 - number of rules
 - number of (root) symbols
 - left-linearity, right-linearity
- Representation of current strategy
 - Counter specifying steps until last occurrence of concrete method



[8] Manuel Schneckeneither, “Dynamic Strategies for TCT ”, 2020

Conclusion

Summary

- Discounted Reinforcement Learning
- Average Reward Adjusted Discounted Reinforcement Learning
- Implementation Idea for Dynamic Strategies in T_C^T



Conclusion

Summary

- Discounted Reinforcement Learning
- Average Reward Adjusted Discounted Reinforcement Learning
- Implementation Idea for Dynamic Strategies in $T_C T$

Thank you for your attention!





Implementation in $T_C T$

Reward Function [9]

Assume current strategy s_t at step t , a timeout of T seconds, a measured execution time of $t' \leq T$ seconds and resulting polynomial complexity C :

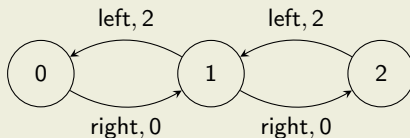
$$r(s_t) = \begin{cases} cT + t' & \text{if } C \text{ is } O(n^c) \text{ and } c < P_{\max} \\ P_{\max} T & \text{otherwise} \end{cases}$$

where P_{\max} is the maximum polynomial degree being considered.

[9] Manuel Schneckener, “Dynamic Strategies for $T_C T$ ”, 2020

Reinforcement Learning (RL)

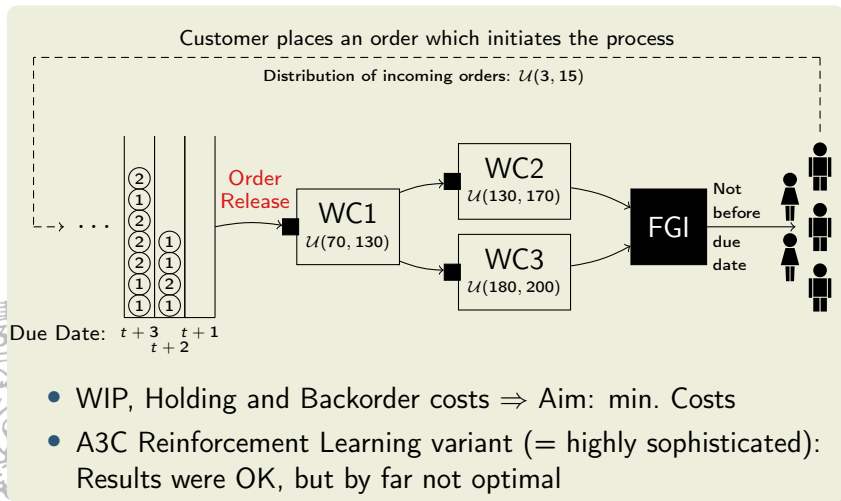
Markov Decision Processes ^[10]



- Based on dynamic programming: states, actions, rewards
- Markov Prop.: transitions with probability $p(s_{t+1}, r_t \mid s_t, a_t)$
- RL processes: *Markov decision processes* (MDPs)

[10] Miller and Veinott, "Discrete Dynamic Programming with a Small Interest Rate", 1969

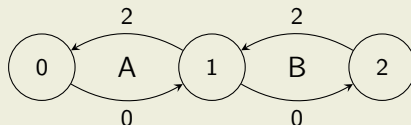
Motivation: More Applications



M. Schneckenreither and Haeussler, "RL Methods for OR Applications", 2019

Add: Normed state values have a higher spread (Slide 13)

Reconsider the Example



- Both policies are have an average reward of 1 (gain-optimal).
- But only π_A is bias-optimal: $V^{\pi^*}(s) - V^{\pi}(s) \geq 0$.

s	$V^{\pi_A}(s)$	$V^{\pi_B}(s)$	$V_{0.99}^{\pi_A}$	$V_{0.99}^{\pi_B}$
0	-0.5	-1.5	99.4975	98.5025
1	0.5	-0.5	100.5025	99.4975
2	1.5	0.5	101.4975	100.5025