

Deep Reinforcement Learning for Integer Programming

Akshay Ajagekar and Alyf Janmohamed

Potential Outline

1. Brief overview of MDP (e.g. our reward)
2. Methods
 - a. Graph Neural Networks
 - b. DQN
 - c. DDQN
 - d. Prioritized Replay
 - e. Training Details
3. Results
 - a. Graph of performance
 - b. Table of performance

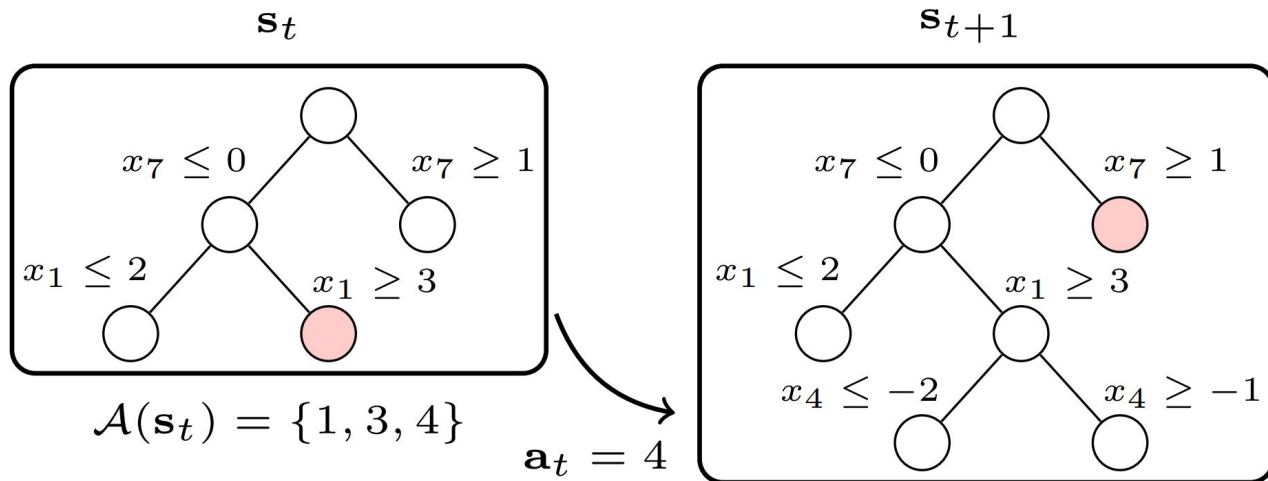
MDP Overview

State

Action

Reward

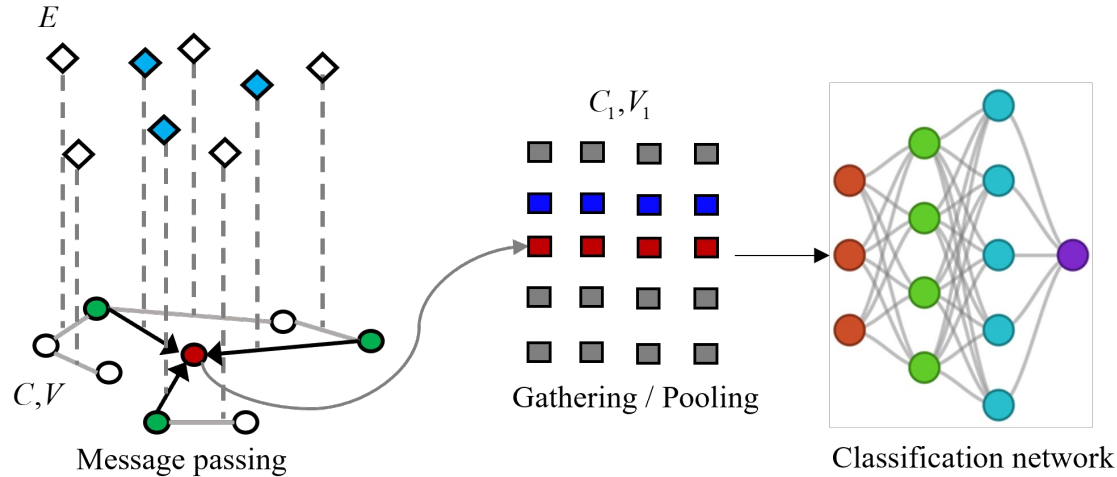
Transition



M. Gasse, D. Ch'etelat, N. Ferroni, L. Charlin, and A. Lodi, "Exact combinatorial optimization with graph convolutional neural networks," arXiv preprint arXiv:1906.01629, 2019.

Methods : Graph Neural Networks

- GNN used to parameterize state-action value function $Q(s,a)$



- Two half convolutions performed
- Layer wise normalization applied

- Edge-conditioned filters used in the MPNN model

$$x'_i = \left(x_i W_{root} + b \right) + \sum_{j \in N(i)} x_j MLP(e_{ji})$$

Methods : Algorithms

Deep Q-Network (DQN)

- $Q(s, a; \theta)$ and $Q(s, a; \theta')$
- $Y_t^{DQN} = r_{t+1} + \gamma \max Q(s_{t+1}, a; \theta')$
- $L_{DQN} = E_{(s,a,r,s')} \left[\left(Q(s_t, a_t; \theta) - Y_t^{DQN} \right)^2 \right]$

- Soft update

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta'$$

Double Deep Q-Network (DDQN)

- $Q(s, a; \theta)$ and $Q(s, a; \theta')$
- $Y_t^{DDQN} = r_{t+1} + \gamma Q(s_{t+1}, \underset{a}{\operatorname{argmax}} Q(s_{t+1}, a; \theta); \theta')$
- $L_{DDQN} = E_{(s,a,r,s')} \left[\left(Q(s_t, a_t; \theta) - Y_t^{DDQN} \right)^2 \right]$

Algorithm 1 Q-Learning

```
1: 1. Initialisation:  
   Load a simulation environment: price series, fill probability;  
   Initialise the value function  $V_0$  and set the parameters:  $\alpha, \epsilon$ ;  
2: 2. Optimisation:  
3: for episode = 1, 2, 3... do  
4:   for t = 1, 2, 3... T do  
5:     Observe current state  $s_t$ ;  
6:     Take an action  $a_t(Q_t, s_t, \epsilon)$ ;  
7:     Observe new state  $s_{t+1}$ ;  
8:     Receive reward  $r_t(s_t, a_t, s_{t+1})$ ;  
9:     Update value function using  $r_t$  and current estimate  $Q_t$ :  
10:    Compute targets and update  $Q(s, a)$   
11:   end for
```

Methods - Training Details

- Implementation :



+



Spektral

- Exploration performed by epsilon greedy strategy where the exploration rate is decayed at rate of 0.995
- Agent's memory buffer of size 5000 and sampled experiences with a batch size of 120.
- $\gamma = 0.99$, $\tau = 0.001$

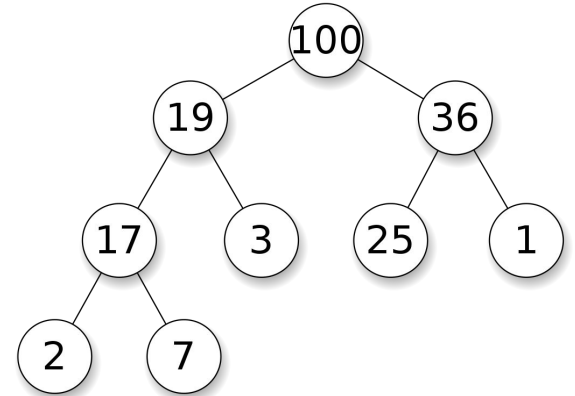
- Data recording :



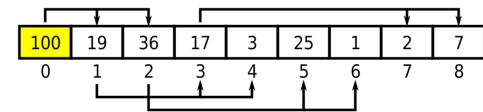
Methods: Prioritized Replay (PER)

- Experience replay smoothes the training distribution over the previous history of the RL agent
- Prioritized replay is a method to weight the training distribution by the TD-error
- Can be implemented using a binary heap

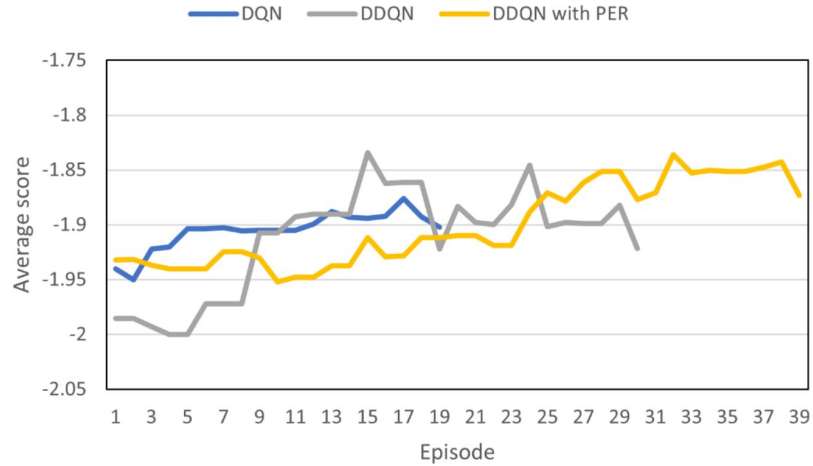
Tree representation



Array representation



Results



Model	Time (s)	Nodes
SCIP	1.21 ± 0.74	$16 \pm 25\%$
DQN	1.12 ± 0.3	$5 \pm 40\%$
DQN with PER	1.41 ± 0.88	$10 \pm 71\%$
DDQN	0.92 ± 0.23	$5 \pm 45\%$
DDQN with PER	1.04 ± 0.35	$6 \pm 60\%$

Further Work

1. Alternate reward functions

- Can consider alternate metrics such as time at the node and information gained (e.g. improvement in bounds)?

2. Rules that generalize between problems

- Can we identify rules that generalize beyond their class of problems? (e.g. set cover, capacitated facility location, max independent set)

3. Can we expand the decision space?

- Node selection, cutting planes, etc.