



程序代写代做 CS编程辅导



We Adversarial Reinforcement Learning

WeChat: cstutorcs

Assignment Project Exam Help

COMP90073
Email: tutorcs@163.com
Security Analytics

QQ: 749389476
Yi Han, CIS

<https://tutorcs.com>
Semester 2, 2021

程序代写代做 CS编程辅导

- 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



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Background on reinforcement learning

- Introduction
 - Q-learning
 - Application in defending against DDoS attacks



- Adversarial attacks against RL models

- Test time attack
 - Training time attack

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

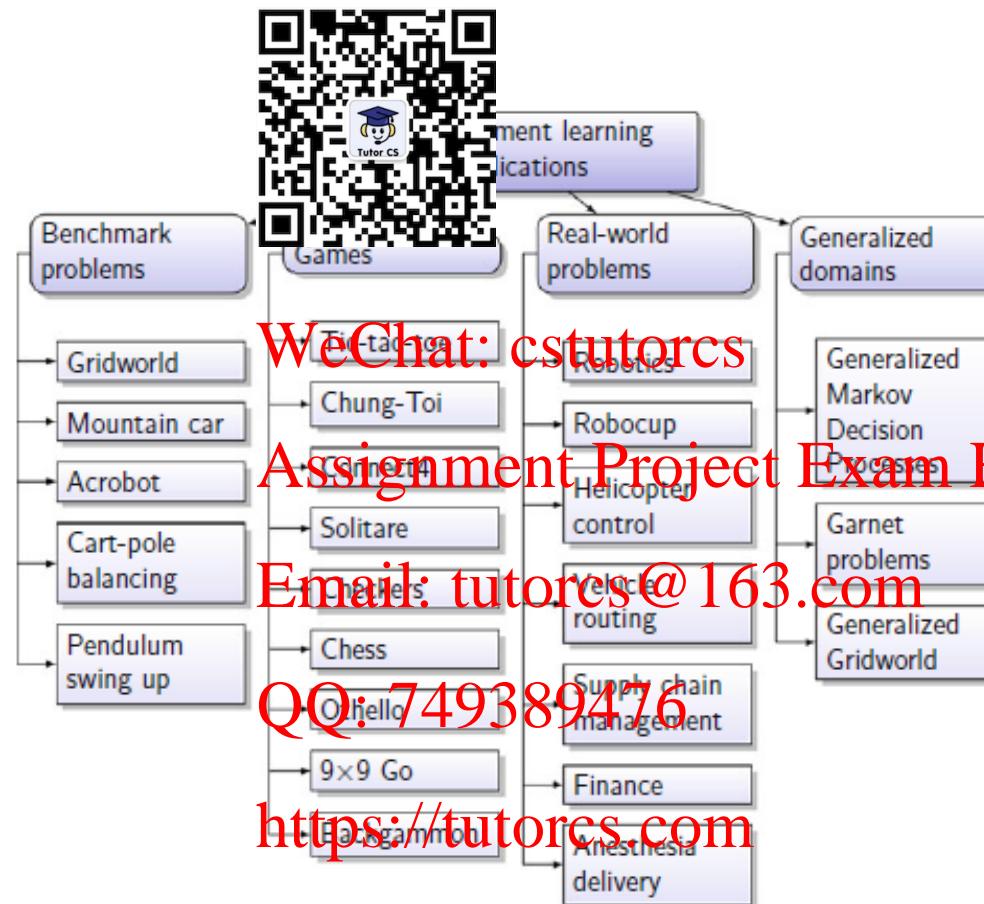
QQ: 749389476

<https://tutorcs.com>

Background on reinforcement learning

- Application

程序代写代做 CS编程辅导



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutores@163.com

QQ: 749389476

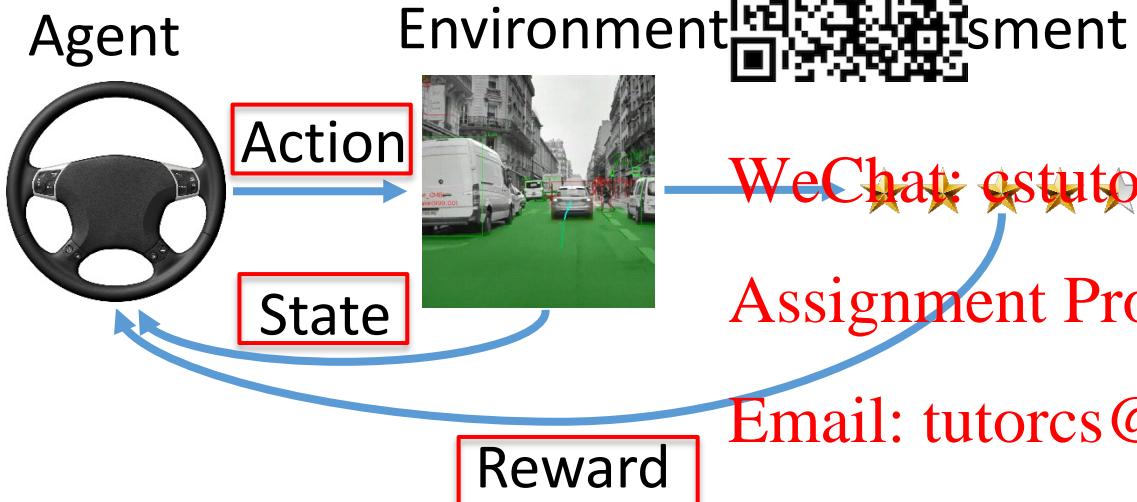
<https://tutores.com>

<https://www.youtube.com/watch?v=9DgXsautpilot-self-driving-car>

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

Background on reinforcement learning

- Introduction



程序代写代做 CS编程辅导



WeChat: cstutorcs

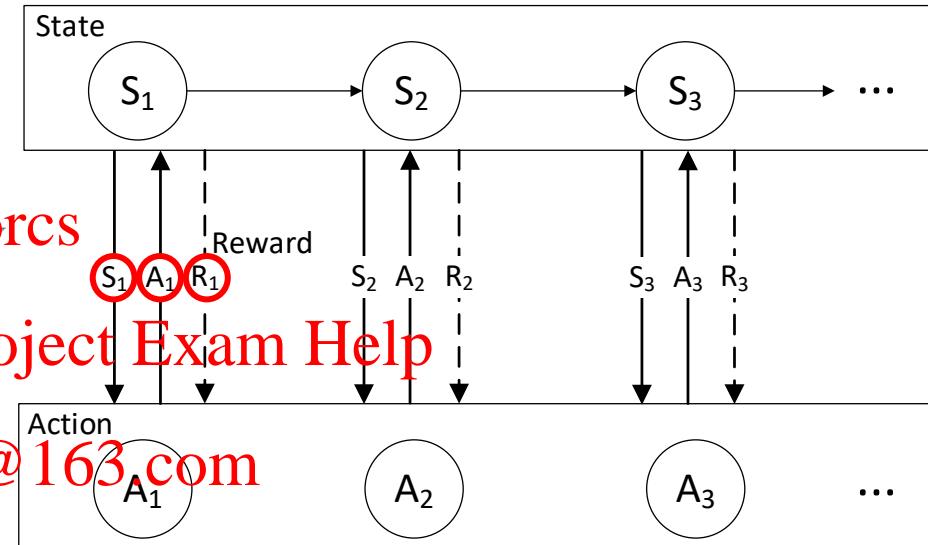
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

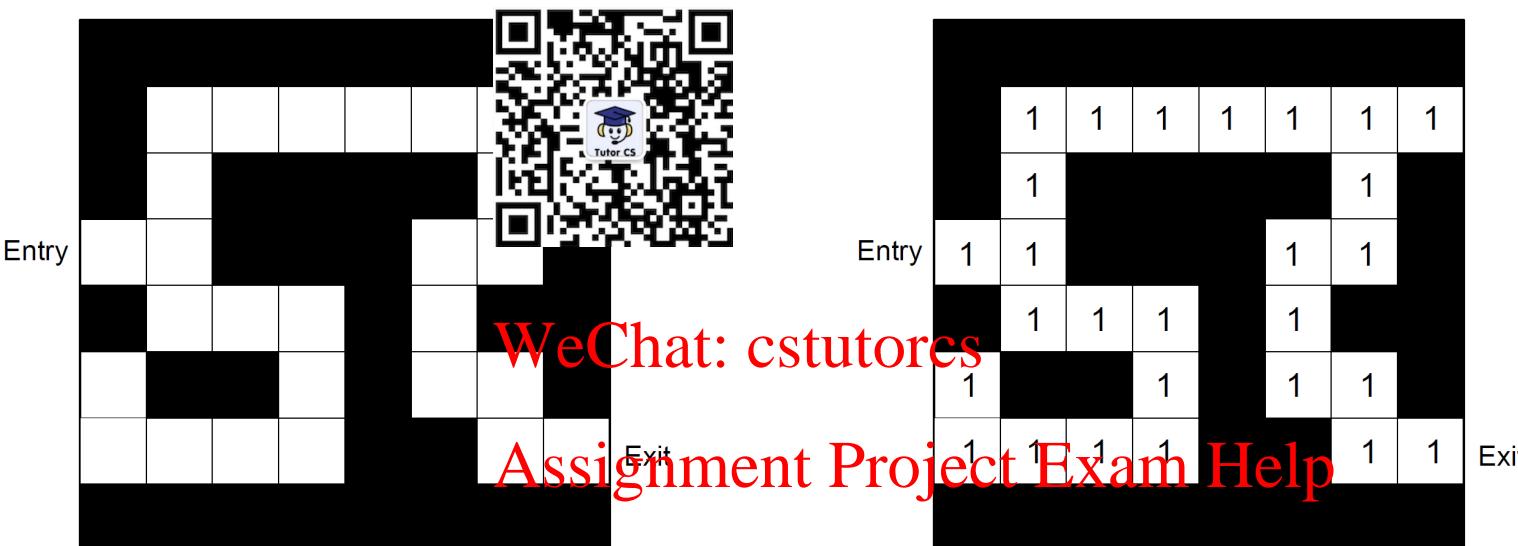
Maximise the discounted cumulative rewards
over the long run: $R_t = \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{t,\tau}$

<https://tutorcs.com>



程序代写代做 CS编程辅导

- State



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

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

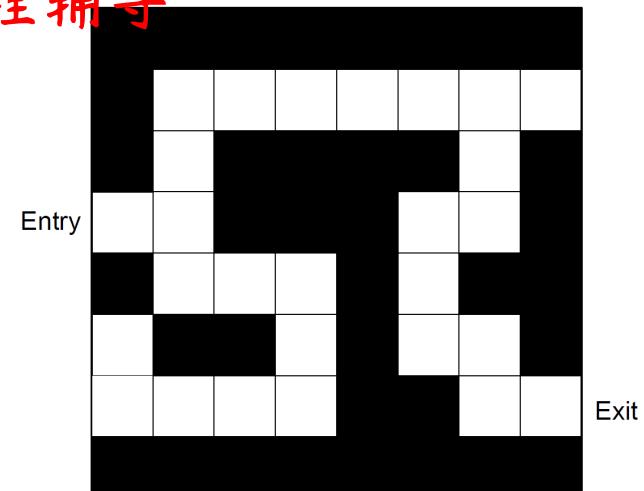
0	2	1	1	1	1	1	1
0	2	0	0	0	0	1	0
2	2	0	0	0	1	1	0
0	1	1	1	0	1	0	0
1	0	0	1	0	1	1	0
1	1	1	1	0	0	1	1

- Action
 - Up
 - Left
 - Down
 - Right

程序代写代做 CS编程辅导



WeChat: cstutorcs



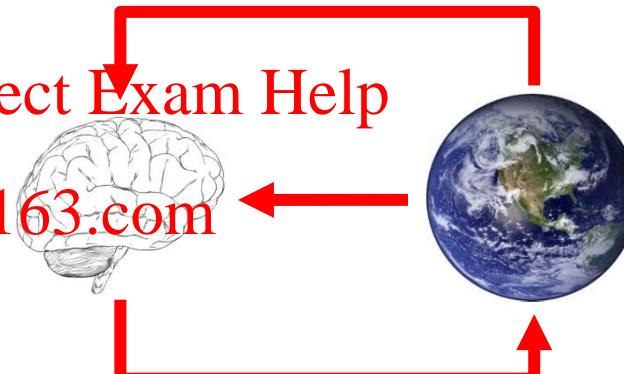
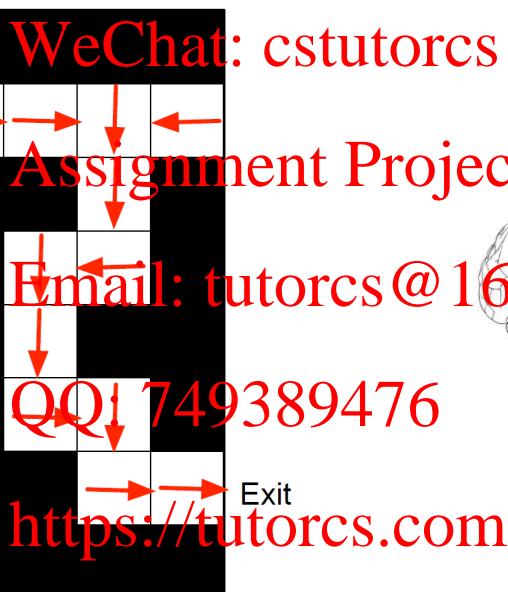
- Reward: an immediate feedback on whether an action is good
 - In the range of $[-1, 1]$
 - 1: reach the exit
 - -0.8: move to a blocked cell
 - -0.3: move to a visited cell
 - -0.05: move to an adjacent cell

Email: tutorcs@163.com

QQ: 749389476

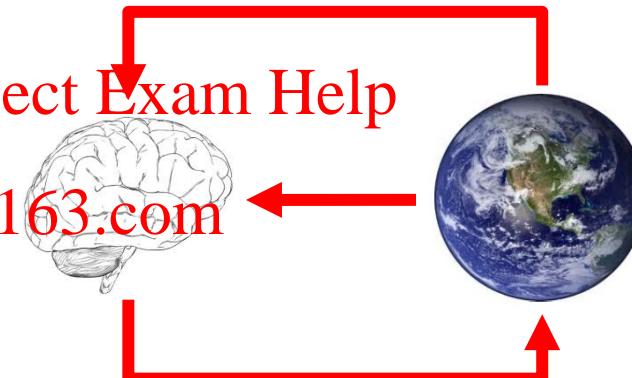
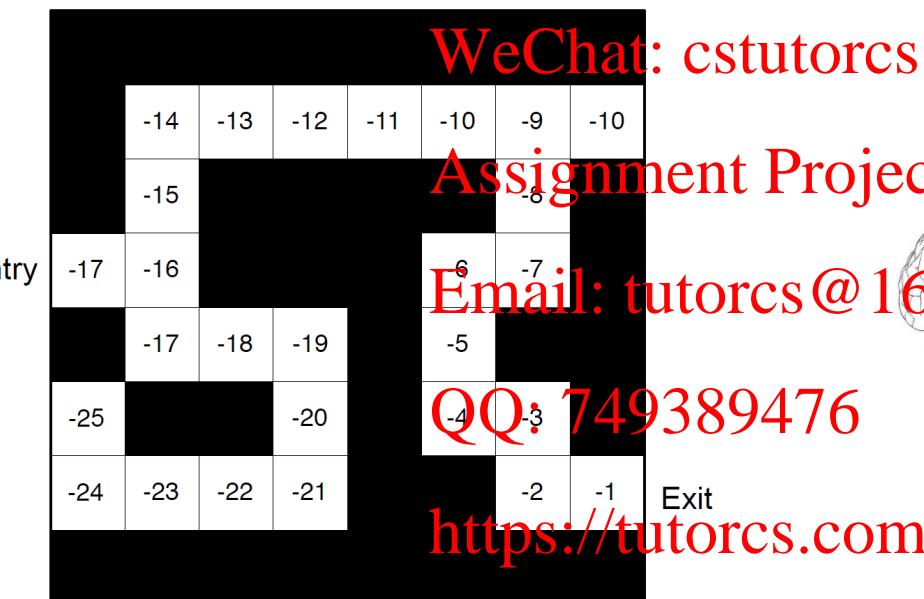
<https://tutorcs.com>

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



程序代写代做 CS编程辅导

- 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 so far observed actions
 - Expected value of following policy π starting from state s



程序代写代做 CS编程辅导

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



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

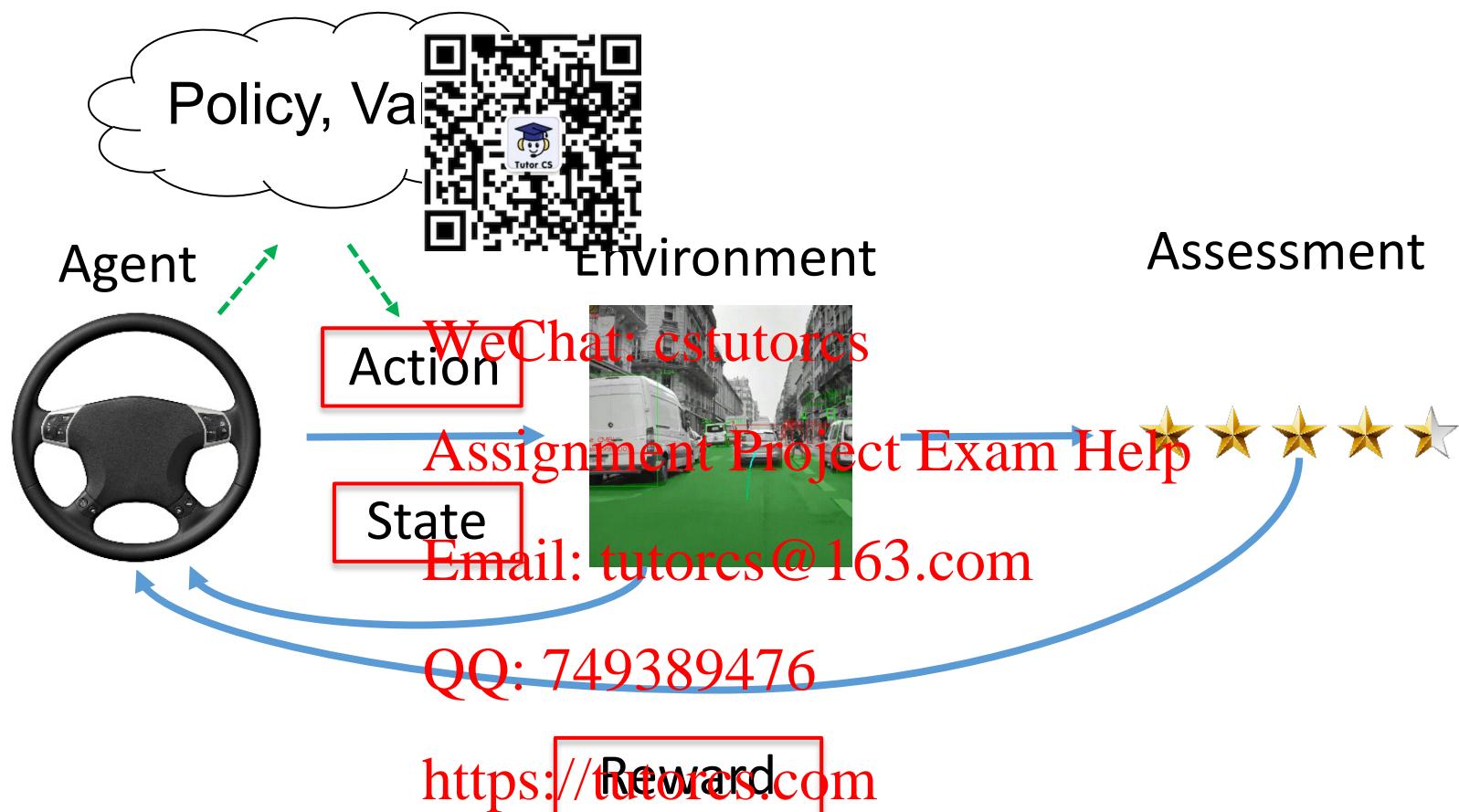
<https://tutorcs.com>



A red rectangular frame surrounds the text "Assignment Project Exam Help" and the globe icon. Red arrows point from the text "Assignment Project Exam Help" to the brain icon and from the globe icon to the text "Assignment Project Exam Help".

Background on reinforcement learning

程序代写代做 CS编程辅导

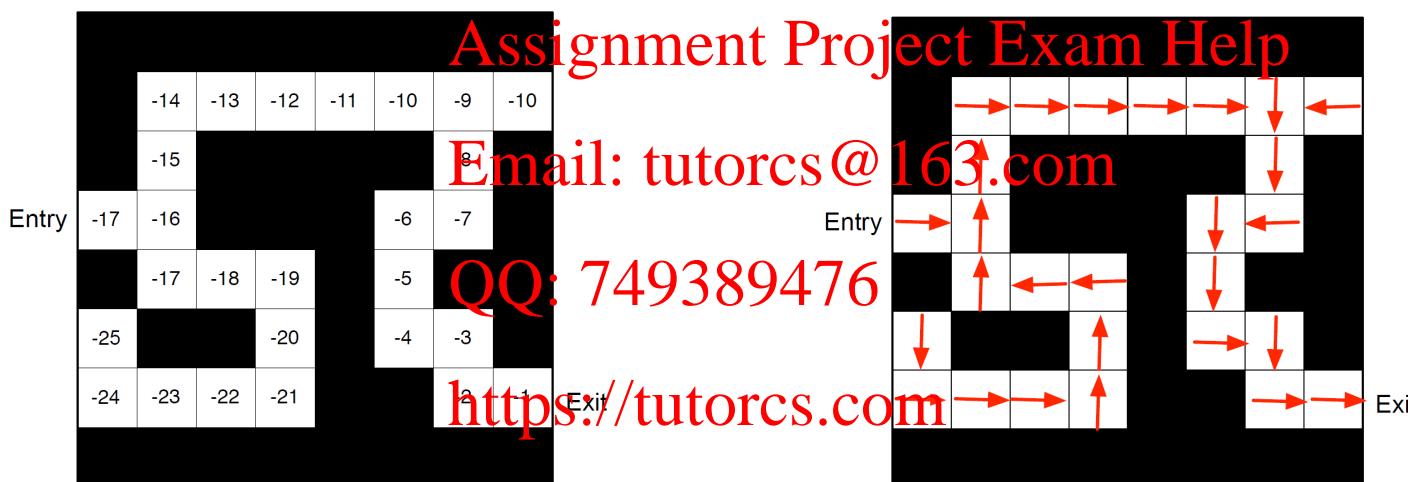


程序代写代做 CS编程辅导

- 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

WeChat: cstutorcs



程序代写代做 CS编程辅导

- Classification

- Model free algorithm directly learns the policy and/or the value function
- Model based algorithm first builds up how the environment works

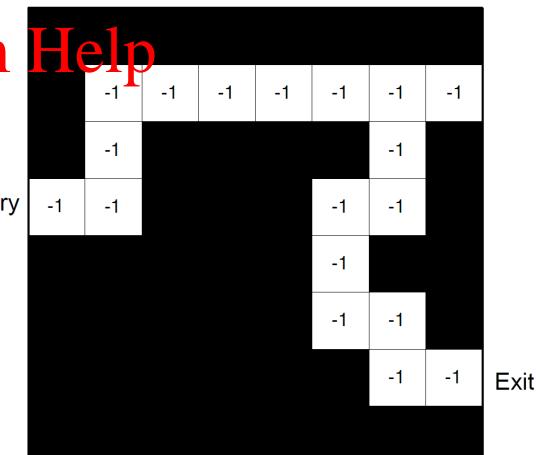
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

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



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

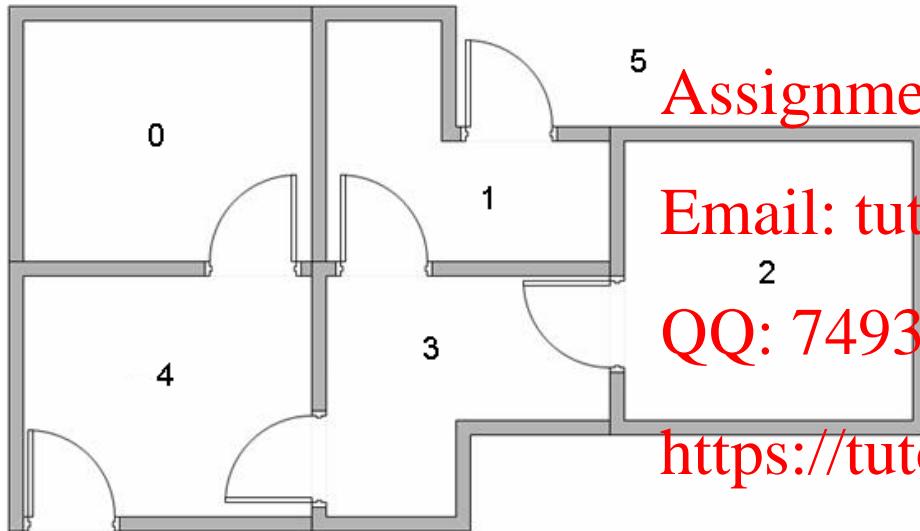
程序代写代做 CS编程辅导

- Q-learning: estimate action-value function $Q(s, a)$
 - Expected value of action a in state s and then following policy π :



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

WeChat: cstutorcs

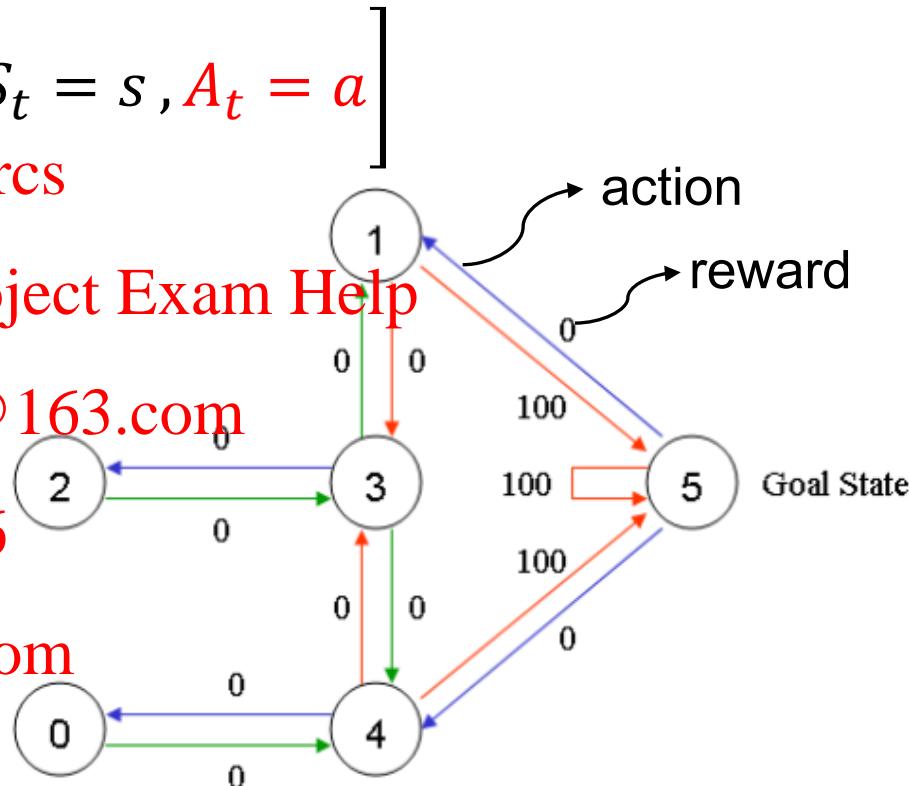


Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



<http://mnemstudio.org/path-finding-q-learning-tutorial.htm>

- Q-learning

程序代写代做 CS编程辅导

$$R = \begin{array}{c} \begin{matrix} & \text{Action} \\ \text{State} & 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \end{matrix} \\ \begin{matrix} 0 & \begin{bmatrix} -1 & -1 & -1 & -1 & 0 \end{bmatrix} \\ 1 & \begin{bmatrix} -1 & -1 & -1 & 0 & -1 \end{bmatrix} \\ 2 & \begin{bmatrix} -1 & -1 & -1 & 0 & -1 \end{bmatrix} \\ 3 & \begin{bmatrix} -1 & 0 & 0 & -1 & 0 \end{bmatrix} \\ 4 & \begin{bmatrix} 0 & -1 & -1 & 0 & -1 \end{bmatrix} \\ 5 & \begin{bmatrix} -1 & 0 & -1 & -1 & 0 \end{bmatrix} \end{matrix} \end{array}$$


 WeChat: cstutorcs
 Assignment Project Exam Help
 Email: tutorcs@163.com

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 1 & 0 & 0 & 0 & 320 & 100 + 0.8 * \max(0, 0, 0) \\ 2 & 0 & 0 & 0 & 320 & 0 \\ 3 & 0 & 400 & 256 & 0 & 400 + 0.8 * \max(0, 100) \\ 4 & 320 & 0 & 0 & 320 & 0 \\ 5 & 0 & 400 & 0 & 0 & 400 \\ & & & & & 500 \end{bmatrix}$$

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

QQ: 749389476

estimate of optimal future value
<https://tutorcs.com>

Background on reinforcement learning

- Exploitation vs. Exploration

ϵ -greedy



程序代写代做 CS编程辅导

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

WeChat: cstutorcs
Assignment Project Exam Help

Target Q
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

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

程序代写代做 CS编程辅导

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



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Background on reinforcement learning

程序代写代做 CS编程辅导

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



- Can RL be applied to throttle flooding DDoS attacks?

- Problem setup [5]

- A mixed set of legitimate & attackers
- Aggregated traffic at
- RL agents decides the drop rates
- No anomaly detection

程序代写代做 CS编程辅导



WeChat: cstutorcs
expensive

Assignment Project Exam Help

R: router
H: host
Email: tutorcs@163.com

QQ: 749389476

Server to protect
<https://tutorcs.com>

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



程序代写代做 CS编程辅导

- RL problem formalisation

- State space

- Aggregated traffic at the router over the last T seconds



- Action set

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

- Reward

WeChat: cstutorcs

- Aggregated traffic at $s > U_s$?

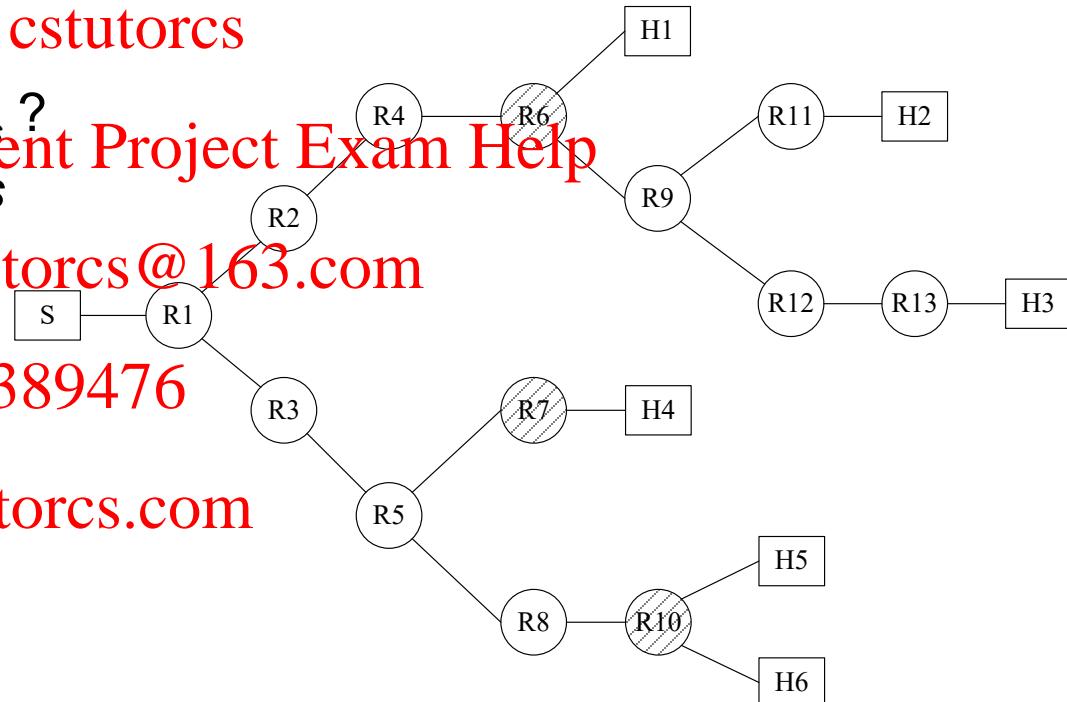
Assignment Project Exam Help

- Legitimate traffic reached s

Email: tutorcs@163.com

QQ: 749389476

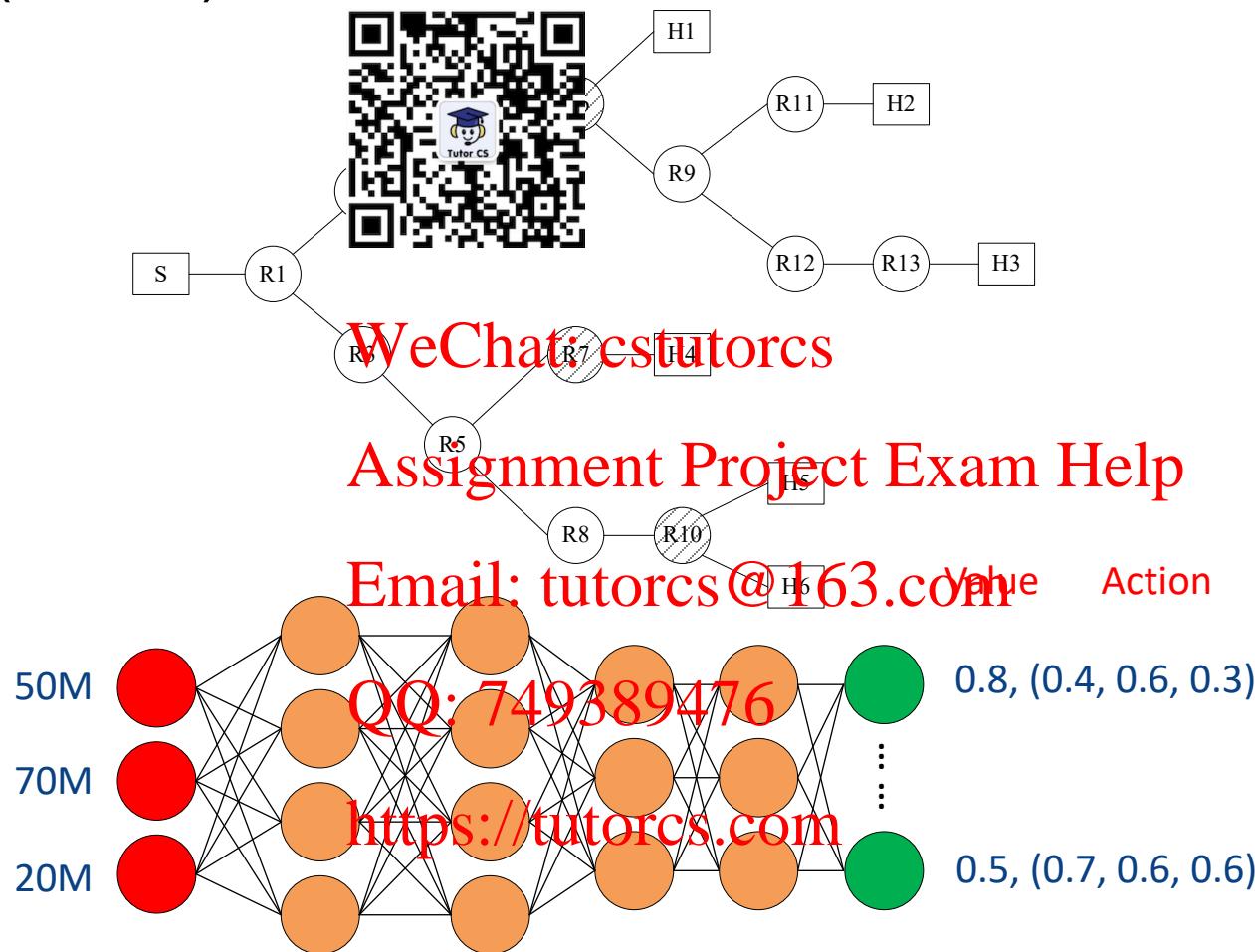
<https://tutorcs.com>



Background on reinforcement learning

- Training (DDQN)

程序代写代做 CS编程辅导



程序代写代做 CS编程辅导

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



WeChat: cstutorcs

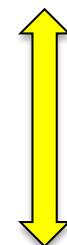


Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

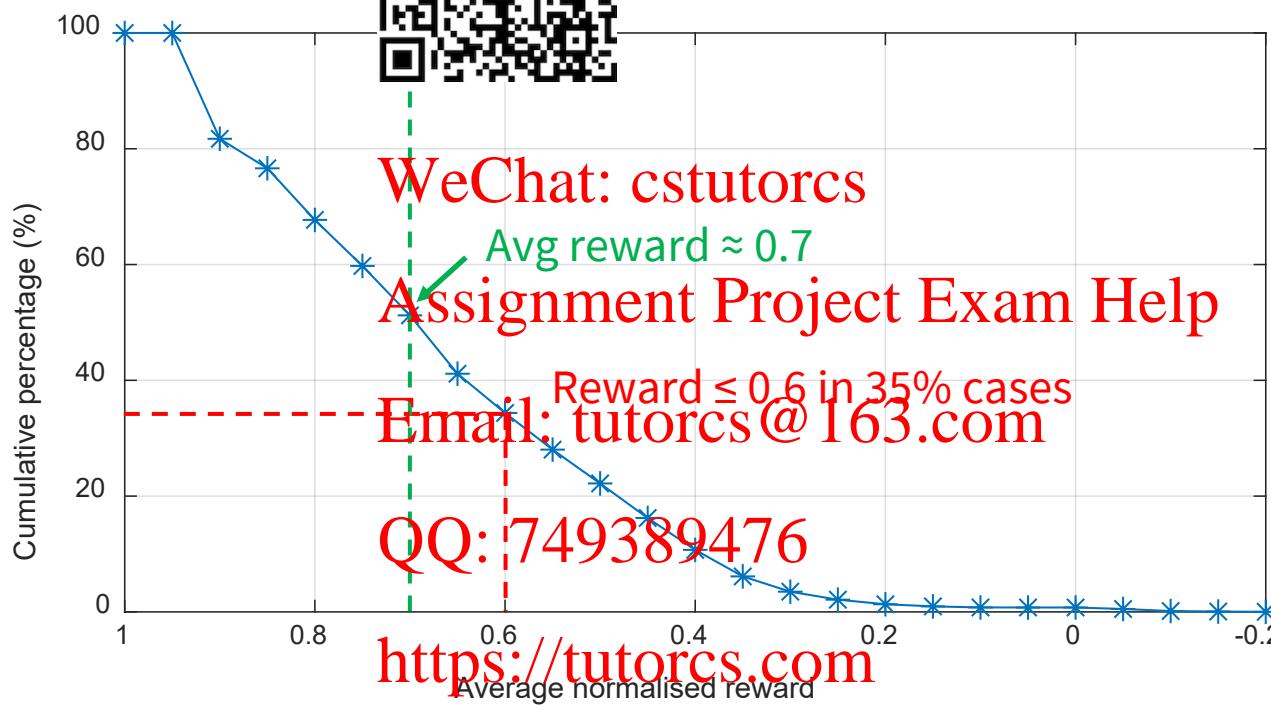


Background on reinforcement learning

程序代写代做 CS编程辅导

- Test

- 10000 cases (many not seen in training)



程序代写代做 CS编程辅导

- 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

Email: tutorcs@163.com

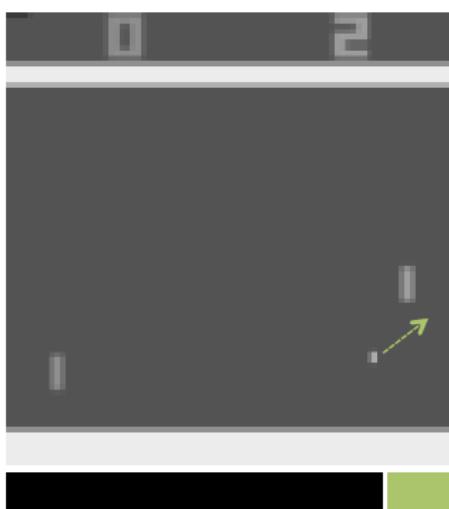
QQ: 749389476

<https://tutorcs.com>

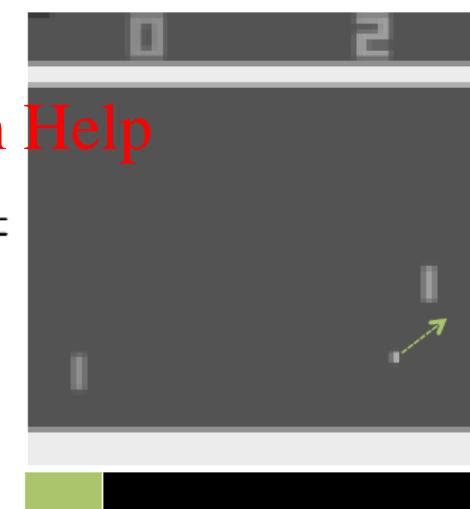
程序代写代做 CS编程辅导

- 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: $\dots, s_t + \delta_t, a_t, r'_t, s_{t+1} + \delta_{t+1}, a'_{t+1}, r'_{t+1}, s_{t+2} + \delta_{t+2}, \dots$



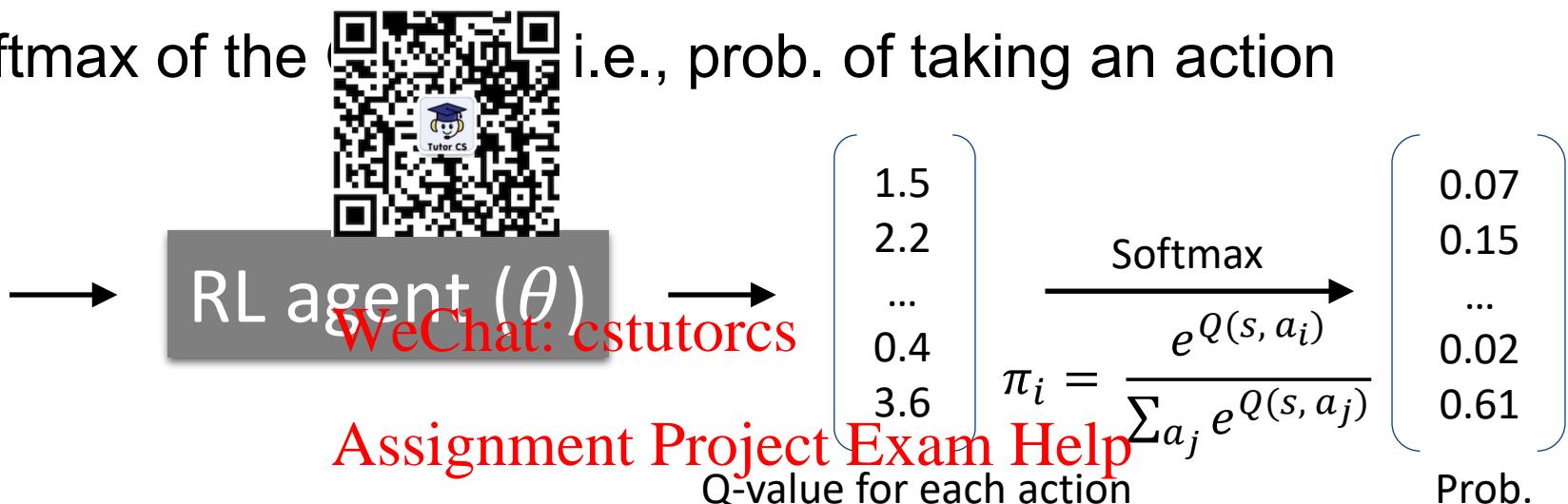
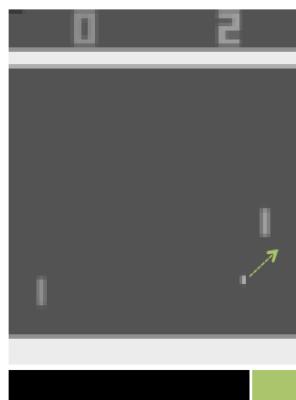
agent observed $\dots, s_t, a_t, r_t, s_{t+1}, a_{t+1}, r_{t+1}, s_{t+2}, \dots$



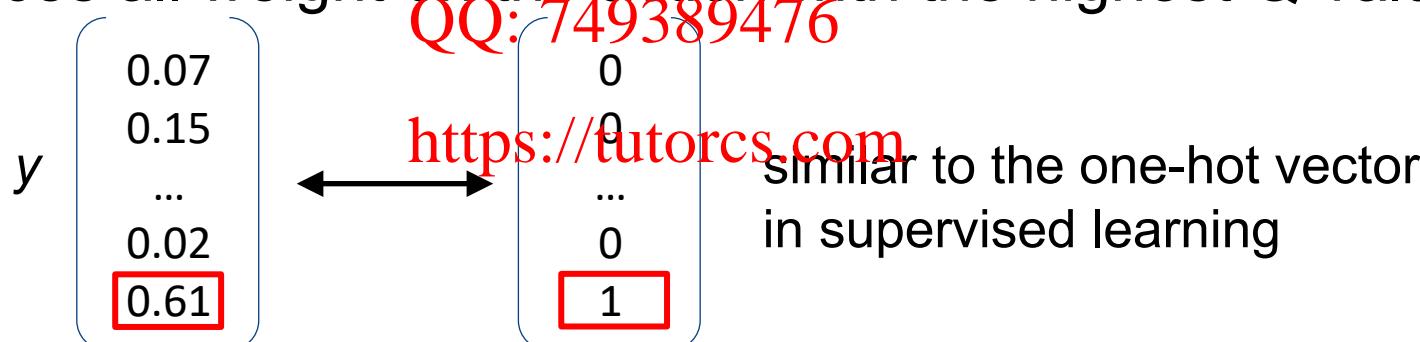
程序代写代做 CS编程辅导

- $J(\theta, x, y)$

- y : softmax of the Q-values i.e., prob. of taking an action



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



程序代写代做 CS编程辅导

- Timing of the attack

- Heuristic method



Attack the attack only when

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

WeChat: cstutorcs



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

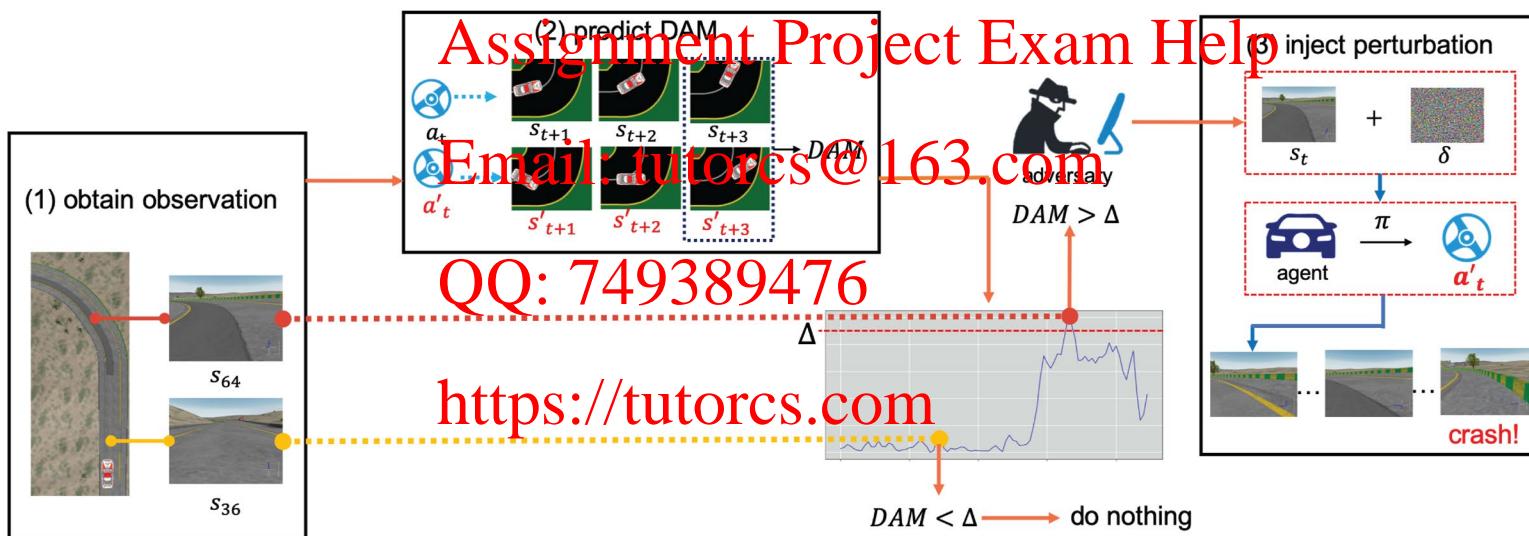
程序代写代做 CS编程辅导

- Timing of the attack [8]

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



WeChat: cstutorcs



程序代写代做 CS编程辅导

- 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

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Timing of the attack [8]

- Train an antagonist



- 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

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

• Black-box attack [9]

- Train a proxy model  learns a task that is *related* to the target agent's policy
- S threat model 
 - Only have access to the states **WeChat: cstutorcs**
 - Approximate an expectation of the state transition function **Assignment Project Exam Help**
 - $\text{psychic}(s_t, \theta_P) \approx \mathbb{E}_{\pi_T}[P(s_{t+1}|s_t)]$
- SR threat model **Email: tutorcs@163.com**
 - Have access to the states and reward **QQ: 749589476**
 - Estimate the value V of a given state under the policy π_T **<https://tutorcs.com>**
 - $\text{assessor}(s_t, \theta_A) \approx \mathbb{E}_{\pi_T}[\sum_{k=0}^{\infty} \gamma_t^{(k)} r_{t+k+1}] = V^{\pi_T}(s_t)$

程序代写代做 CS编程辅导

- Black-box attack [9]

- SA threat model

- Have access to states and actions
 - Approximate the target's policy π_T
 - $imitator(s_t, \theta_I) \sim \pi_T(s_t)$



- SRA threat model

Assignment Project Exam Help

- Have access to states, actions and rewards
 - Action-conditioned psychic (AC-psychic):
 $AC-psychic(s_t, \theta_P) \sim \mathbb{E}_{s_{t+1}}[P(s_{t+1} | s_t, a_t)]$
 - Combine assessor and AC-psychic to decide whether to perturb the state

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Black-box attack [9]

- SRA threat mode



Algorithm 1: Strategically-timed attack

Input: Trained assessor, 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

WeChat: cstutorcs

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

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Adversarial attacks against RL models

程序代写代做 CS编程辅导

- Black-box attack [9]

- Surrogate: assumes 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.

程序代写代做 CS编程辅导

- 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



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Training time attack

- Without attack: $\dots, (s_t, a_t, r_t), (s_{t+1}, a_{t+1}, s_{t+2}, r_{t+1}), \dots$
- With attack: $\dots, (s_t, a_t, r_t), (s_{t+1} + \delta_{t+1}, a'_{t+1}, s_{t+2} + \delta_{t+2}, r'_{t+1}), \dots$
- Purpose: generate δ_t such that the agent will not take the next action a_{t+1}

– Cross entropy loss: $J = -\sum_i p_i \log \pi_i$

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

WeChat: cstutorcs
Assignment Project Exam Help

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

What about targeted attacks?
Email: tutorcs@163.com

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

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

<https://tutorcs.com>

QQ: 749389476

COMP90073 Security Analysis

程序代写代做 CS编程辅导

- 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



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Adversarial training [7]

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



Untampered state → WeChat: [tutorcs](http://tutorcs.com) Potentially non-optimal action
→ explore more
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Reinforcement learning
 - State, action, reward
 - Value function, policy
 - 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



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

References

- [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.



WeChat: estutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

References

程序代写代做 CS编程辅导

- [8] Jianwen Sun and Tianwei Zhang and Xiaofei Xie and Lei Ma and Yan Zheng and Kangjie Chen  Liu, "Stealthy and Efficient Adversarial Attacks against Deep Reinforcement Learning," AAAI 2020: 5883-5891
- [9] Matthew Inkawich, Yi Zhou, 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.

WeChat: cstutors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导



Adversarial Reinforcement Learning in Autonomous Cyber Defence

WeChat: cstutorcs

Assignment Project Exam Help

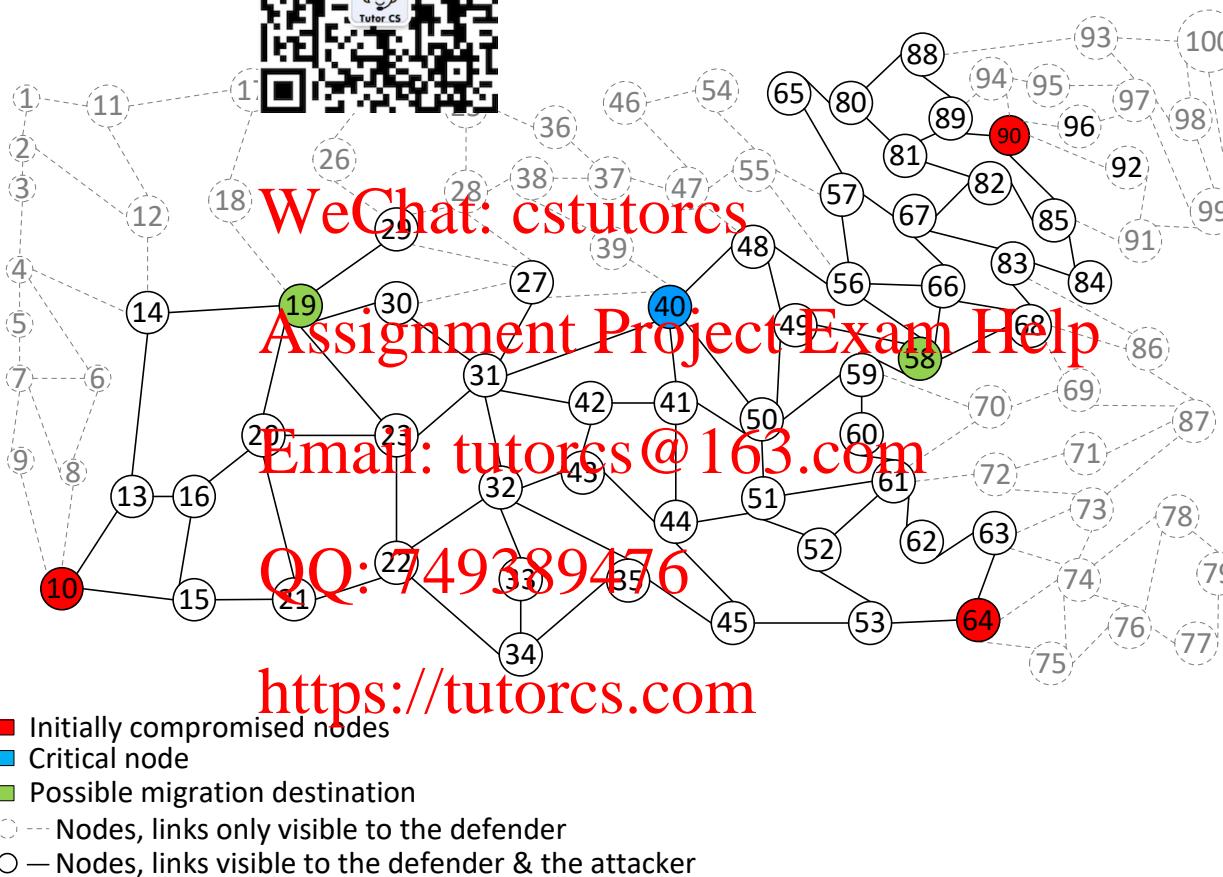
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

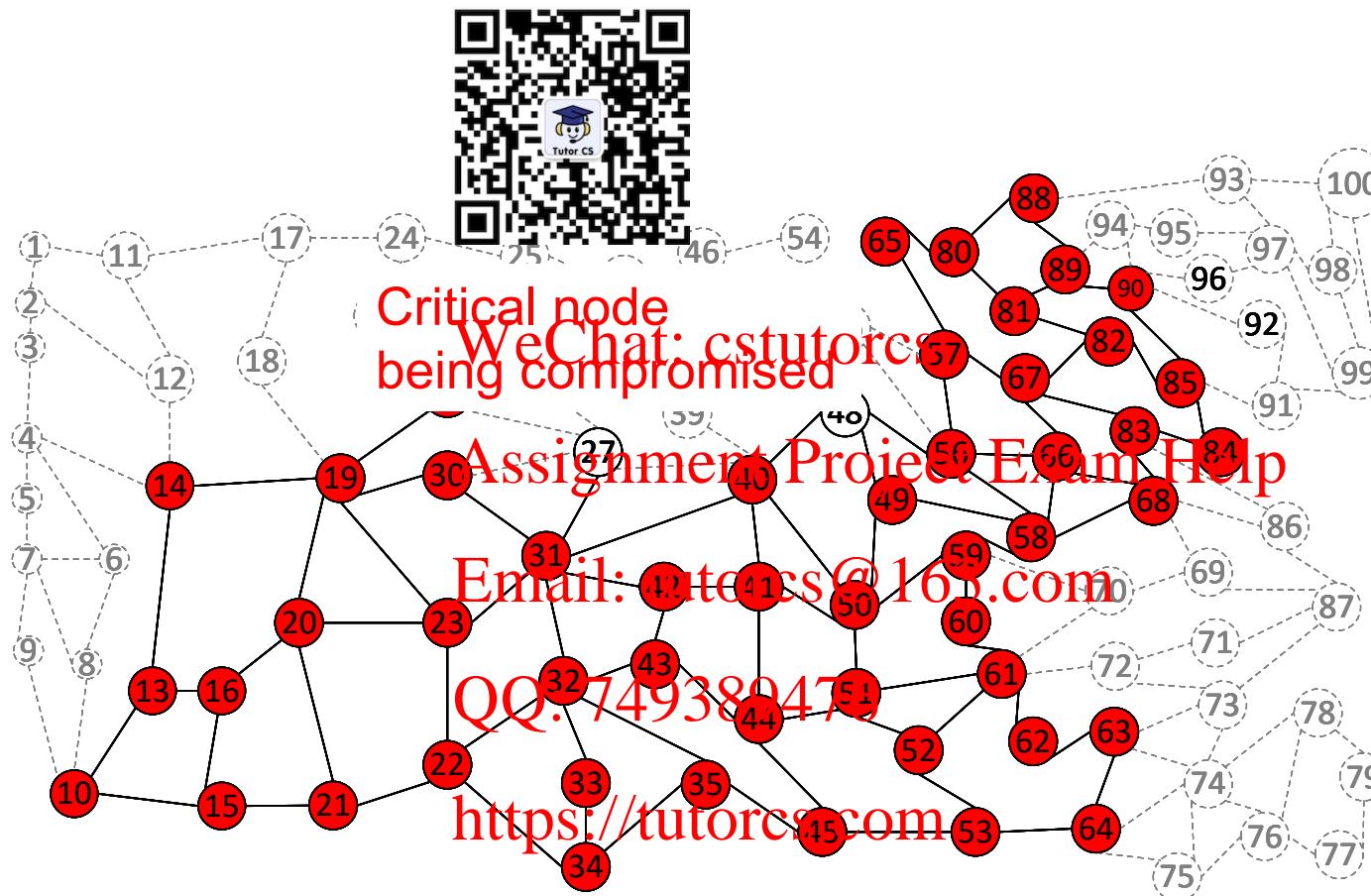
Adversarial attacks against RL models

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



Adversarial attacks against RL models

- Attacker: partial observability of the network topology



• Problem definition

- State: [0, 0, ..., 0, 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:

Email: tutorcs@163.com

- -1: (1) critical node is compromised or isolated, (2) invalid action

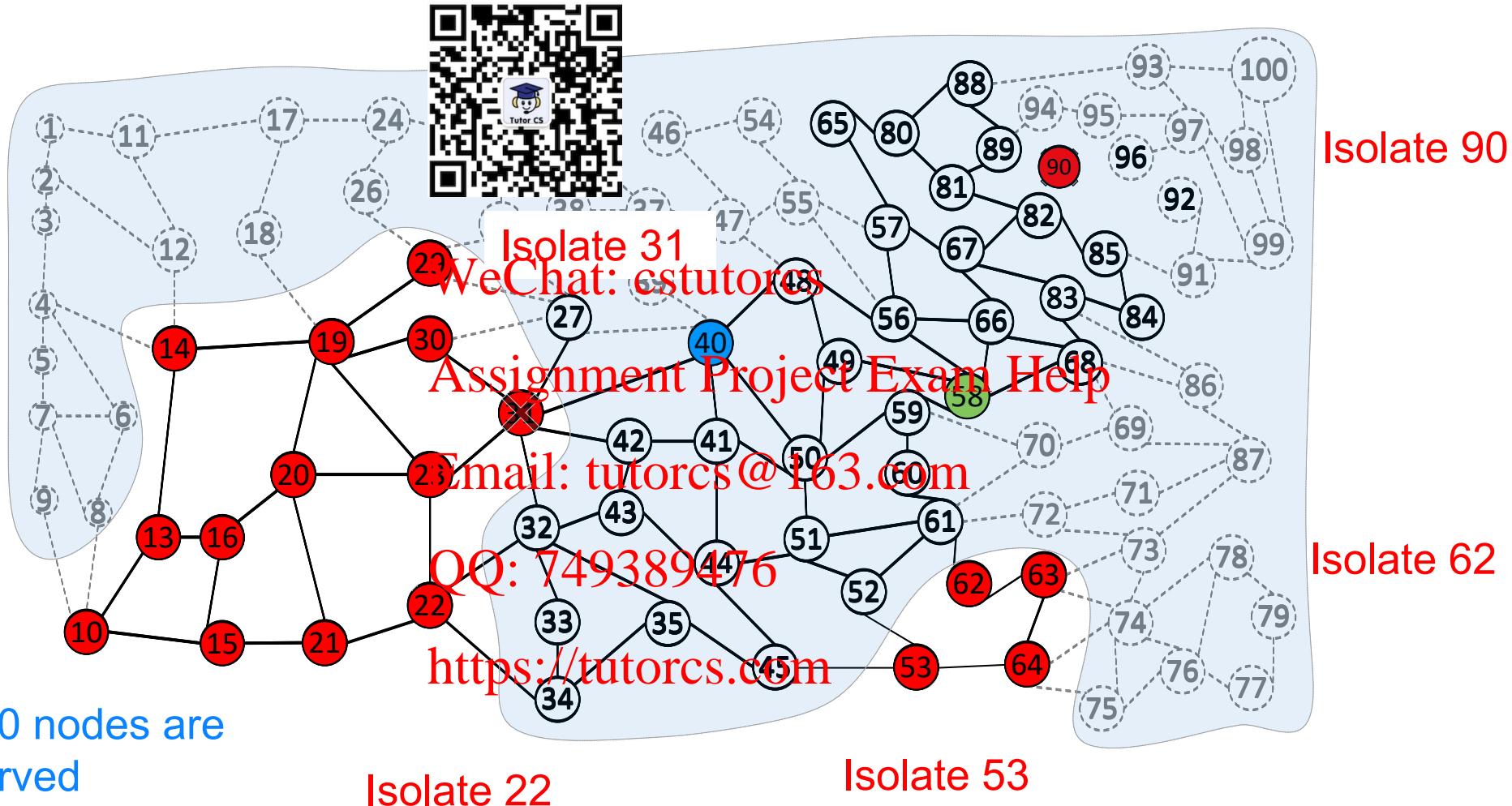
- Proportional to number of uncompromised nodes that can still reach the critical node

<https://tutorcs.com>

- Attacker can only compromise a node x if there is a visible link between x and any compromised node

Adversarial attacks against RL models

- Without training-time attacks



- Training-time attack: manipulate states to prevent agents from taking optimal actions

- $(s_t, a_t, s_{t+1}, r_t) \rightarrow (s_t, \delta_{t+1}, r'_t)$



- Binary state → 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

- $\min Q(s_{t+1} + \delta_{t+1}, a_{t+1})$: the optimal action for s_{t+1} that has been learned so far

- Loop through $L_{FP} \cup L_{FN}$ 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

程序代写代做 CS编程辅导

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 false positive (false negative) by the L_{FP} (L_{FN});

The main DQN, Q ;

Limit on the number of FPs and FNs per time, LIMIT

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



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

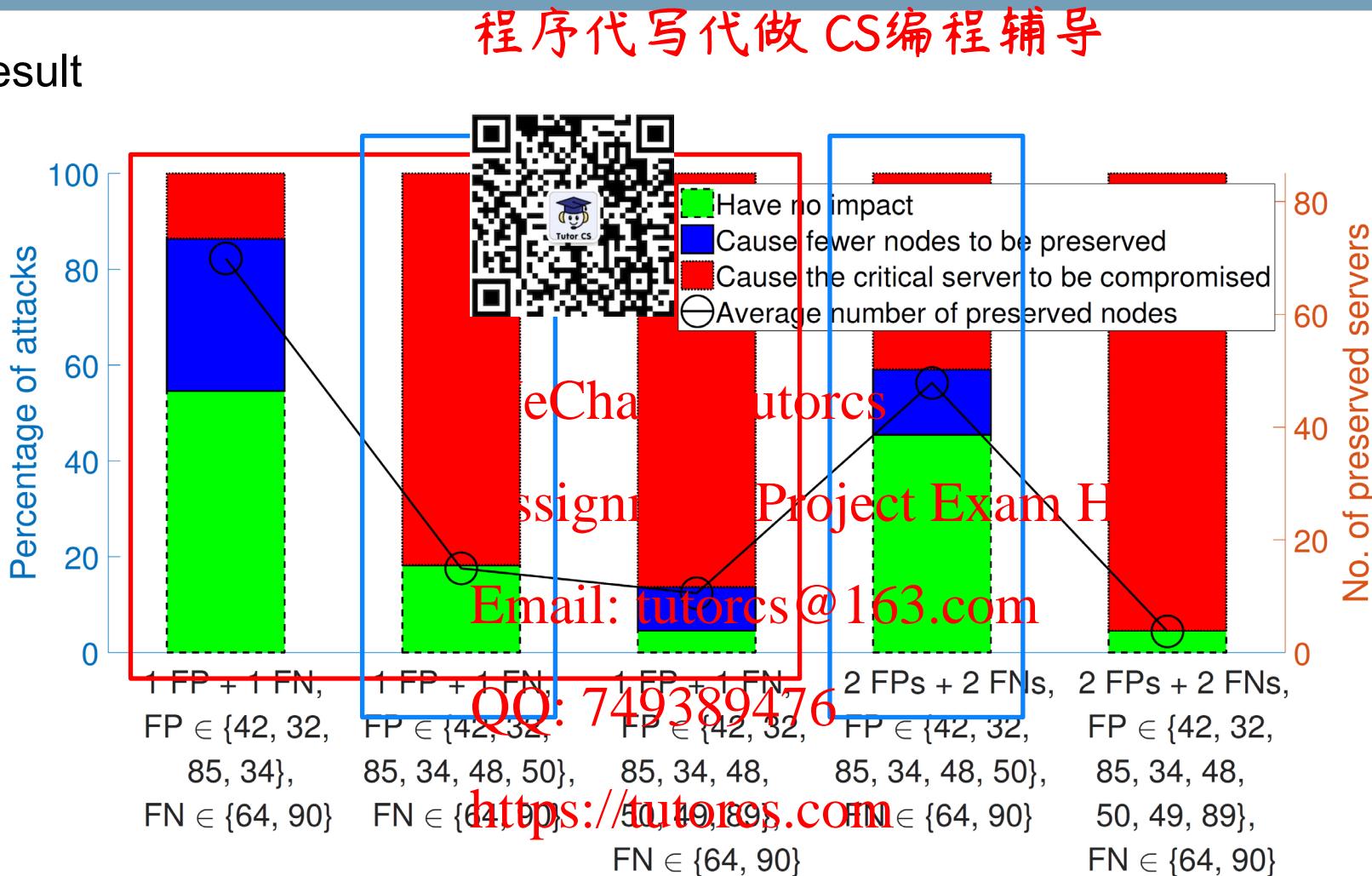
```

1  $FN = FP = \{\};$ 
2  $minQ_{FN} = minQ_{FP} = \{\};$ 
3  $a' = \text{argmax}_{a^*} Q(s', a^*);$ 
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  $|FN| > \text{LIMIT}$  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  $|FP| > \text{LIMIT}$  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  $(s, a, s' + \delta, r')$ 

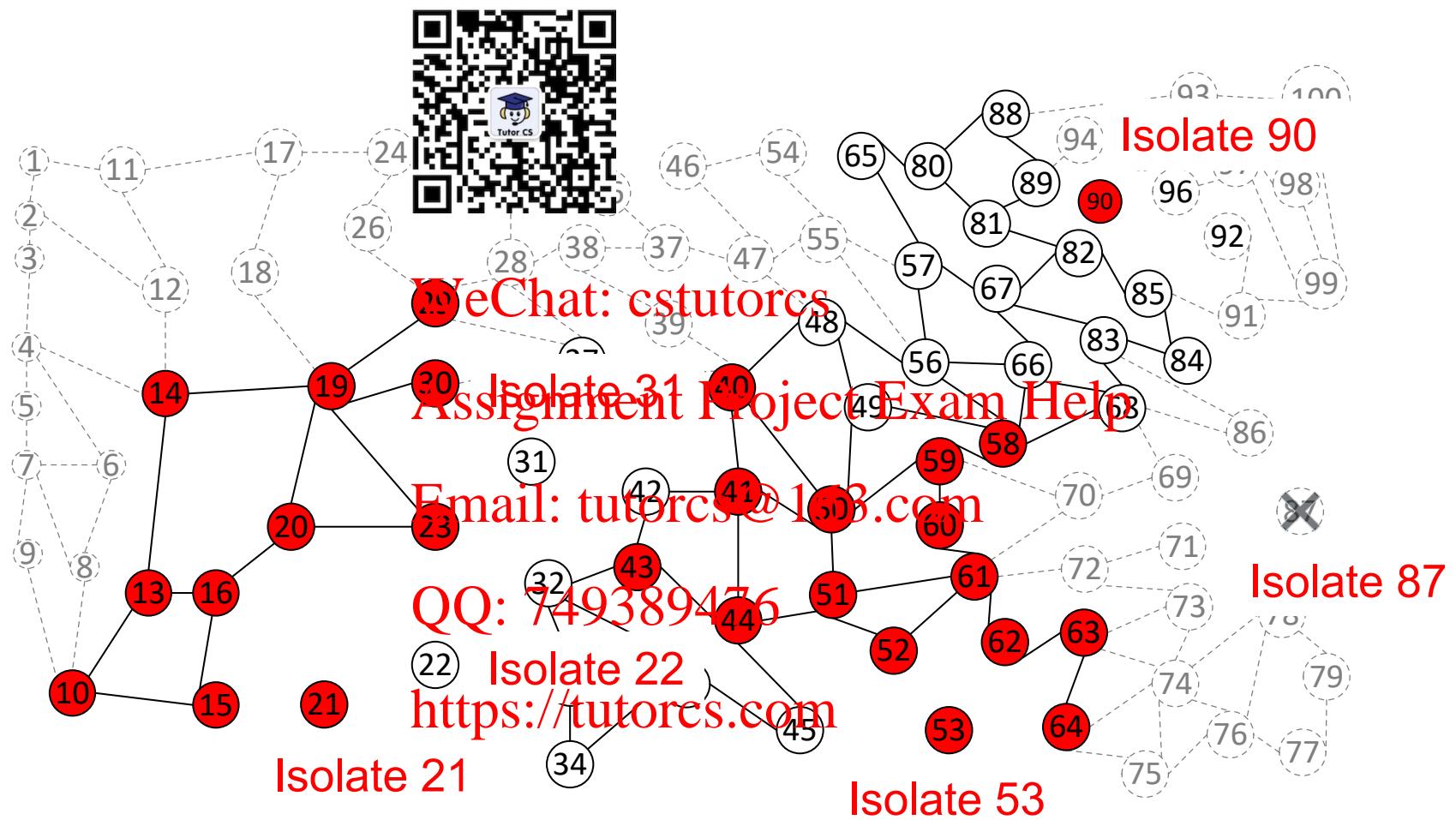
```

Adversarial attacks against RL models

- Result



- After training-time attacks



程序代写代做 CS编程辅导

- 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}	WeChat: cstutorcs Loop through all nodes
Flip K nodes	Assignment Project Exam Help Flip K' nodes

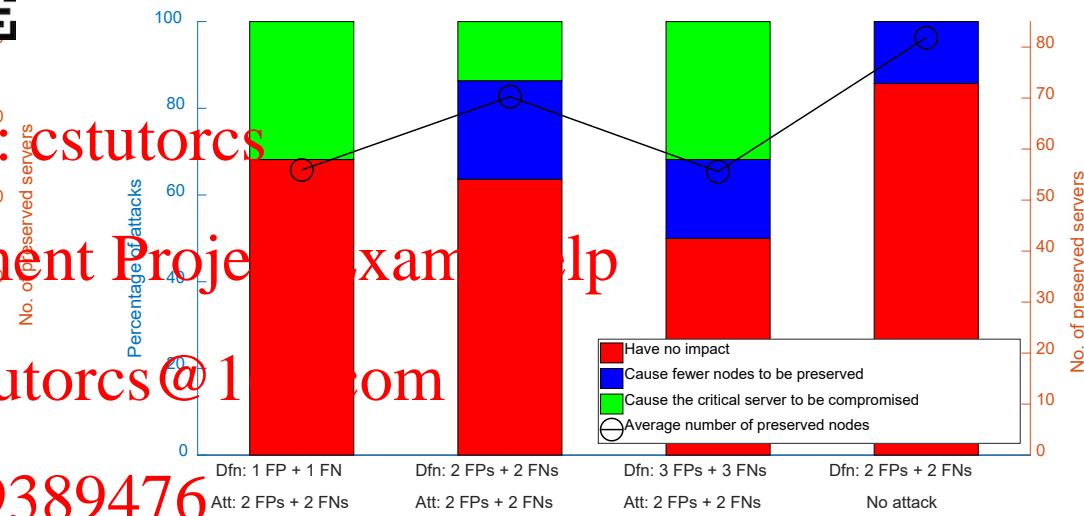
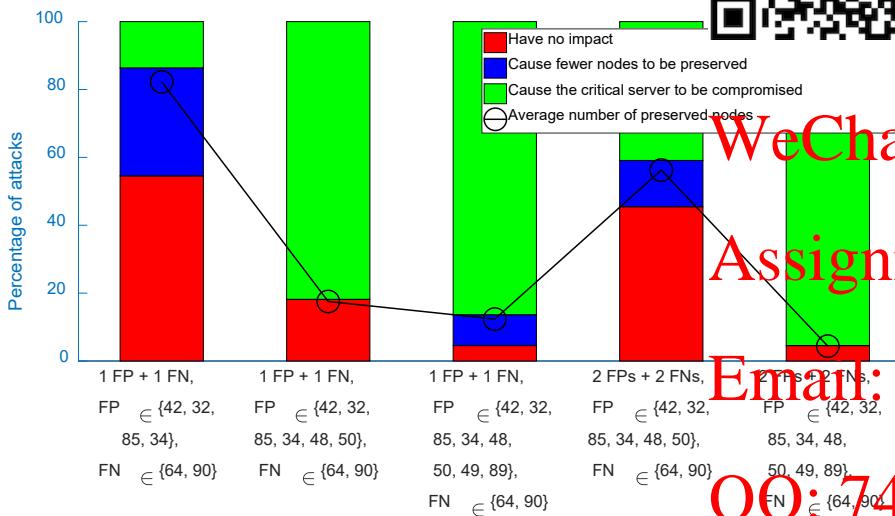
- Effective even if $K' \neq K$ Email: tutorcs@163.com
- Minimum impact on normal training process (i.e., $K = 0, K' > 0$) QQ: 749389476

<https://tutorcs.com>

Inversion defence method

程序代写代做 CS编程辅导

- Before & after the defence method is applied



QQ: 749389476

<https://tutorcs.com>