



Week 11: Adversarial Reinforcement Learning

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

COMP90073
Security Analytics

Yi Han, CIS

Semester 2, 2021

- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

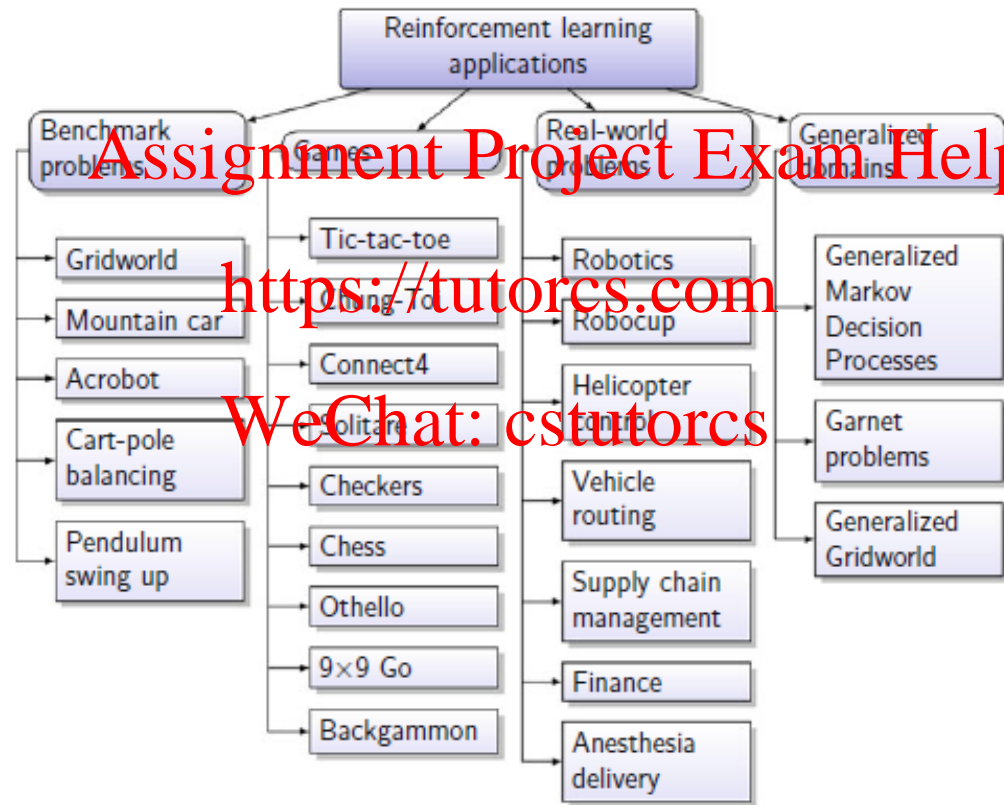
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

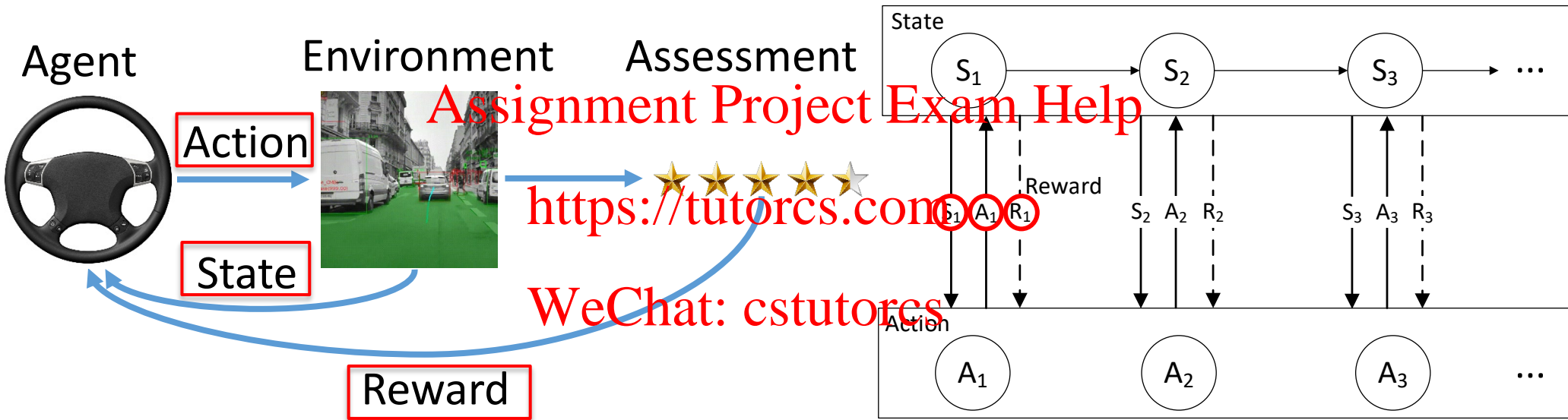
Background on reinforcement learning

- Application



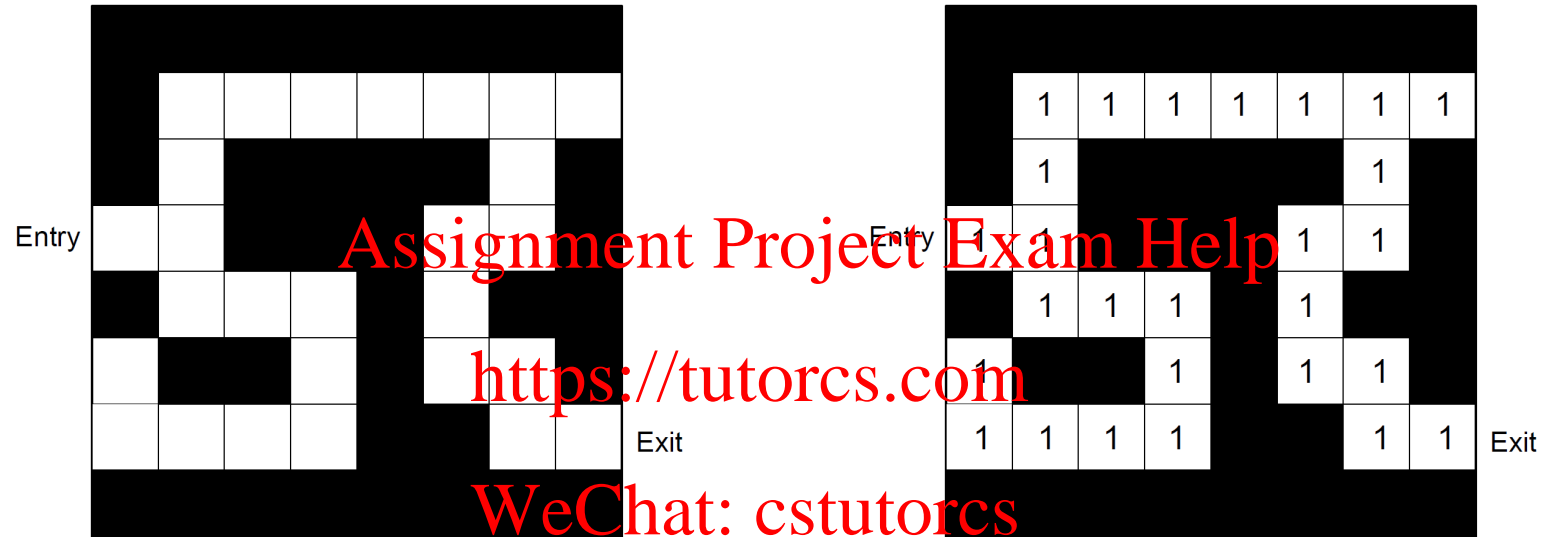
<https://www.myrealfacts.com/2019/05/applications-of-reinforcement-learning.html>

- Introduction



Maximise the discounted cumulative rewards over the long run: $R_t = \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau}, \gamma \in (0, 1]$

- State



$$\begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 2 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 2 & 0 & 0 & 0 & 0 & 1 & 0 \\ 2 & 2 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 \end{bmatrix}$$

Background on reinforcement learning

- Action

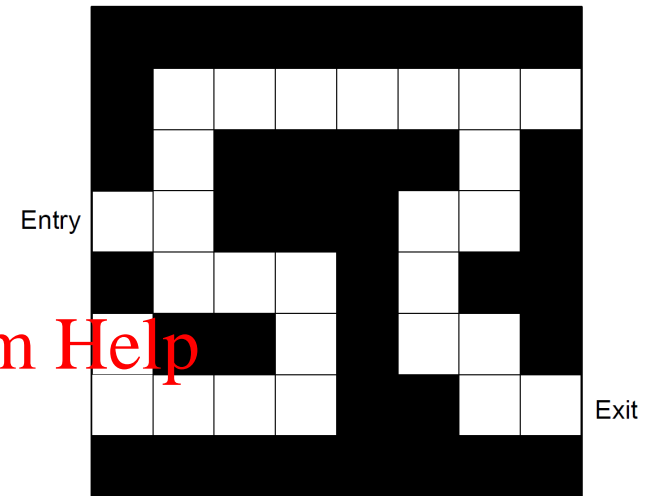
- Up
- Left
- Down
- Right

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

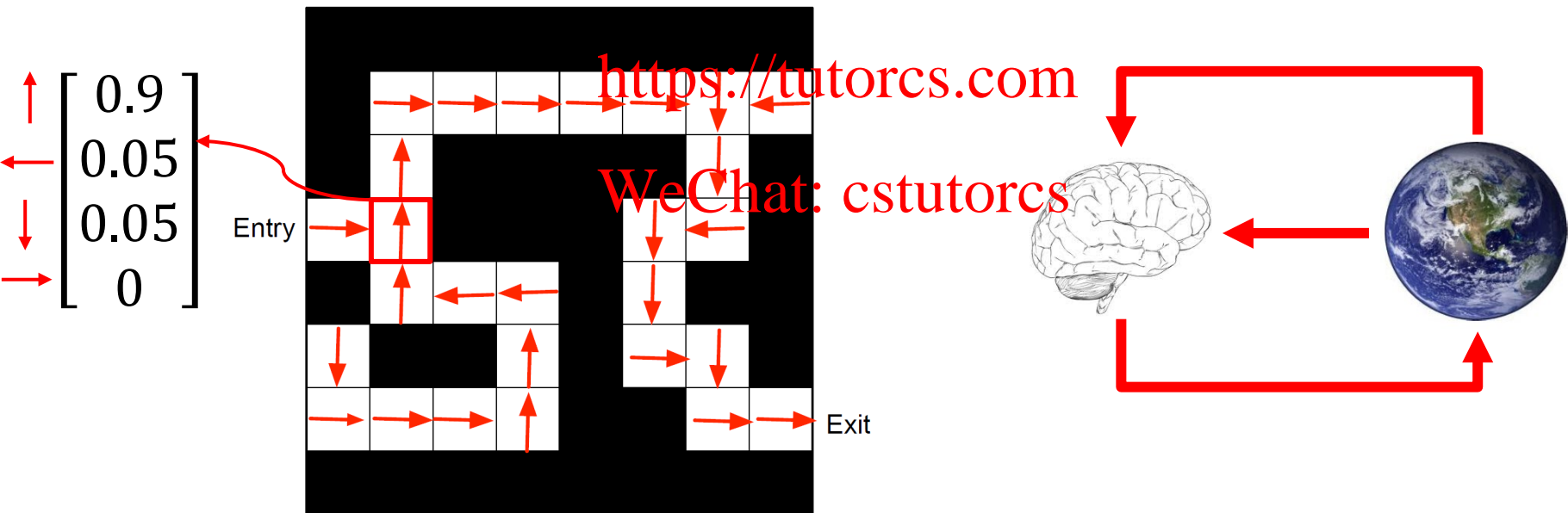
- Reward: an immediate feedback on whether an action is good
 - In the range of $[-1, 1]$
 - 1: reach the exit
 - -0.8: move to a blocked cell
 - -0.3: move to a visited cell
 - -0.05: move to an adjacent cell



Background on reinforcement learning

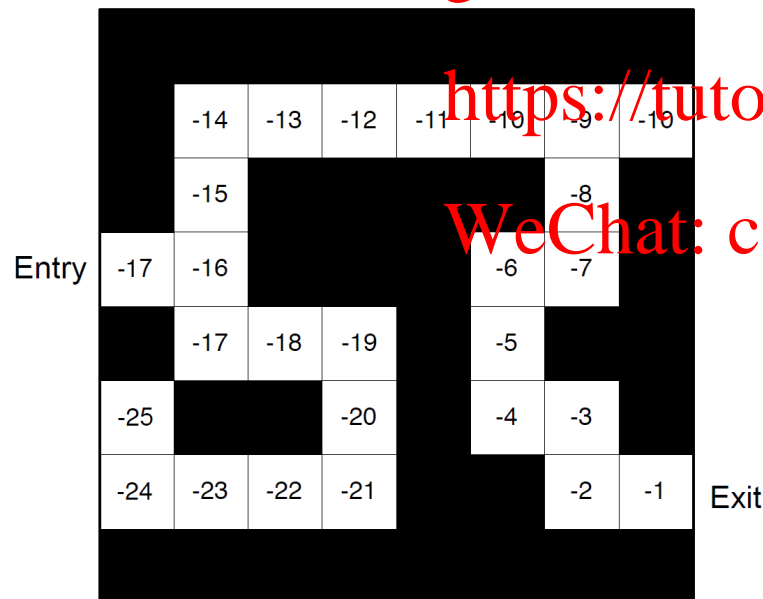
- Policy (π): a mapping from state to action, i.e. $a = \pi(s)$, it tells the agent what to do in a given state

Assignment Project Exam Help

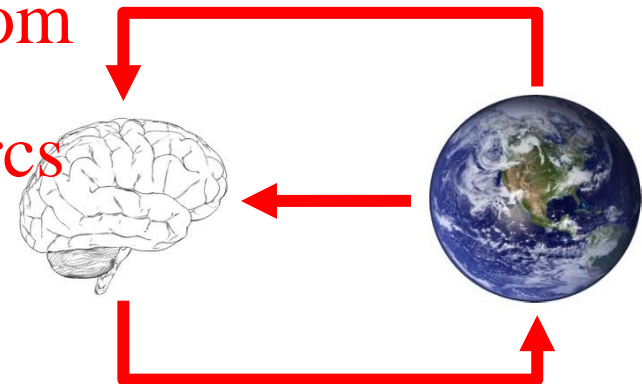


Background on reinforcement learning

- Value function: the future, long term reward of a state
 - Value of state s **under policy** π : $V^\pi(s) = \mathbb{E}[\sum_{i=1}^T \gamma^{i-1} r_i | S_t = s]$
 - Conditional on some policy π
 - Expected value of following policy π starting from state s



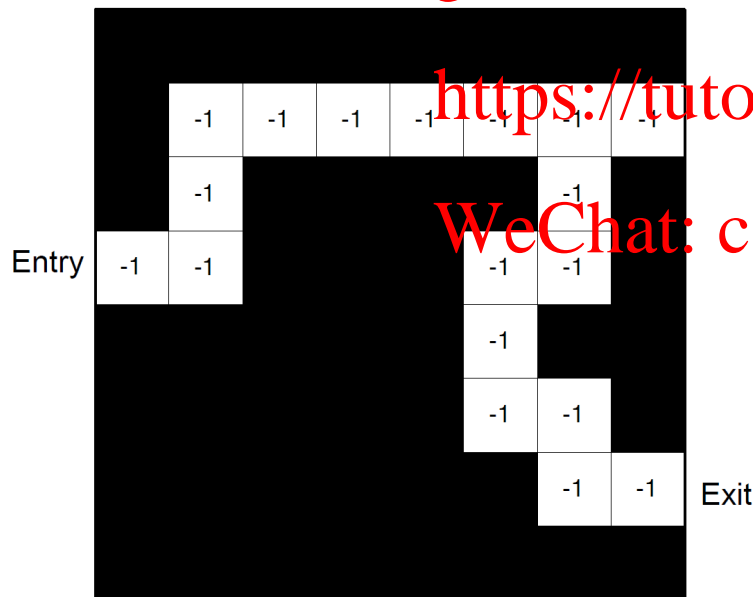
<https://tutorcs.com>
WeChat: cstutorcs



Background on reinforcement learning

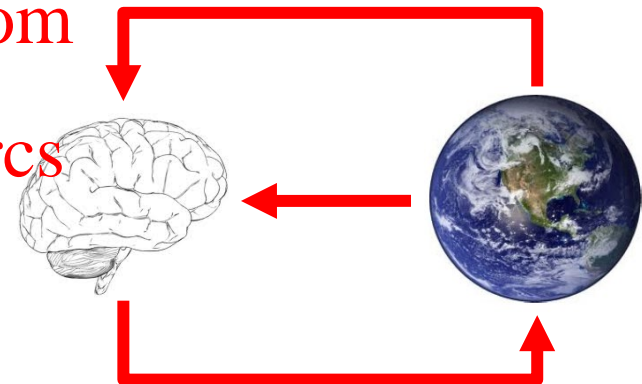
- Model of the environment: mimic the behaviour of the environment, e.g., given a state & action, what the next state & reward might be.

Assignment Project Exam Help

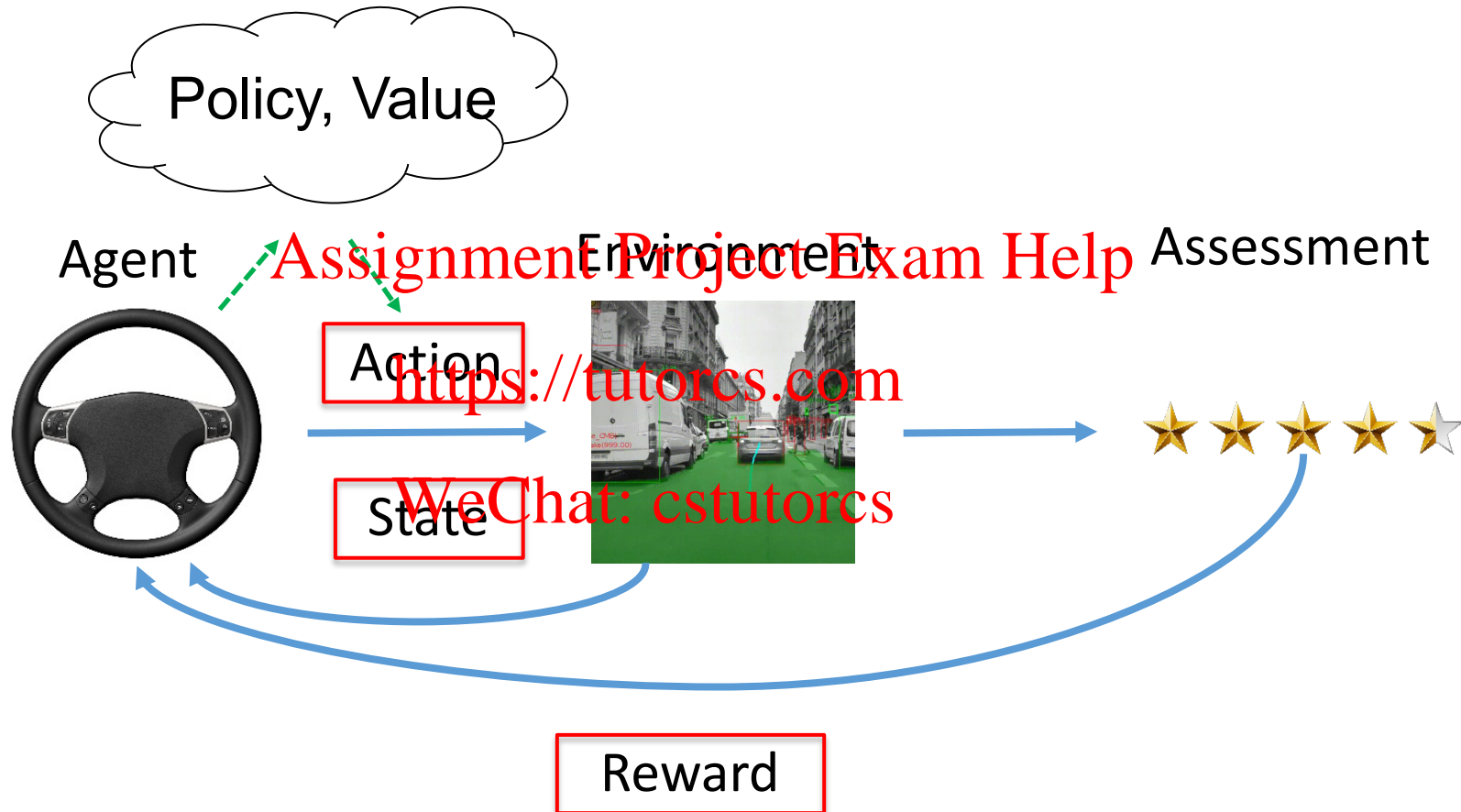


<https://tutorcs.com>

WeChat: cstutorcs



Background on reinforcement learning

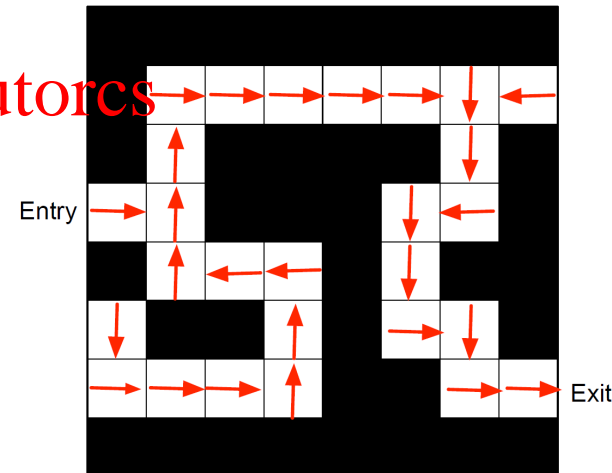
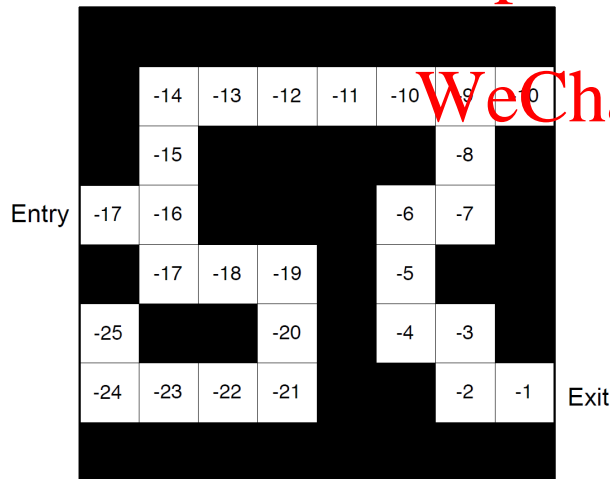


- Classification

- **Value-based** algorithm estimates the value function
- **Policy-based** algorithm learns the policy directly
- **Actor-critic**: critic updates action-value function, actor updates policy

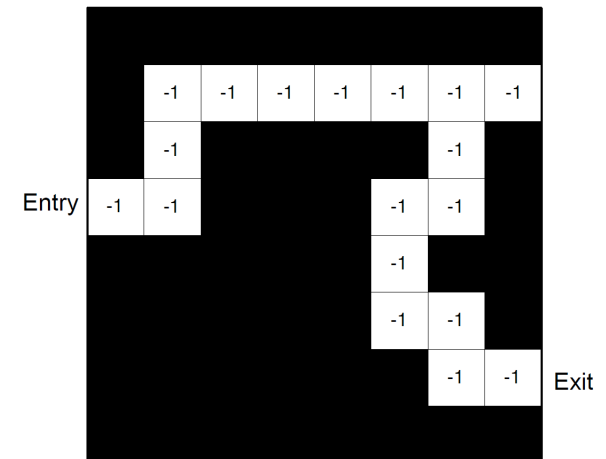
<https://tutorcs.com>

WeChat: cstutorcs



- # Assignment Project Exam Help

WeChat: cstutorcs



- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

Assignment Project Exam Help

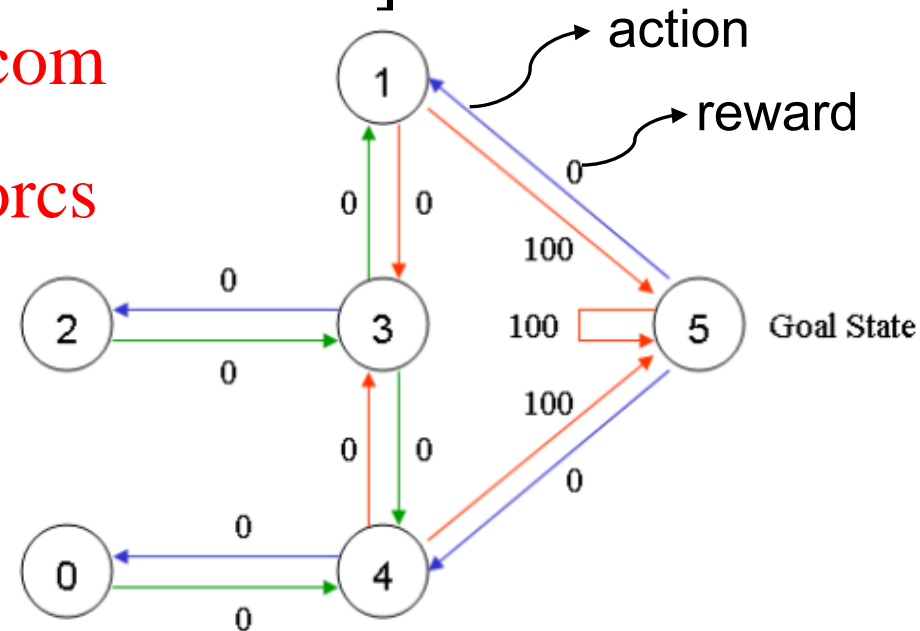
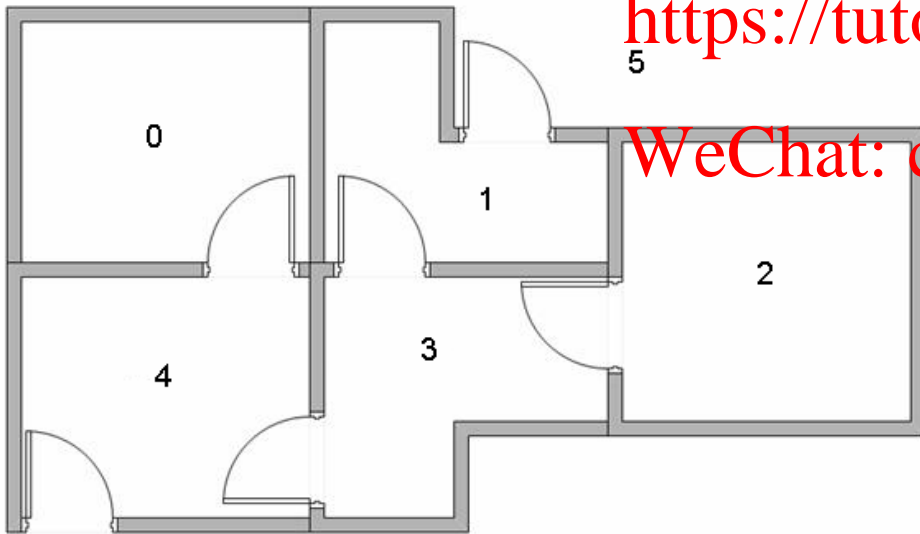
<https://tutorcs.com>

WeChat: cstutorcs

- $$Q^\pi(s, a) = \mathbb{E} \left[\sum_{i=1}^T \gamma^{i-1} r_i \mid S_t = s, A_t = a \right]$$

<https://tutorcs.com>

WeChat: cstutorcs



COMP90073 Security Analysis

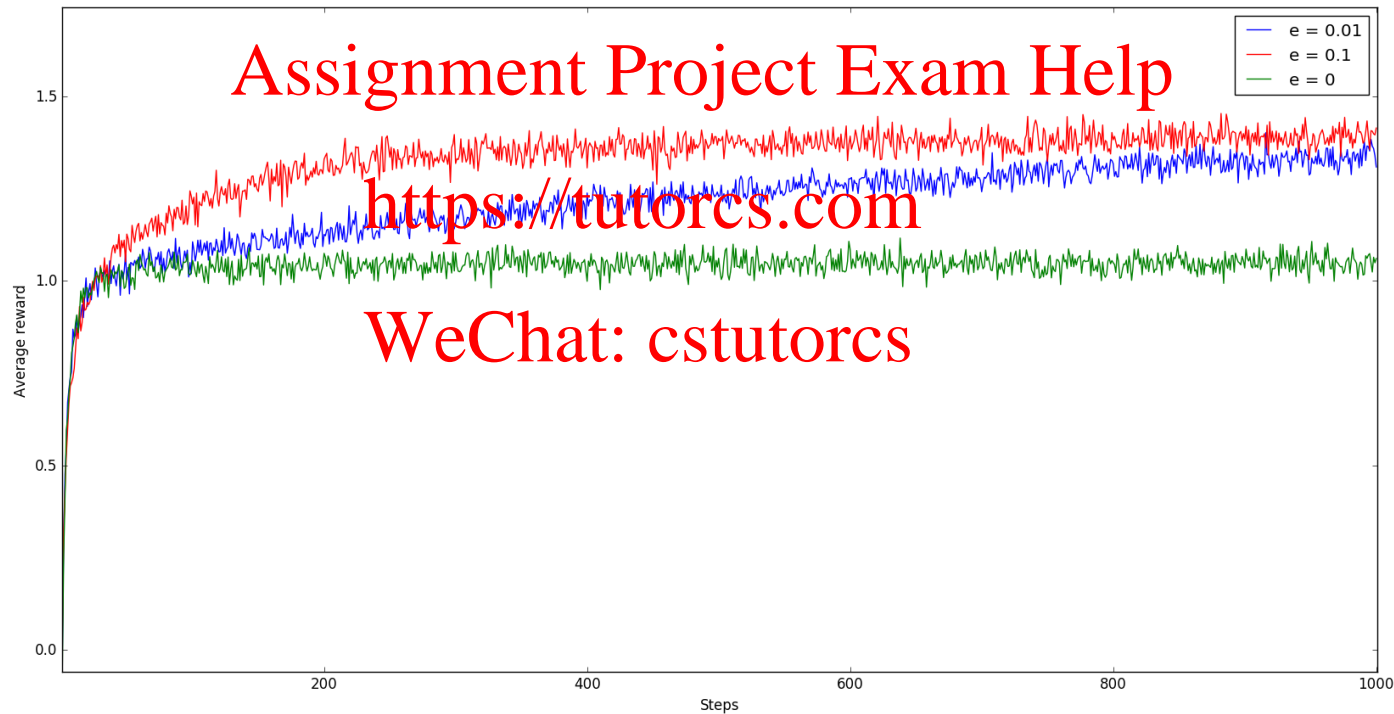
- Q-learning

Assignment Project Exam Help
https://tutorcs.com
WeChat: cstutorcs

State	Action												
	0	1	2	3	4	5		0	1	2	3	4	5
0	-1	-1	-1	-1	0	-1	0	0	0	0	0	400	0
1	-1	-1	-1	0	-1	100	1	0	0	0	320	100+0.8*max(0,0,0)	500
2	-1	-1	-1	0	-1	-1	2	0	0	0	320	0	0
3	-1	0	0	-1	0	-1	3	0	400	256	0	0+0.8*max(0,100)	400
4	0	-1	-1	0	-1	100	4	320	0	0	320	0	500
5	-1	0	-1	-1	0	100	5	0	400	0	0	400	500

$$Q(s_t, a_t) \leftarrow \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}}$$

- Exploitation vs. Exploration
 ϵ -greedy



- The tabular version does not scale with the action/state space
- Classical Q Network [1]
 - Function approximation
 - Approximate $Q(s, a)$ via a neural network: $Q(s, a) \approx Q^*(s, a, \theta)$
 - $L(\theta) = \mathbb{E} \left[\left(r + \gamma \max_a Q(s', a'; \theta) - Q(s, a; \theta) \right)^2 \right]$
 - Unstable

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

target Q

- Deep Q Network (DQN) [2]
 - **Experience replay**: draw randomly from a buffer of (s, a, s', r)
 - $Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a', \theta^-)$ (s', a', s'', r')
 - $L(\theta) = \mathbb{E} \left[\left(r + \gamma \max_{a'} Q(s', a', \theta^-) - Q(s, a; \theta) \right)^2 \right]$
 - Reward clipped to $[-1, 1]$
- Double DQN (DDQN) [3]
 - Separate action selection from action evaluation
 - $Q_1(s, a) \leftarrow r + \gamma Q_2(s', \arg \max_{a'} Q_1(s', a'))$
 - $L(\theta) = \mathbb{E} \left[\left(r + \gamma \max_{a'} Q_2(s', \arg \max_{a'} Q_1(s', a')) - Q_1(s, a; \theta) \right)^2 \right]$

- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Background on reinforcement learning

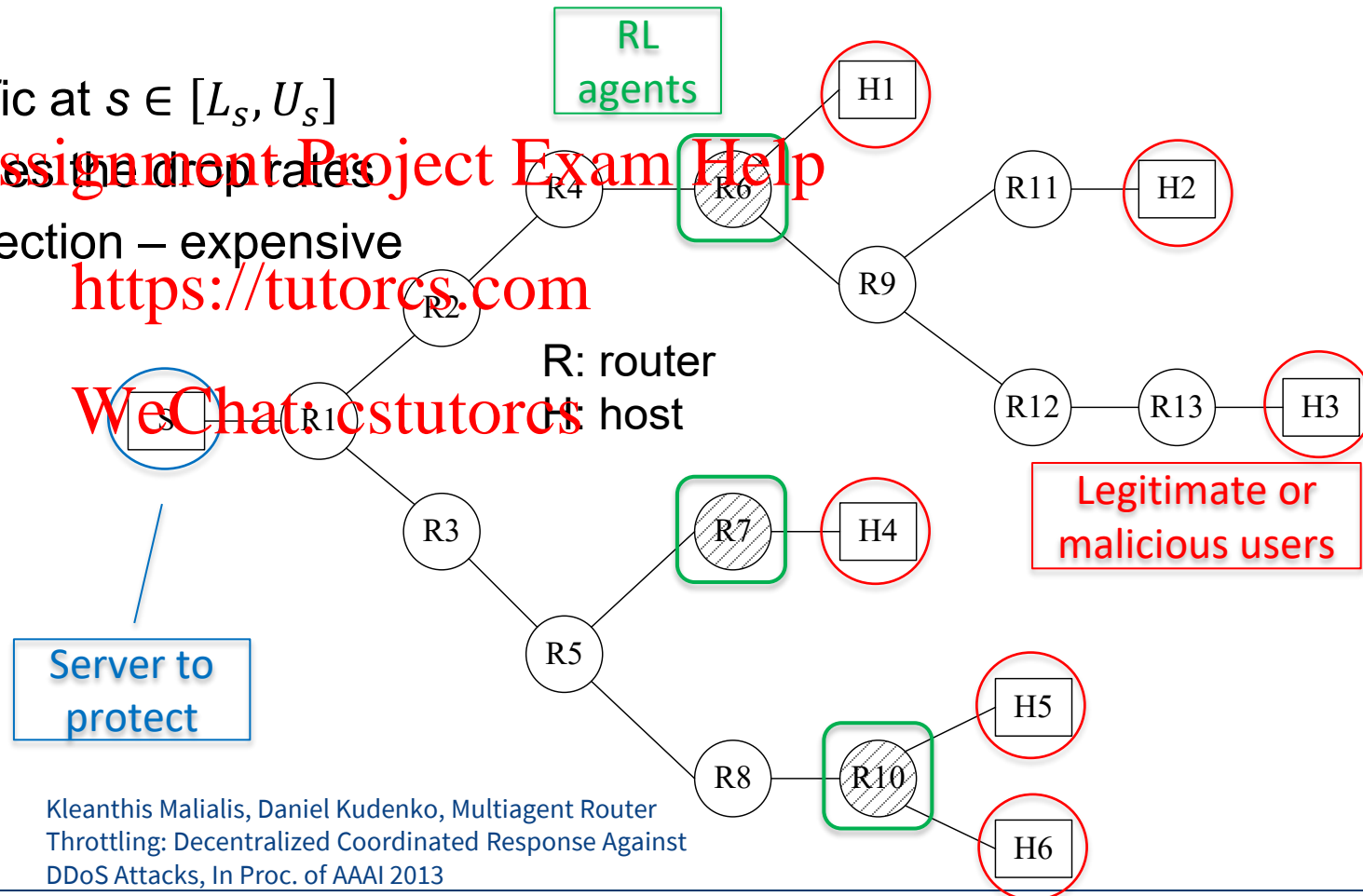
- Distributed Denial-of-Service (DDoS) attacks still occur almost every hour globally
 - <http://www.digitalattackmap.com/>
 - Statistics are gathered by Arbor's Active Threat Level Analysis System from 330+ ISP customers with 130 Tbps of global traffic



- Can RL be applied to throttle flooding DDoS attacks?

- Problem setup [5]

- A mixed set of legitimate users & attackers
- Aggregated traffic at $s \in [L_s, U_s]$
- RL agents decide the drop rates
- No anomaly detection – expensive



Kleanthis Malialis, Daniel Kudenko, Multiagent Router Throttling: Decentralized Coordinated Response Against DDoS Attacks, In Proc. of AAAI 2013

- RL problem formalisation

- State space

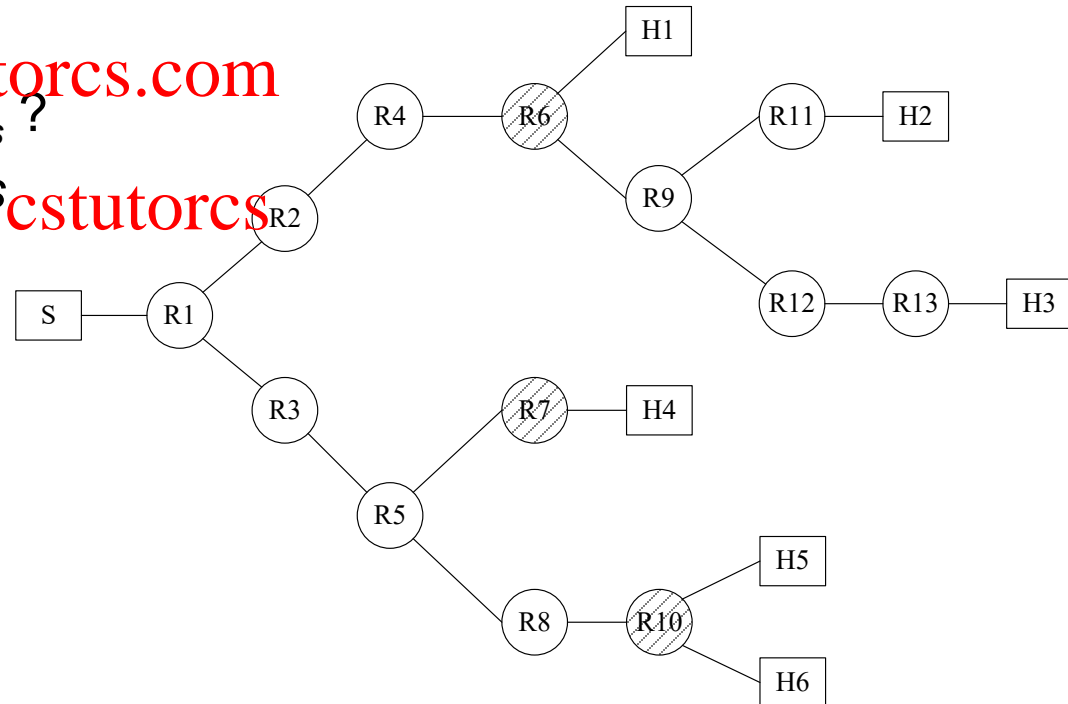
- Aggregated traffic arrived at the router over the last T seconds

- Action set

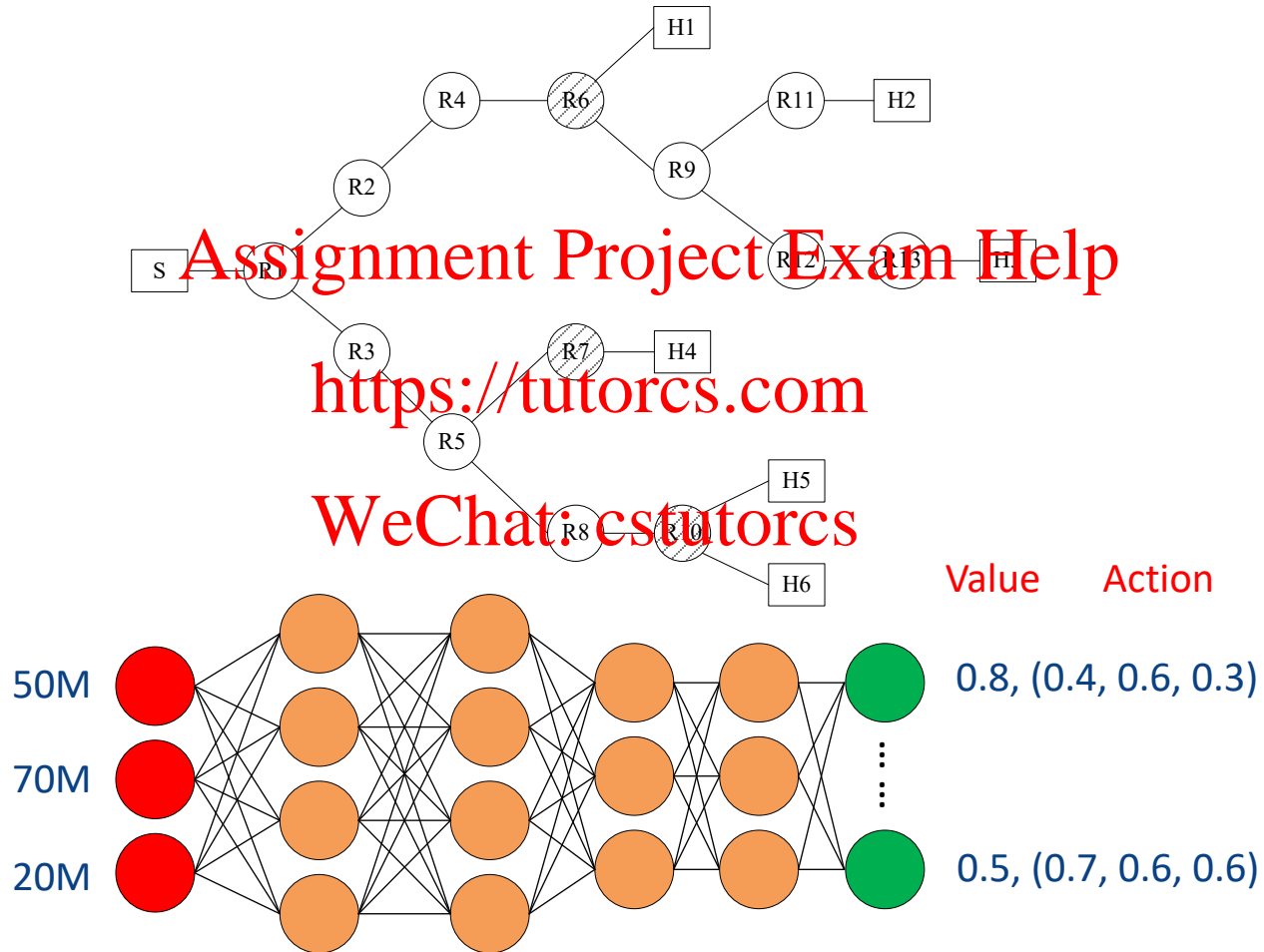
- Percentage of traffic to drop: 0, 10%, 20%, 30%, ... 90%

- Reward

- Aggregated traffic at $s > U_s$?
- Legitimate traffic reached s



- Training (DDQN)



$$L(\theta) = \mathbb{E} \left[\left(r + \gamma Q_2(s', \arg \max_{a'} Q_1(s', a')) - Q_1(s, a; \theta) \right)^2 \right]$$

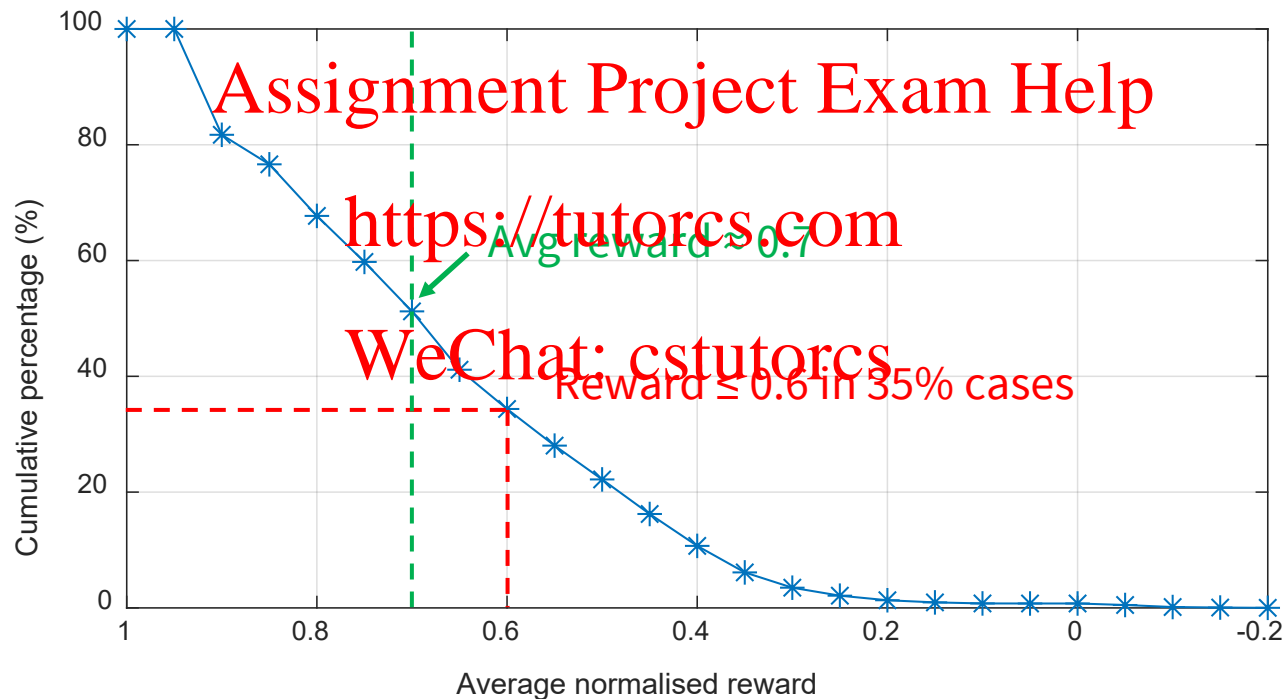
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



- Test
 - 10000 cases (may not be seen in training)



- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

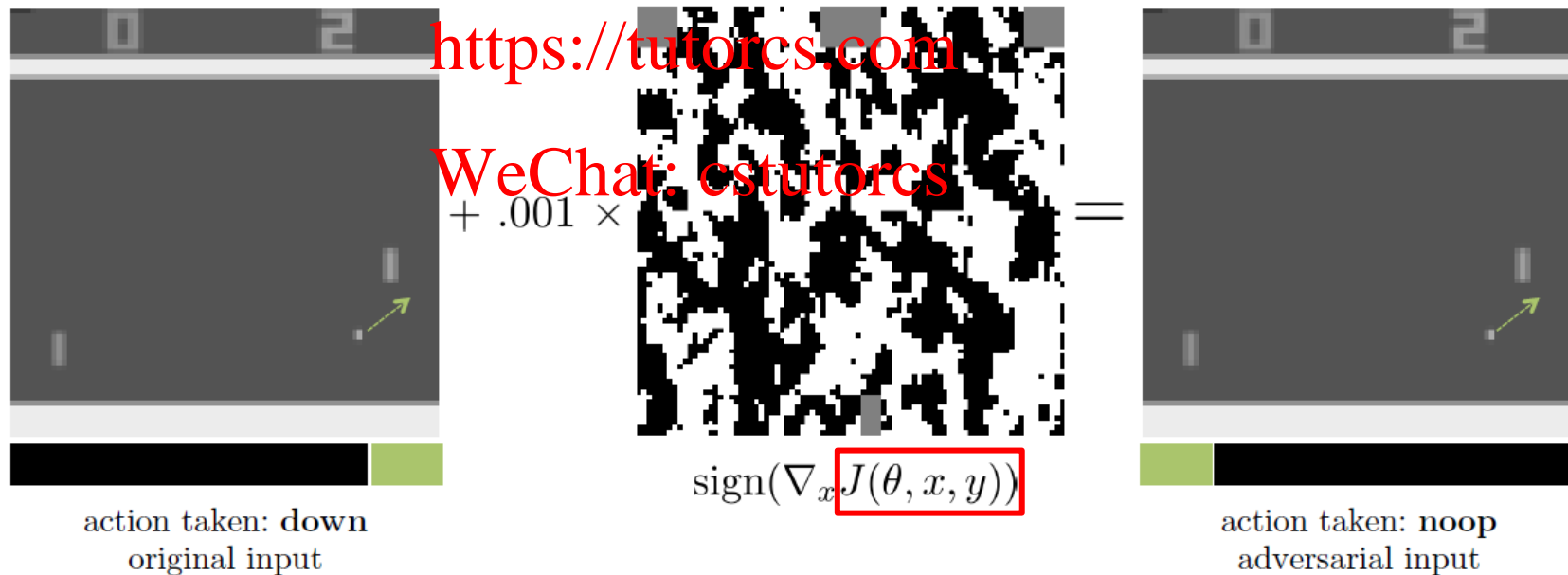
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

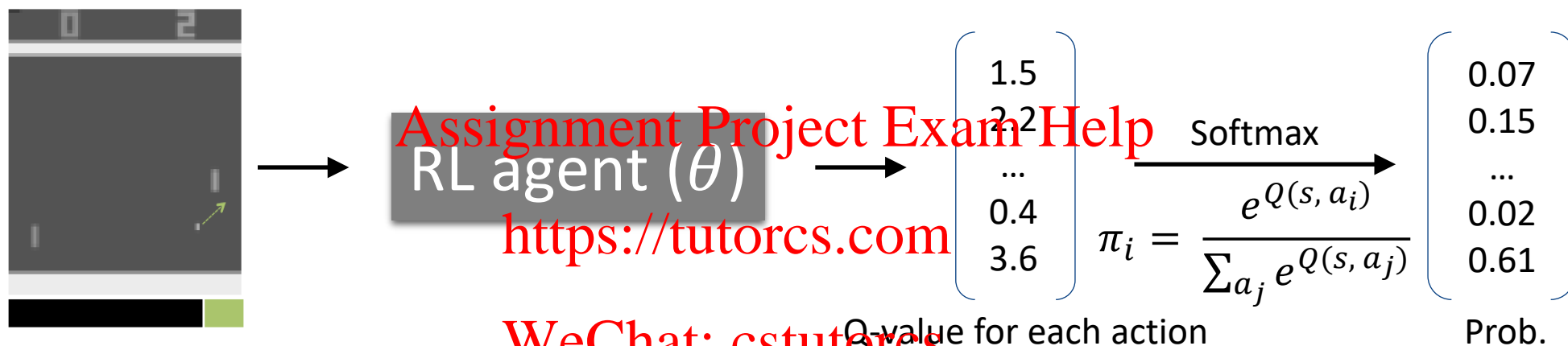
Adversarial attacks against RL models

- Test time attacks
 - Manipulate the environment observed by the agent [5]
 - Without attack: $\dots, s_t, a_t, r_t, s_{t+1}, a_{t+1}, r_{t+1}, s_{t+2}, \dots$
 - With attack: $\dots, s_t, a_t, r_t + \delta_t, s_{t+1}, a_{t+1}, r_{t+1} + \delta_{t+1}, s_{t+2} + \delta_{t+2}, \dots$

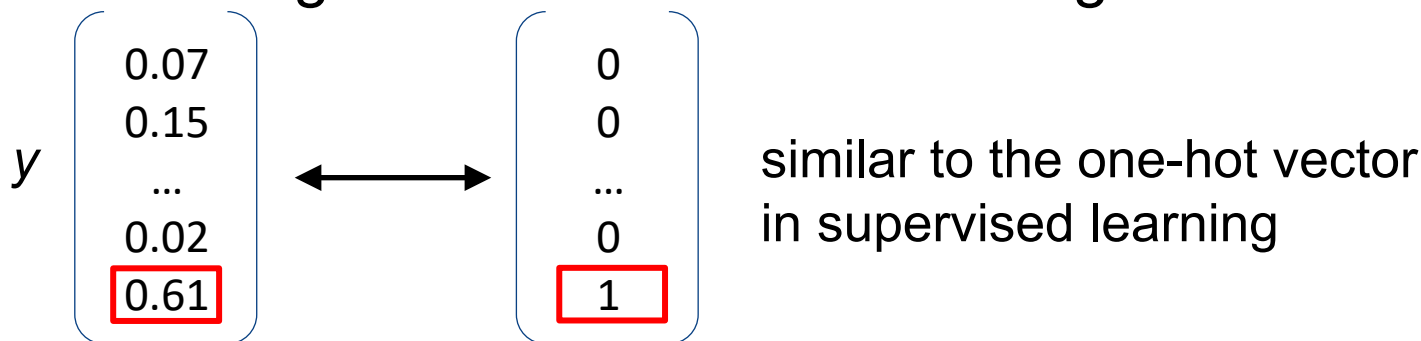


Adversarial attacks against RL models

- $J(\theta, x, y)$
 - y : softmax of the Q-value, i.e., prob. of taking an action



- J : cross-entropy loss between y and the distribution that places all weight on the action with the highest Q-value



Adversarial attacks against RL models

- Timing of the attack
 - Heuristic method [6]: launch the attack only when

$$c(s) = \max_a \frac{e^{\frac{Q(s,a)}{T}}}{\sum_{a_k} e^{\frac{Q(s,a_k)}{T}}} - \min_a \frac{e^{\frac{Q(s,a)}{T}}}{\sum_{a_k} e^{\frac{Q(s,a_k)}{T}}} > \beta$$

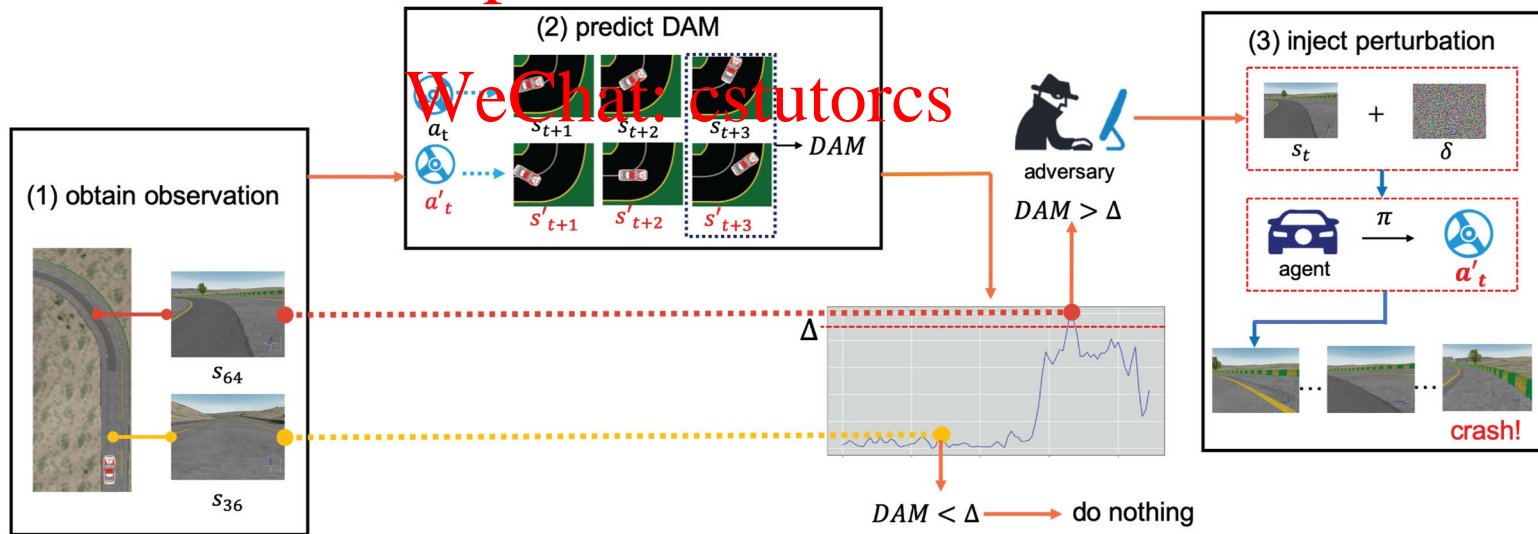
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



- Timing of the attack [8]
 - “Brute-force” search
 - Consider all possible N consecutive perturbations
 - Evaluate the attack damage at step $t + M$ ($M \geq N$)
 - $s_t, a_t, s_{t+1}, a_{t+1}, \dots, s_{t+N-1}, \dots, s_{t+M}$



- Timing of the attack [\[8\]](#)
 - “Brute-force” search
 - Train a prediction model: $(s_t, a_t) \rightarrow s_{t+1}$
 - Predict the subsequent states and actions, $\{(s_t, a_t), (s_{t+1}, a_{t+1}), \dots (s_{t+M}, a_{t+M})\}$
 - Assess the potential damage of all possible strategies
 - Danger Awareness Metric (DAM):
$$DAM = |T(s'_{t+M}) - T(s_{t+M})|$$
$$T$$
: domain-specific definition, e.g., distance between the car and the centre of the road

- Timing of the attack [8]
 - Train an antagonist model
 - Learn the optimal attack strategy automatically without any domain knowledge
 - Maintain a policy: $s_t \rightarrow (p_t, a'_t)$
 - If $p_t > 0.5$, add the perturbation to trigger a'_t
 - Take the original action a_t
 - Reward: negative of the agent's reward

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- Black-box attack [9]
 - Train a proxy model that learns a task that is *related* to the target agent's policy
 - S threat model
 - Only have access to the states
 - Approximate an expectation of the state transition function
 - $psychic(s_t, \theta_P) \approx \mathbb{E}_{\pi_T} [P(s_{t+1} | s_t)]$
 - SR threat model
 - Have access to the states and reward
 - Estimate the value V of a given state under the policy π_T
 - $assessor(s_t, \theta_A) \approx \mathbb{E}_{\pi_T} [\sum_{k=0}^{\infty} \gamma_t^{(k)} r_{t+k+1}] = V^{\pi_T}(s_t)$

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- Black-box attack [9]
 - SA threat model
 - Have access to states and actions
 - Approximate the target's policy π_T
 $imitator(s_t, \theta_I) \approx \pi_T(s_t)$
 - SRA threat model
 - Have access to states, actions and rewards
 - Action-conditioned psychic (AC-psychic):
 $AC-psychic(s_t, \theta_P) \approx \mathbb{E}_{\pi_T}[P(s_{t+1}|s_t, a_t)]$
 - Combine *assessor* and *AC-psychic* to decide whether to perturb the state

- Black-box attack [9]
 - SRA threat model

Algorithm 1: Strategically-timed snooping attack

Input: Trained assessor, trained AC-psychic, trained proxy \mathcal{M}_κ , trained target agent \mathcal{T} , β

for $t = 1, T$ **do**

 Initialize empty list q ;

foreach $a \in \mathcal{A}$ **do**

 Predict s_{t+1}^H with AC-psychic(s_t, a)

 Estimate V^H with assessor(s_{t+1}^H)

 Append V^H to q ;

end

$c(s_t) = \max [\text{Softmax}(q)] - \min [\text{Softmax}(q)]$

if $c(s_t) \geq \beta$ **then**

 | Perturb s_t using $\nabla_x J_{\mathcal{M}_\kappa}$

end

 Feed s_t to target \mathcal{T} for action decision;

end

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$\mathcal{K} \in \{S, SR, SA\}$

- Black-box attack [9]
 - Surrogate: assume the adversary has access to the target agent's environment and can train an identical model

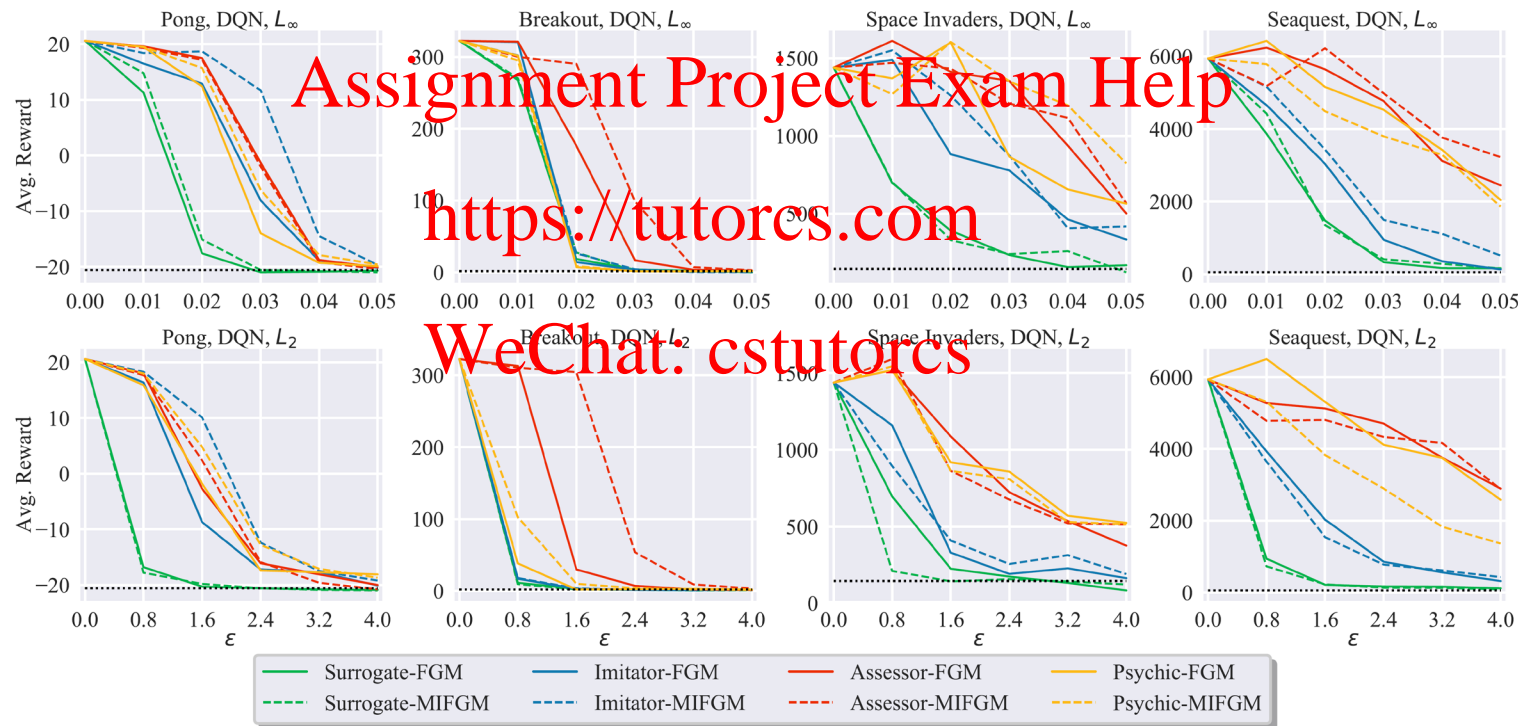


Figure 5: Performance reduction of DQN agents due to L_∞ and L_2 bounded perturbations. The black dotted line represents a random-guess policy.

- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- Training time attack

- Without attack: $\dots, (s_t, a_t, s_{t+1}, r_t), (s_{t+1}, a_{t+1}, s_{t+2}, r_{t+1}), \dots$
- With attack: $\dots, (s_t, a_t, s_{t+1} + \delta_{t+1}, r'_t), (s_{t+1} + \delta_{t+1}, a'_{t+1}, s_{t+2} + \delta_{t+2}, r'_{t+1}), \dots$
- Purpose: generate δ_{t+1} so that the agent will not take the next action a_{t+1}
- Cross entropy loss: $J = -\sum_i p_i \log \pi_i$

- $\pi_i = \frac{e^{Q(s, a_i)}}{\sum_{a_j} e^{Q(s, a_j)}} \rightarrow \text{prob. of taking action } a_i$

- $p_i = \begin{cases} 1, & \text{if } a_i = a_{t+1} \\ 0, & \text{otherwise} \end{cases}$

- Maximise $J = -\log \pi_{t+1} \rightarrow$ minimise the prob. of taking a_{t+1}

- $\delta = \alpha \cdot \text{Clip}_\epsilon \left(\frac{\partial J}{\partial s} \right)$

- Background on reinforcement learning
 - Introduction
 - Q-learning
 - Application in defending against DDoS attacks
- Adversarial attacks against RL models
 - Test time attack
 - Training time attack
- Defence

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- Adversarial training [7]

- Calculate δ using the attacker's strategy: $(s_t, a_t, s_{t+1} + \delta_{t+1}, r'_t)$
- $a'_{t+1} = \arg \max_a Q(s_{t+1} + \delta_{t+1}, a)$
- Generate experience $(s_{t+1}, a'_{t+1}, s'_{t+2}, r'_{t+1})$ for the agent to train on

Untampered state

Potentially non-optimal action

<https://tutorcs.com>
→ explore more

WeChat: cstutorcs

- Reinforcement learning
 - State, action, reward
 - Value function, policy, model
 - Q-learning → Q-network → DQN → DDQN
- Adversarial reinforcement learning
 - Test time attack
 - Timing of the attack
 - Black-box attack
 - Training time attack
- Defence – adversarial training

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- [1] R. S. Sutton and A. G. Barto, *Introduction to Reinforcement Learning*, First. Cambridge, MA, USA: MIT Press, 1998.
- [2] V. Mnih *et al.*, “Playing Atari with Deep Reinforcement Learning,” *CoRR*, vol. abs/1312.5602, 2013.
- [3] H. V. Hasselt, A. Guez, and D. Silver, “Deep Reinforcement Learning with Double Q-learning,” *eprint arXiv:1509.06461*, Sep. 2015.
- [4] K. Malialis and D. Kudenko, “Multiagent Router Throttling: Decentralized Coordinated Response Against DDoS Attacks,” in *Proc. of the 27th AAAI Conference on Artificial Intelligence*, Washington, 2013, pp. 1551–1556.
- [5] **S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, “Adversarial Attacks on Neural Network Policies,” eprint arXiv:1702.02284, 2017.**
- [6] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun, “Tactics of Adversarial Attack on Deep Reinforcement Learning Agents,” *eprint arXiv:1703.06748*, Mar. 2017.
- [7] **A. Pattanaik, Z. Tang, S. Liu, G. Bommannan, and G. Chowdhary, “Robust Deep Reinforcement Learning with Adversarial Attacks,” arXiv:1712.03632 [cs], Dec. 2017.**

- [8] Jianwen Sun and Tianwei Zhang and Xiaofei Xie and Lei Ma and Yan Zheng and Kangjie Chen and Yang Liu, “Stealthy and Efficient Adversarial Attacks against Deep Reinforcement Learning,” AAAI 2020: 5883-5891
- [9] Matthew Inkawhich, Yiran Chen, and Hai Li. 2020. Snooping Attacks on Deep Reinforcement Learning. In Proceedings of the 19th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS '20). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 557–565.

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



Adversarial Reinforcement Learning in Autonomous Cyber Defence

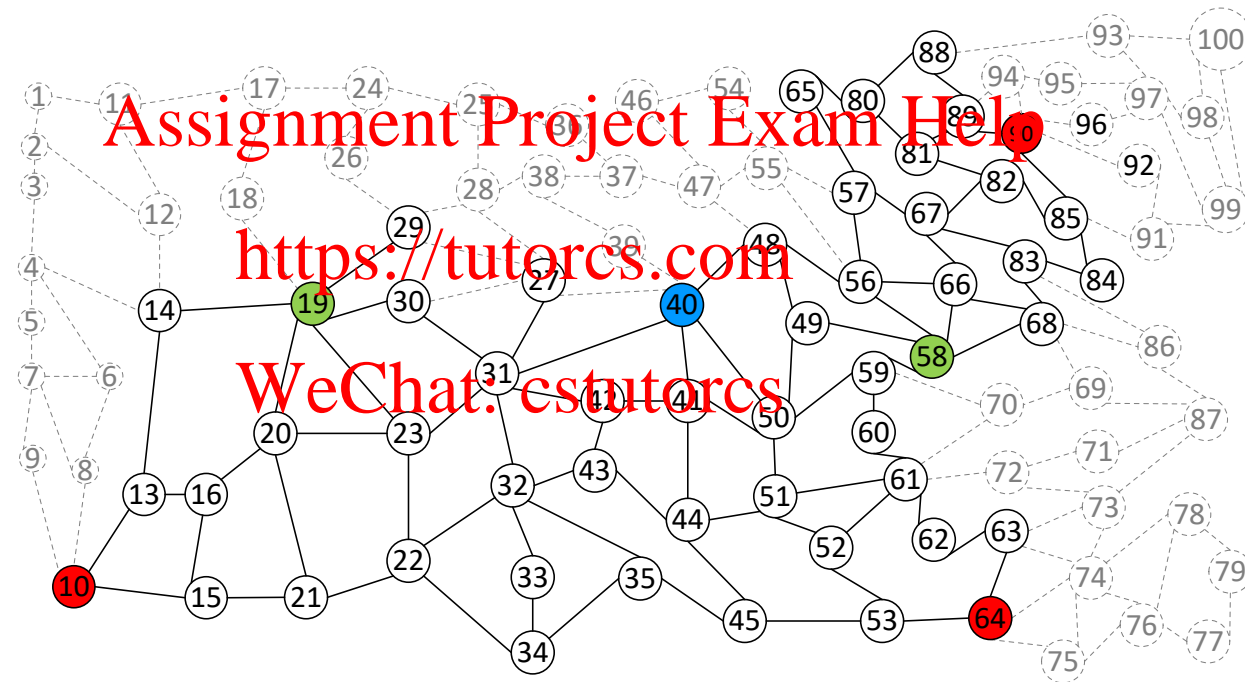
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Adversarial attacks against RL models

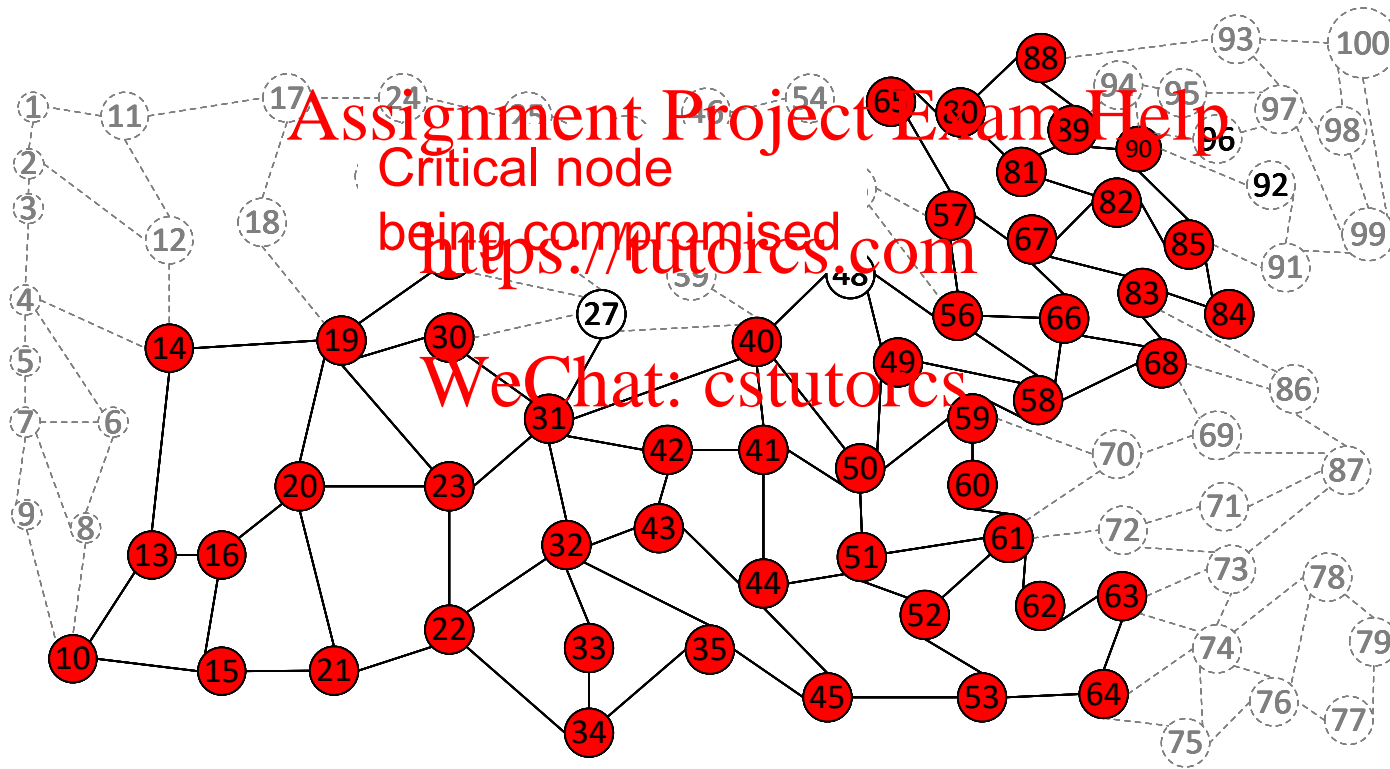
- Attacker: propagates through the network to compromise the critical server
- The defender applies RL to prevent the critical server from compromise, and preserve as many nodes as possible



- Initially compromised nodes
- Critical node
- Possible migration destination
- --- Nodes, links only visible to the defender
- — Nodes, links visible to the defender & the attacker

Adversarial attacks against RL models

- Attacker: partial observability of the network topology



Adversarial attacks against RL models

- Problem definition
 - State: $[0, 0, \dots, 0, 0, 0, \dots, 0]$

State of each link, 0: on, 1: off
 - State of each node, 0: uncompromised, 1: compromised
 - Action:
 - Action $0 \sim N-1$: isolate & patch a node $i \in [0, N-1]$
 - Action $N \sim 2N-1$: reconnect a node $i \in [0, N-1]$
 - Action $2N \sim 2N+M-1$: migrate the critical node to one of the M destinations
 - Action $2N+M$: take no action
 - Reward:
 - -1: (1) critical node is compromised or isolated, (2) invalid action
 - Proportional to number of uncompromised nodes that can still reach the critical node
 - Attacker can only compromised a node x if there is a visible link between x and any compromised node

- 82/100 nodes are preserved



- Training-time attack: manipulate states to prevent agents from taking optimal actions
 - $(s_t, a_t, s_{t+1}, r_t) \rightarrow (s_t, a_t, s_{t+1} + \delta_{t+1}, r'_t)$
 - Binary state \rightarrow cannot use gradient-descent based method
 - δ : false positives & false negatives
 - The attacker cannot manipulate the states of all the observable nodes
 - L_{FP} : nodes that can be perturbed as false positive
 - L_{FN} : nodes that can be perturbed as false negative
 - $\min Q(s_{t+1} + \delta_{t+1}, a_{t+1})$, a_{t+1} : the optimal action for s_{t+1} that has been learned so far
 - Loop through L_{FP} (L_{FN}) and flip the state of one node per time
 - Rank all nodes based on ΔQ (decrease of Q-value by flipping state)
 - Flip the states of the top K nodes

Adversarial attacks against RL models

Algorithm 1: Causative attack against DDQN via state perturbation

Input : The original experience, (s, a, s', r) ;
The list of observable nodes, N_O ;
The list of nodes that can be perturbed as false positive (false negative) by the attacker, L_{FP} (L_{FN});
The main DQN, Q ;
Limit on the number of FPs and FNs per time, $LIMIT$

Output: The tampered experience $(s, a, s' + \delta, r')$

```

1   $FN = FP = \{\}$ ;
2   $minQ_{FN} = minQ_{FP} = \{\}$ ;
3   $a' = \operatorname{argmax}_{a^*} Q(s', a^*)$ ;
4  for node  $n$  in  $N_O$  do
5      if  $n$  is compromised and  $n$  in  $L_{FN}$  then
6          mark  $n$  as uncompromised;
7          if  $Q(s' + \delta, a') < \text{any value in } minQ_{FN}$  then
8              //  $\delta$  represents the FP and/or FN readings
              insert  $n$  and  $Q(s' + \delta, a')$  into appropriate
              positions in  $FN$  and  $minQ_{FN}$ ;
9          if  $|FN| > LIMIT$  then
10             remove extra nodes from  $FN$  and
              $minQ_{FN}$ ;
11         restore  $n$  as compromised;
12     else if  $n$  is uncompromised and  $n$  in  $L_{FP}$  then
13         mark  $n$  as compromised;
14         if  $Q(s' + \delta, a') < \text{any value in } minQ_{FP}$  then
15             insert  $n$  and  $Q(s' + \delta, a')$  into appropriate
             positions in  $FP$  and  $minQ_{FP}$ ;
16         if  $|FP| > LIMIT$  then
17             remove extra nodes from  $FP$  and
              $minQ_{FP}$ ;
18     restore  $n$  as uncompromised;
19 Change nodes in  $FN$  to uncompromised;
20 Change nodes in  $FP$  to compromised;
21 return  $(s, a, s' + \delta, r')$ 

```

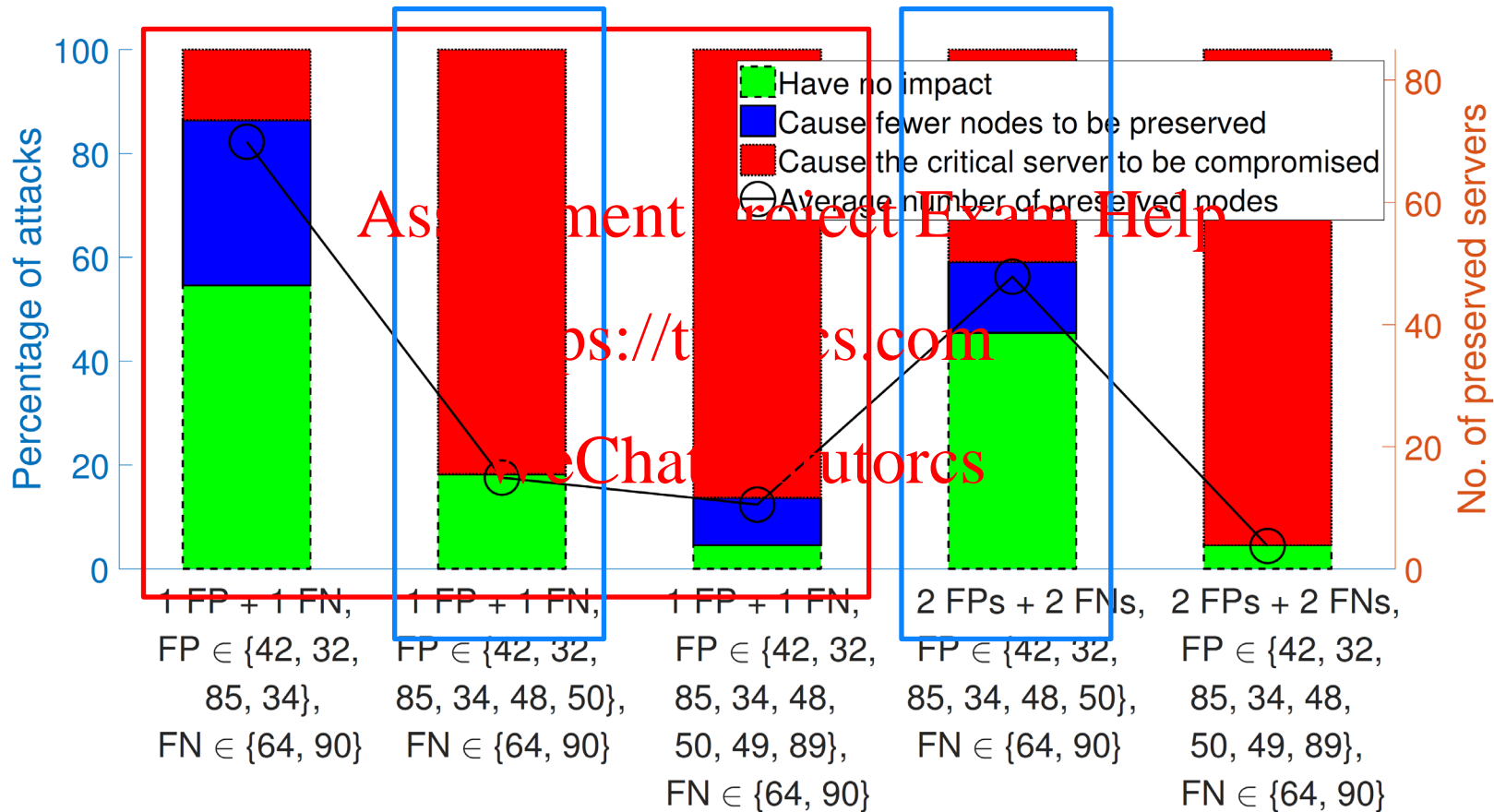
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

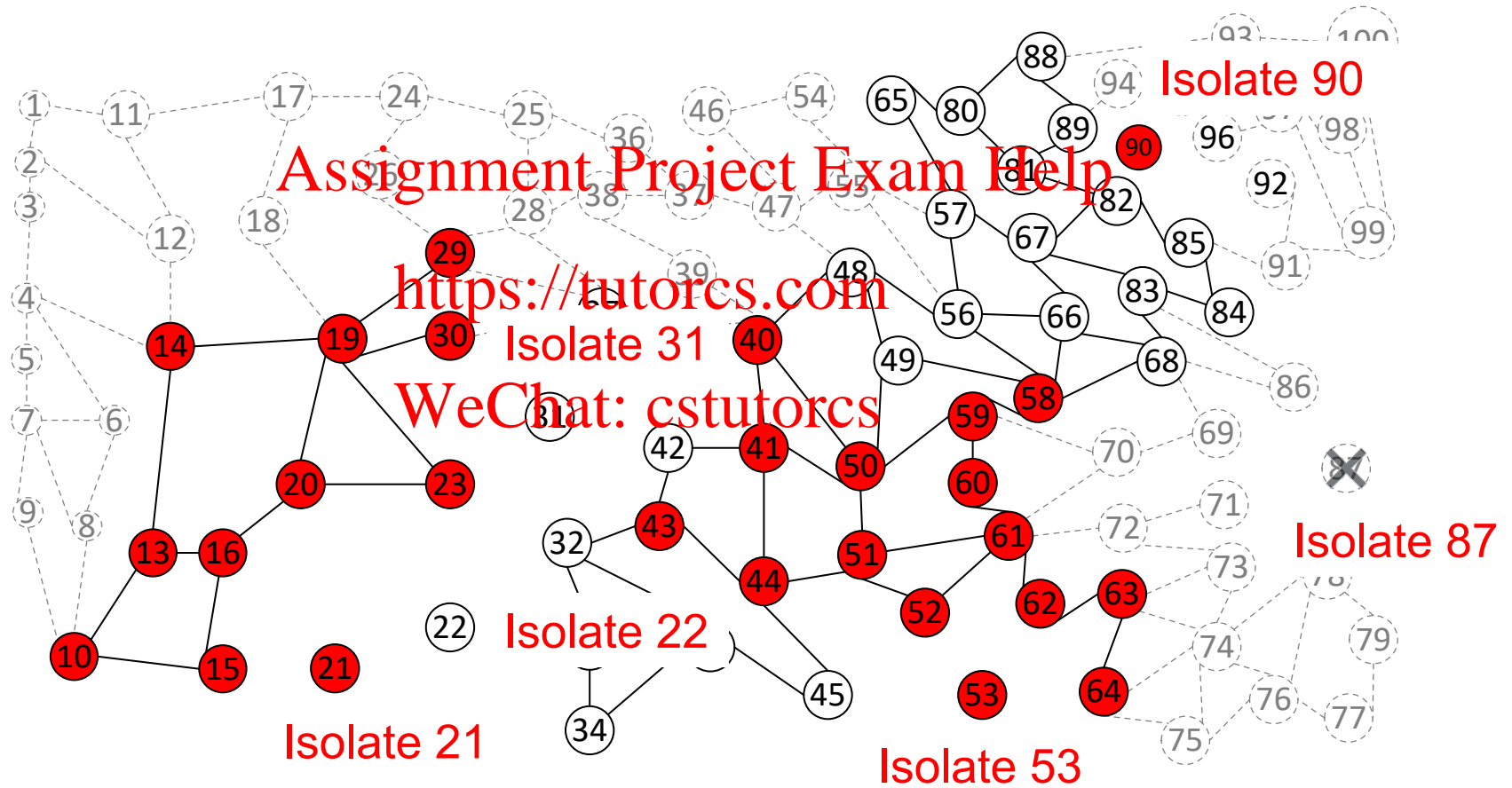
Adversarial attacks against RL models

- Result



Adversarial attacks against RL models

- After training-time attacks



Inversion defence method

- Aim to revert the perturbation/false readings

Attacker	Defender
$s_{t+1} \rightarrow s_{t+1} + \delta_{t+1}$ minimise $Q(s_{t+1} + \delta_{t+1}, a_{t+1})$	$s_{t+1} + \delta_{t+1} \rightarrow s_{t+1} + \delta_{t+1} + \delta'_{t+1}$ maximise $Q(s_{t+1} + \delta_{t+1} + \delta'_{t+1}, a_{t+1})$
Loop through L_{FP} and L_{FN}	Loop through all nodes
Flip K nodes	Flip K' nodes

- Effective even if $K' \neq K$
- Minimum impact on normal training process (i.e., $K = 0$, $K' > 0$)

Inversion defence method

- Before & after the defence method is applied

