

Assignment Project Exam Help Announcements

Add WeChat powcoder

Reminder: pset5 out, due midnight today

- pset5 self-grading form out Monday, due 11/16 (1 week)

<https://powcoder.com>

- pset 6 out next week 11/12

Add WeChat powcoder



Reinforcement Learning

Assignment Project Exam Help
Deep Mind's bot playing Atari Breakout

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<https://www.youtube.com/watch?v=TmPfTpjtdgg>



Reinforcement Learning

- Plays Atari video games
- Beats human champions at Poker and Go
- Robot learns to pick up, stack blocks
- Simulated quadruped learns to run

<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

What is reinforcement learning?

Reinforcement Learning

Assignment Project Exam Help

Types of learning

Add WeChat powcoder



Supervised



Unsupervised



Reinforcement

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

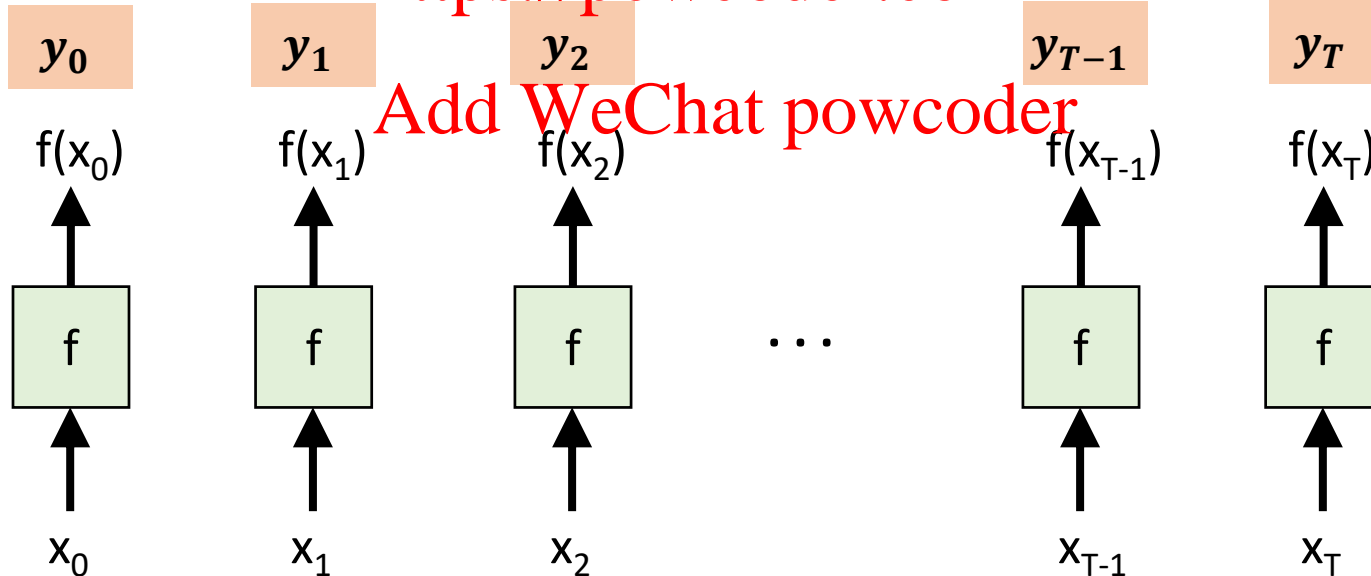
Assignment Project Exam Help

Supervised learning

Add WeChat powcoder

- model f receives input x
- also gets correct output y
- predictions do not change future inputs

Supervised learning: (in arbitrary order of examples)



Assignment Project Exam Help

Add WeChat powcoder
This is not how humans learn!

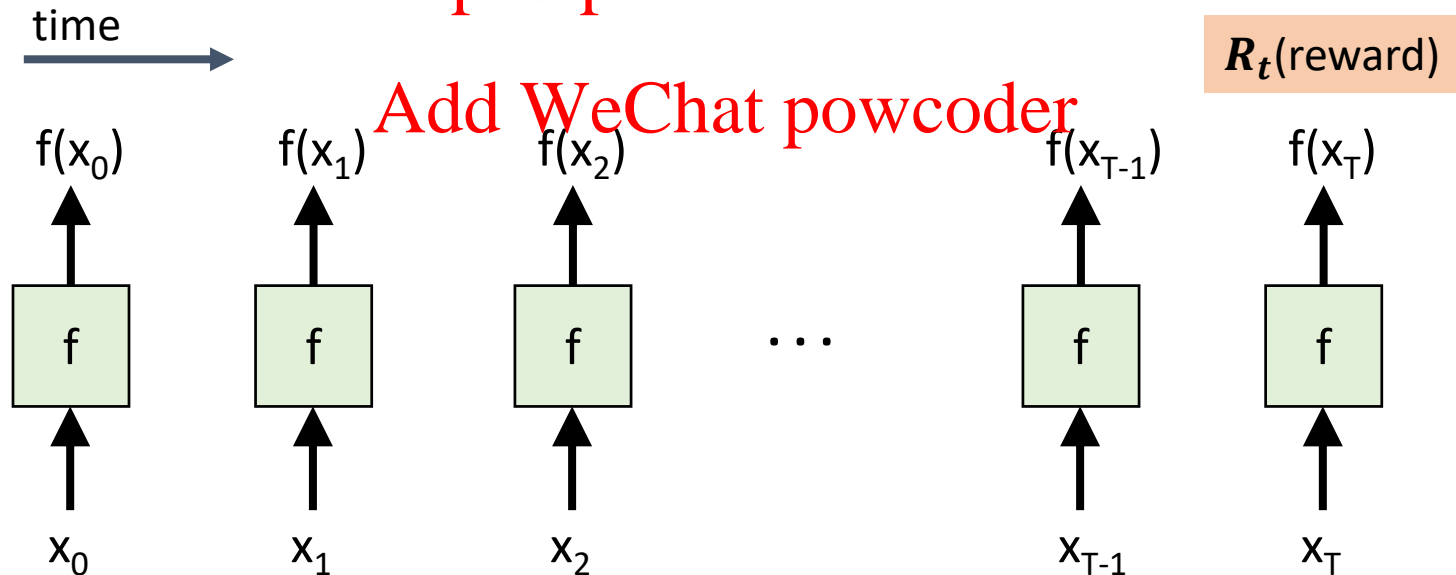


Reinforcement learning

Assignment Project Exam Help
Add WeChat powcoder

- agent receives input x , chooses action
- gets R (reward) after T time steps
- actions affect the next input (state)

Reinforcement learning:



Assignment Project Exam Help

Input is the “world’s” state

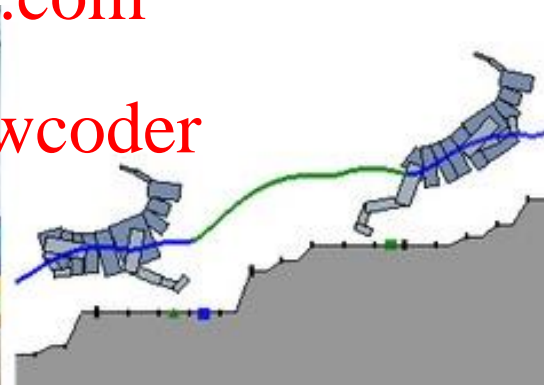
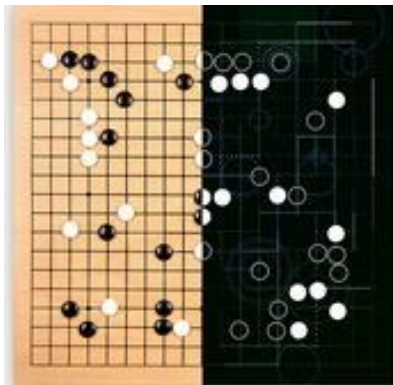
Add WeChat powcoder

- Current game board layout
- Picture of table with blocks
- Quadriped position and orientation

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

Output is an action

Add WeChat powcoder

- Which game piece to move where (discrete)
- Orientation and position of robot arm (continuous)
- Joint angles of quadruped legs (continuous)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Actions affect state!

Assignment Project Exam Help
action → reward

Add WeChat powcoder



Assignment Project Exam Help

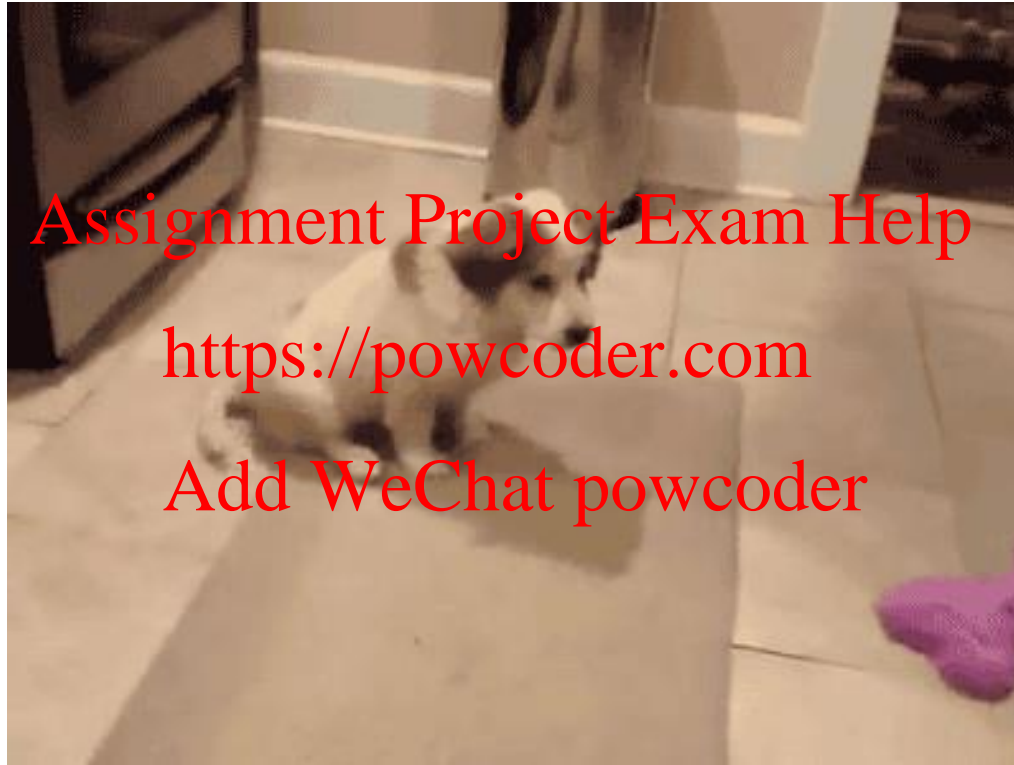
<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Add WeChat powcoder.

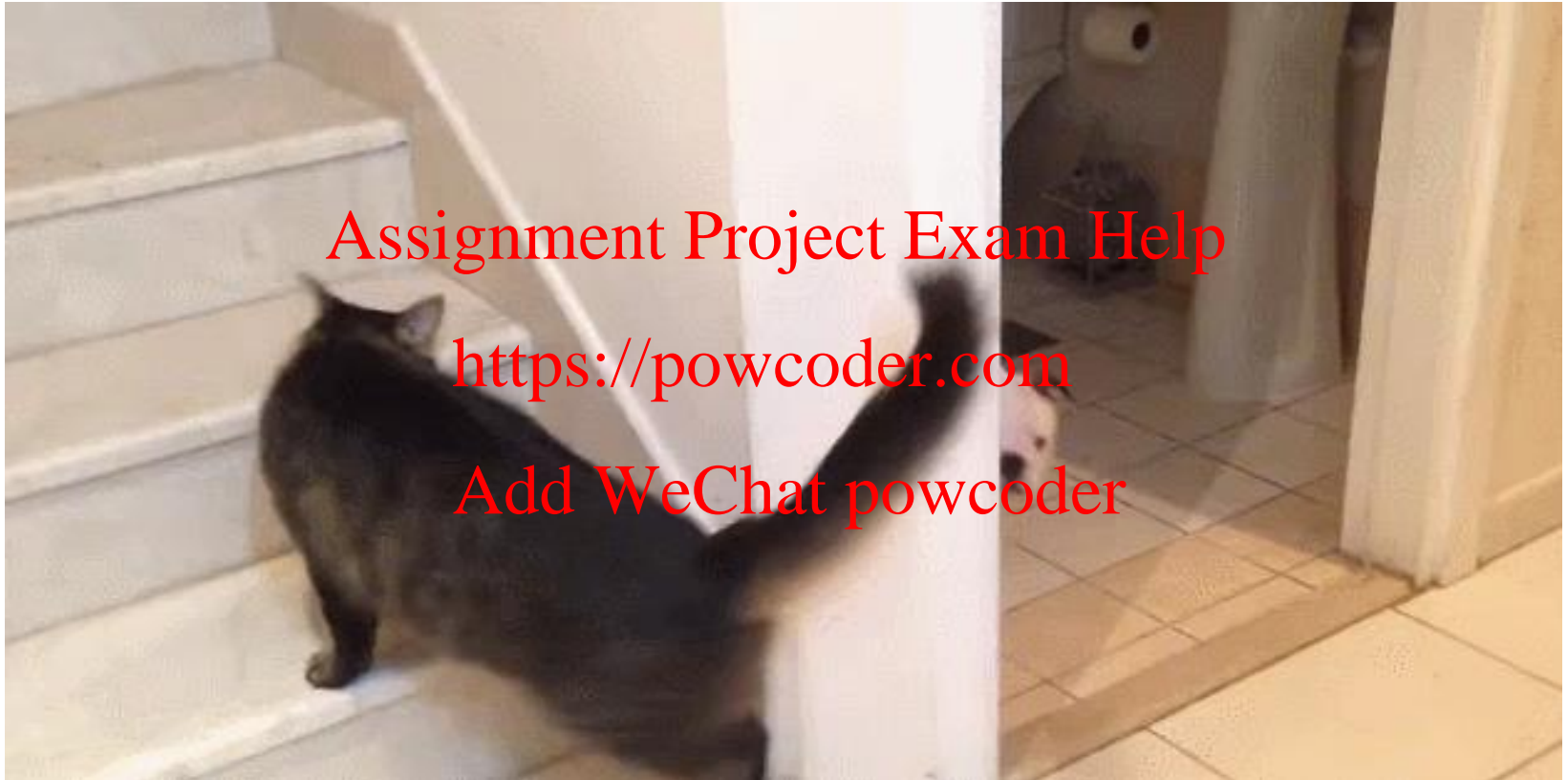
Only some actions lead to rewards



Assignment Project Exam Help

Some rewards are negative

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Reward examples

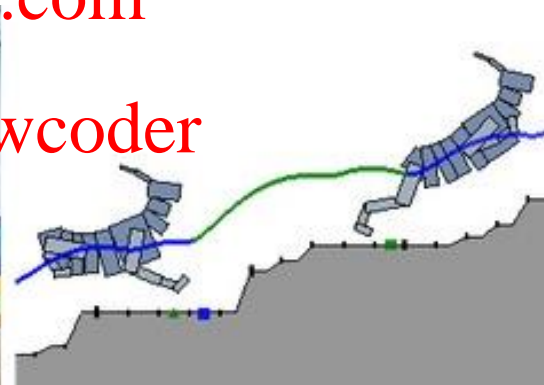
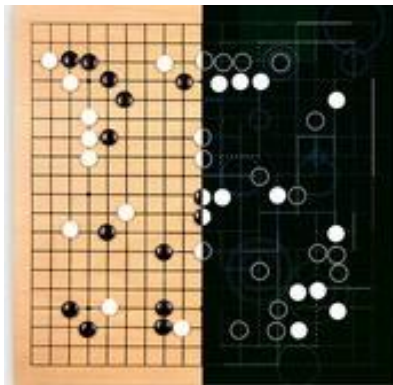
Add WeChat powcoder

- Wining the game (positive)
- Successfully picking up block (positive)
- Falling (negative)

Assignment Project Exam Help

<https://powcoder.com>

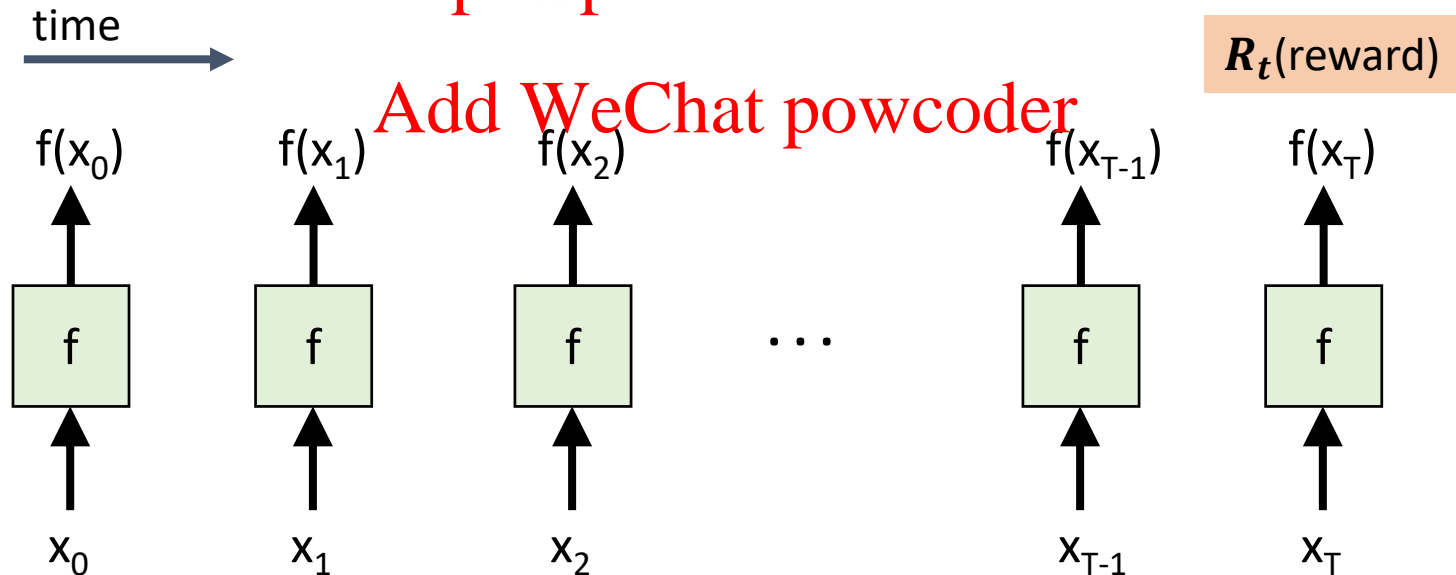
Add WeChat powcoder



Goal of reinforcement learning

- Learn to predict actions that maximize future rewards
- Need a new mathematical framework

Reinforcement learning:





Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Markov Decision Process

Reinforcement Learning

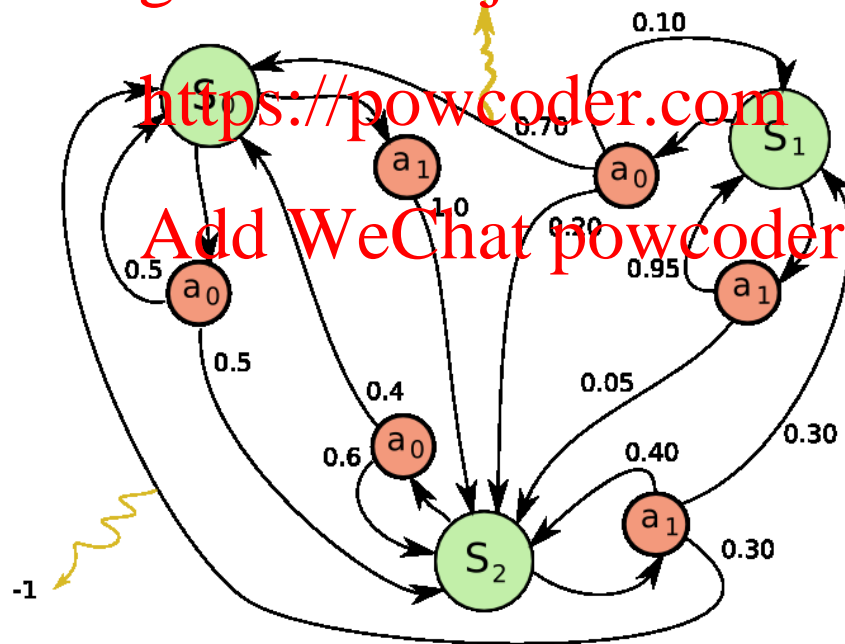
Assignment Project Exam Help

Markov Decision Process (MPD)

Add WeChat powcoder

Definition: a mathematical framework for modeling decision making in situations where outcomes are partly random and partly under the control of a decision maker.

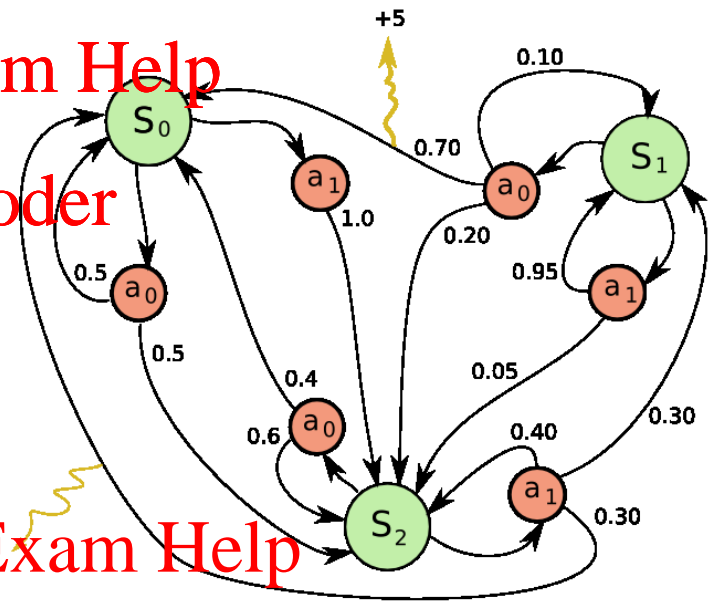
Assignment Project Exam Help



https://en.wikipedia.org/wiki/Markov_decision_process

MDP notation

Assignment Project Exam Help
Add WeChat powcoder



- S – set of States
 - A – set of Actions
 - $R: S \rightarrow \mathbb{R}$ (Reward)
 - P_{sa} – transition probabilities ($p(s, a, s') \in \mathbb{R}$)
 - γ – discount factor
- Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

$$\text{MDP} = (S, A, R, P_{sa}, \gamma)$$

Assignment Project Exam Help

Discount factor γ

Add WeChat powcoder

- discount factor prevents the total reward from going to infinity ($0 \leq \gamma \leq 1$)

Assignment Project Exam Help

- makes the agent prefer immediate rewards to rewards that are potentially received far away in the future

<https://powcoder.com>

Add WeChat powcoder

- E.g., two paths to the goal state: 1) longer path but gives higher reward 2) shorter path with smaller reward; the γ value controls which the path the agent should prefer

MDP (Simple example)

Assignment Project Exam Help
Add WeChat powcoder



Assignment Project Exam Help


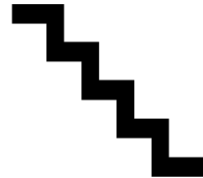

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

MDP (Simple example)

Add WeChat powcoder

	1	2	3	4
1				
2				
3				

Assignment Project Exam Help

MDP (Simple example)

- States S = locations
- Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$

	1	2	3	4
1				
2				
3				



<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

MDP (Simple example)

- States S = locations
- Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$
- Reward $R: S \rightarrow \mathbb{R}$

	1	2	3	4
1	0	0	0	+1
2	0		0	-1
3	0	0		0



<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

MDP (Simple example)

- States S = locations
- Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$
- Reward $R: S \rightarrow \mathbb{R}$

	1	2	3	4
1	-0.02	-0.02	-0.02	+1
2	-0.02		-0.02	-1
3	-0.02	-0.02		-0.02

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

MDP (Simple example)

- States S = locations
- Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$
- Reward $R: S \rightarrow \mathbb{R}$
- Transition P_{sa}

	1	2	3	4
1	-0.02	-0.02	-0.02	+1
2	-0.02		-0.02	-1
3	-0.02	-0.02		-0.02

$$P_{(3,3),\uparrow}((2,3)) = 0.8$$

$$P_{(3,3),\uparrow}((3,4)) = 0.1$$

$$P_{(3,3),\uparrow}((3,2)) = 0.1$$

$$P_{(3,3),\uparrow}((1,3)) = 0$$

$$\vdots$$

Assignment Project Exam Help

MDP - Dynamics

Add WeChat powcoder

- Start from state S_0
- Choose action A_0
- Transit to $S_1 \sim P_{S_0 A_0}$
- Continue...

	1	2	3	4
1	-0.02	-0.02	-0.02	+1
2	-0.02		-0.02	-1
3				-0.02

<https://powcoder.com>

-0.02

-0.02

-0.02

Add WeChat powcoder

- Total payoff:

-0.02

-0.02

-0.02

$$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$$

Assignment Project Exam Help

How do we choose good actions?



Choosing actions in MDP

Assignment Project Exam Help
Add WeChat powcoder

States S = locations

Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$

Reward $R: S \rightarrow \mathbb{R}$

Transition P_{sa}

	1	2	3	4
1	-.02	-.02	-.02	+1
2	-.02		-.02	-1
3	-.02	-.02		-.02

<https://powcoder.com>

- Goal - Choose actions as to maximize expected total payoff:

$$E [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

- In our example:

R – get to charge station, avoid stairs

γ – discourage long paths, how much to delay reward

MDP – Policy π


Add WeChat powcoder

States S = locations

Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$

Reward $R: S \rightarrow \mathbb{R}$

Transition P_{sa}

	1	2	3	4
1				+1
2				-1
3				

<https://powcoder.com>

Add WeChat powcoder

- Goal - Choose actions as to maximize expected total payoff:

$$E [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

- Policy is a function $\pi: S \rightarrow A$



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Policy Value and Q functions

Reinforcement Learning

MDP – Policy value function

Assignment Project Exam Help

Add WeChat powcoder

States S = locations

Actions $A = \{\uparrow, \rightarrow, \leftarrow, \downarrow\}$

Reward $R: S \rightarrow \mathbb{R}$

Transition P_{sa}

Policy $\pi: S \rightarrow A$

	1	2	3	4
1				+1
2				-1
3				

Assignment Project Exam Help

<https://powcoder.com>

- Value function for policy $\pi: S \rightarrow A$

Add WeChat powcoder

$$V^\pi(s) = \mathbb{E} [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \mid s_0 = s, \pi]$$

Expected sum of discounted rewards

MDP – Policy value function

Assignment Project Exam Help

Add WeChat powcoder

$$V^\pi(s) = \mathbb{E} [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi]$$

$$\Rightarrow V^\pi(s) = E[R(s_0)] + E[\gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

Assignment Project Exam Help

this is recursion!

Bellman's equation:

<https://powcoder.com>

$$V^\pi(s) = R(s) + \gamma E_{s' \sim P_{s, \pi(s)}} [V(s')]$$

Add WeChat powcoder



expectation over values of next state

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Assignment Project Exam Help

MDP – Policy value function

Add WeChat powcoder

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Assignment Project Exam Help

solve
|S|
eqs. {

$$\begin{aligned} V_1 &= R(1) + \gamma P_{1,\downarrow}(2) V_2 + \dots \\ V_2 &= \dots \\ &\dots \\ V_{10} &= R(10) + \gamma (P_{10,\uparrow}(6) V_6 + P_{10,\uparrow}(9) V_9 + P_{10,\uparrow}(11) V_{11}) \\ &\dots \end{aligned}$$

<https://powcoder.com>

Add WeChat powcoder

	1	2	3	4
1	1	2	3	4
2	5		6	7
3	8	9	10	11

	1	2	3	4
1	↓	→	→	+1
2	←		↑	-1
3	→	←	↑	↓

Policy π

Assignment Project Exam Help

Optimal value function

Add WeChat powcoder

If the agent uses a given policy π to select actions, the corresponding **value function** is given by:

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{\pi}(s'|s) V^\pi(s')$$

Assignment Project Exam Help

There exists an **optimal value function** that has higher value than other functions for all states

<https://powcoder.com>

Add WeChat powcoder

$$V^*(s) = \max_{\pi} V^\pi(s) \quad \forall s \in \mathcal{S}$$

The **optimal policy** π^* is the policy that corresponds to the optimal value function

$$\pi^* = \arg \max_{\pi} V^\pi(s) \quad \forall s \in \mathcal{S}$$

Assignment Project Exam Help

Value function vs. Q-function

Add WeChat powcoder

For convenience, RL algorithms introduce the **Q-function**, which takes a state-action pair and returns a real value

Assignment Project Exam Help

The **optimal Q-function** $Q^*(s, a)$ means the highest expected total reward received by an agent starting in s and choosing action a which maximizes value over

<https://powcoder.com>

Add WeChat powcoder

$$V^*(s) = \max_a Q^*(s, a) \quad \forall s \in \mathcal{S}$$

$Q^*(s, a)$ is an indication for how good it is for an agent to pick action a while being in state s

Optimal Q-function

Assignment Project Exam Help
Add WeChat powcoder

If we know the optimal Q-function $Q^*(s, a)$, the optimal policy π^* can be easily extracted by choosing the action a that gives maximum $Q^*(s, a)$ for state s :

$$\pi^*(s) = \arg \max_a Q^*(s, a) \quad \forall s \in \mathcal{S}$$

Assignment Project Exam Help
<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

RL Approaches

Reinforcement Learning

Assignment Project Exam Help

Reinforcement learning approaches

Add WeChat powcoder

Value-based RL

- Estimate the optimal value function
- *i.e.*, the maximum value achievable under any policy
- Guaranteed to converge to optimum

Policy-based RL

- Search directly for the optimal policy
- re-define the policy at each step and compute the value according to this new policy until the policy converges
- Guaranteed to converge, often faster than value

Q-learning

- Search for the optimal Q-function
- No prior knowledge of MDP required

Assignment Project Exam Help

Value Iteration Algorithm

Add WeChat powcoder

Given $P_{s,a}(s') = p(s'|s, a)$, iteratively compute the Q and value functions using Bellman's equation.

Assignment Project Exam Help

```
Initialize  $V(s)$  to arbitrary values
Repeat
  For all  $s \in S$ 
    For all  $a \in \mathcal{A}$ 
       $Q(s, a) \leftarrow E[r|s, a] + \gamma \sum_{s' \in S} P(s'|s, a)V(s')$ 
     $V(s) \leftarrow \max_a Q(s, a)$ 
Until  $V(s)$  converge
```

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Policy Iteration Algorithm

Add WeChat powcoder

Given $P_{s,a}(s') = p(s'|s, a)$, π , iteratively compute the policy's value function and improve the policy to maximize it

Assignment Project Exam Help

Initialize a policy π' arbitrarily

Repeat

- $\pi \leftarrow \pi'$
- Compute the values using π by solving the linear equations
$$V^\pi(s) = E[r|s, \pi(s)] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, \pi(s)) V^\pi(s')$$
- Improve the policy at each state
$$\pi'(s) \leftarrow \arg \max_a (E[r|s, a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s'))$$

Until $\pi = \pi'$

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Reinforcement learning approaches

Add WeChat powcoder

Optimal value function

need P_{sa}

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Bellman: $V^*(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$

Optimal policy

<https://powcoder.com>

need π, P_{sa}

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

Add WeChat powcoder

Optimal state-action value function Q

easier!

Define $Q: S \times A \rightarrow \mathbb{R}$

$$\text{Bellman: } Q^*(s, a) = R(s) + \gamma \max_{a \in A} Q^*(s', a)$$



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Q-Learning (discrete)

Reinforcement Learning

Q-value function

Add WeChat powcoder

- ▶ A **value function** is a prediction of future reward
 - ▶ “How much reward will I get from action a in state s ?”
- ▶ **Q-value function** gives expected total reward
 - ▶ from state s and action a
 - ▶ under policy π
 - ▶ with discount factor γ

$$Q^{\pi}(s, a) = \mathbb{E} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

- ▶ Value functions decompose into a Bellman equation

$$Q^{\pi}(s, a) = \mathbb{E}_{s', a'} [r + \gamma Q^{\pi}(s', a') \mid s, a]$$

Assignment Project Exam Help

Optimal Q-value function

Add WeChat powcoder

- ▶ An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$$

- ▶ Once we have Q^* , we can act optimally.

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

<https://powcoder.com>

- ▶ Optimal value maximises over all decisions. Informally:

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

Add WeChat powcoder

- ▶ Formally, optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

Assignment Project Exam Help

Q-learning algorithm

Add WeChat powcoder

The agent interacts with the environment, updates Q recursively

```
initialize  $Q$  with zeros, run_actions, arbitrarily  
observe initial state  $s$   
repeat  
    select and carry out an action  $a$   
    observe reward  $r$  and new state  $s'$   
     $Q[s, a] = Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'] - Q[s, a])$   
     $s = s'$   
until terminated
```

current value

learning rate

discount

largest increase over all
possible actions in new state

Q-learning example

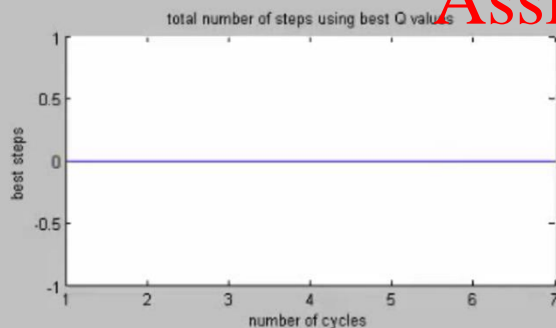
Assignment Project Exam Help

Goal: get from bottom left to top right

Add WeChat powcoder

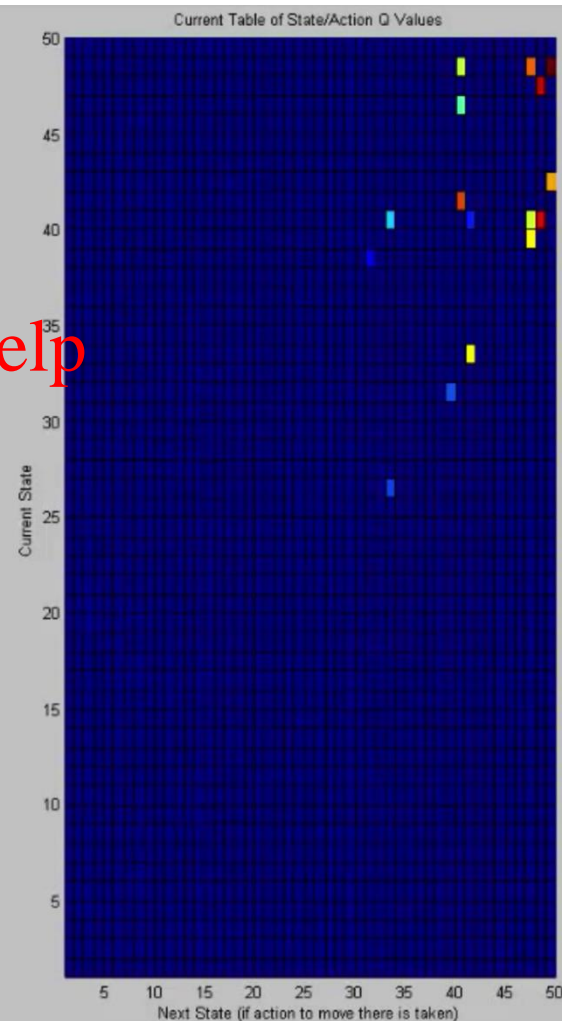
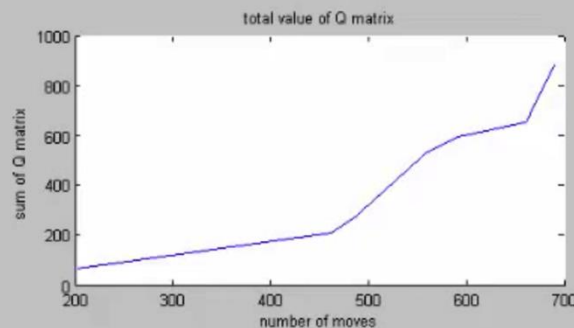
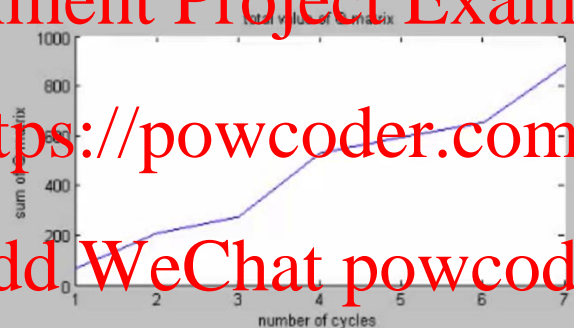
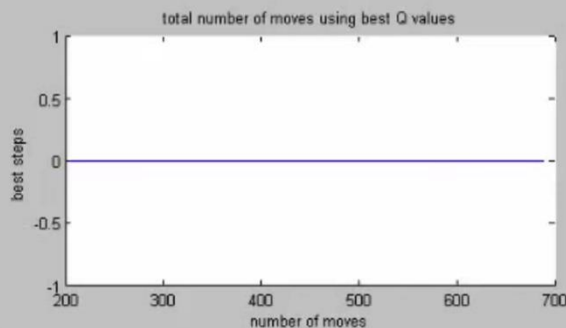


Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder



<https://www.youtube.com/watch?v=R88CiN7dTZc>

Assignment Project Exam Help

Exploration vs exploitation

Add WeChat powcoder

- How does the agent select actions during learning? Should it trust the learned values of $Q(s, a)$ to select actions based on it? or try other actions hoping this may give it a better reward?
- This is known as the exploration vs exploitation dilemma
- Simple ϵ -greedy approach: at each step with small probability ϵ , the agent will pick a random action (explore) or with probability $(1-\epsilon)$ the agent will select an action according to the current estimate of Q-values
- The ϵ value can be decreased overtime as the agent becomes more confident with its estimate of Q-values



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Continuous state

Reinforcement Learning

Assignment Project Exam Help

Continuous state - Pong

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<https://www.youtube.com/watch?v=YOW8m2YGtRg>

MDP for Pong

Add WeChat powcoder



In this case, what are these?

- S – set of States
- A – set of Actions
- $R: S \rightarrow \mathbb{R}$ (Reward)
- P_{sa} – transition probabilities ($p(s, a, s') \in \mathbb{R}$)

Add WeChat powcoder

Can we learn Q-value?

- Can discretize state space, but it may be too large
- Can simplify state by adding domain knowledge (e.g. paddle, ball), but it may not be available
- Instead, use a neural net to learn good features of the state!



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Deep RL

Reinforcement Learning

Assignment Project Exam Help

Deep RL playing DOTA

Add WeChat powcoder



https://www.youtube.com/watch?v=eHipy_j29Xw

Assignment Project Exam Help

Deep RL

Add WeChat powcoder

- V , Q or π can be approximated with deep network

- Deep Q-Learning

- Input: state, action
- Output: Q-value

Cover today

- Alternative: learn a Policy Network

- Input: state
- Output: distribution over actions

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

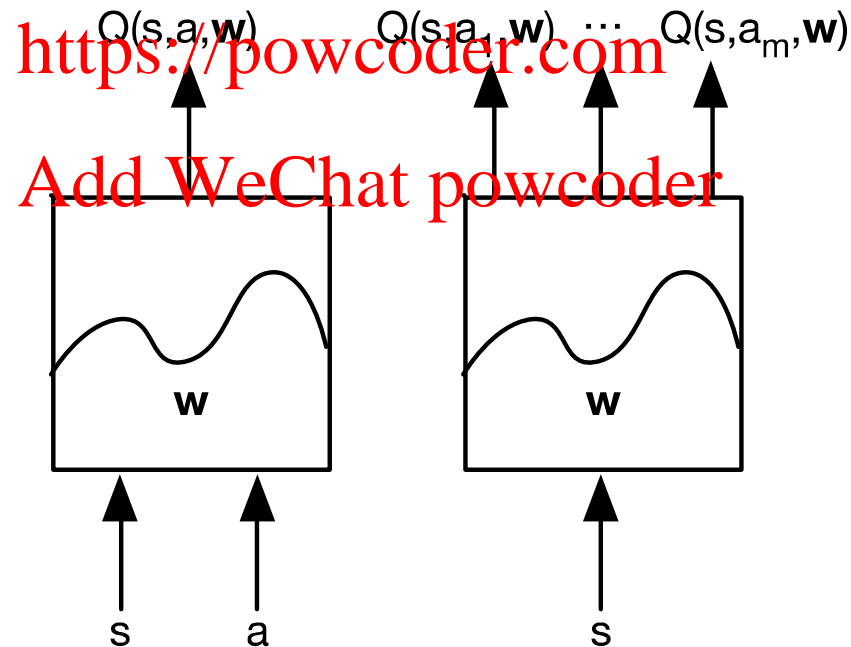
Q-value network

Add WeChat powcoder

Represent value function by **Q-network** with weights **w**

$$Q(s, a, \mathbf{w}) \approx Q^*(s, a)$$

Assignment Project Exam Help



Assignment Project Exam Help

Q-value network

Add WeChat powcoder

- ▶ Optimal Q-values should obey Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q(s', a')^* \mid s, a \right]$$

Assignment Project Exam Help

- ▶ Treat right-hand side $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$ as a target
- ▶ Minimise MSE loss by stochastic gradient descent

Add WeChat powcoder

$$l = \left(r + \gamma \max_a Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

- ▶ Converges to Q^* using table lookup representation
- ▶ But **diverges** using neural networks due to:
 - ▶ Correlations between samples
 - ▶ Non-stationary targets

Assignment Project Exam Help

Deep Q-network (DQN)

Add WeChat powcoder

To remove correlations, build data-set from agent's own experience

	s_1, a_1, r_2, s_2
	s_2, a_2, r_3, s_3
	s_3, a_3, r_4, s_4
	\vdots
$s_t, a_t, r_{t+1}, s_{t+1}$	$\rightarrow s_t, a_t, r_{t+1}, s_{t+1}$

<https://powcoder.com>

Add WeChat powcoder

Sample experiences from data-set and apply update

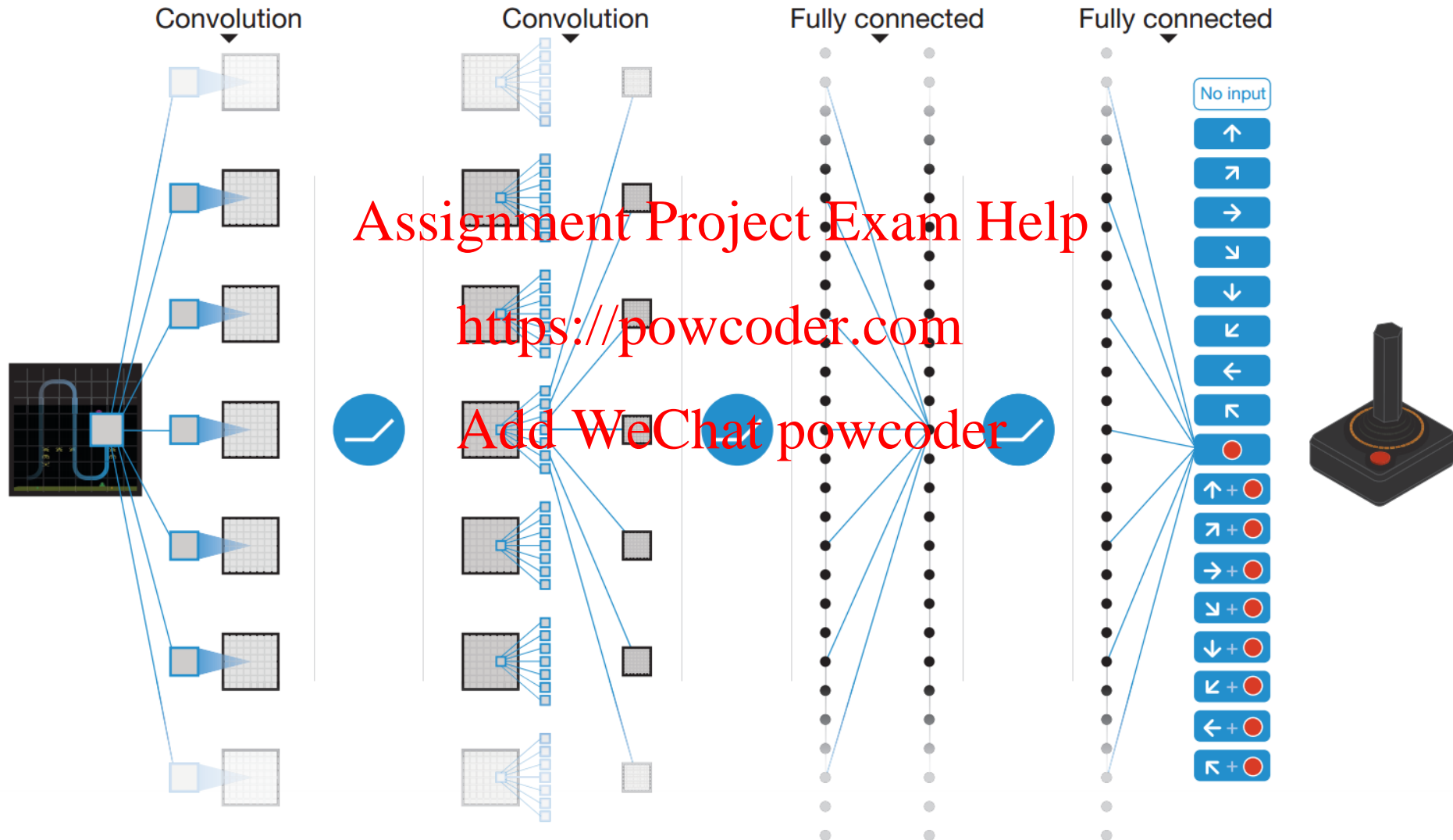
$$l = \left(r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

To deal with non-stationarity, target parameters \mathbf{w}^- are held fixed

Assignment Project Exam Help

DQN - Playing Atari

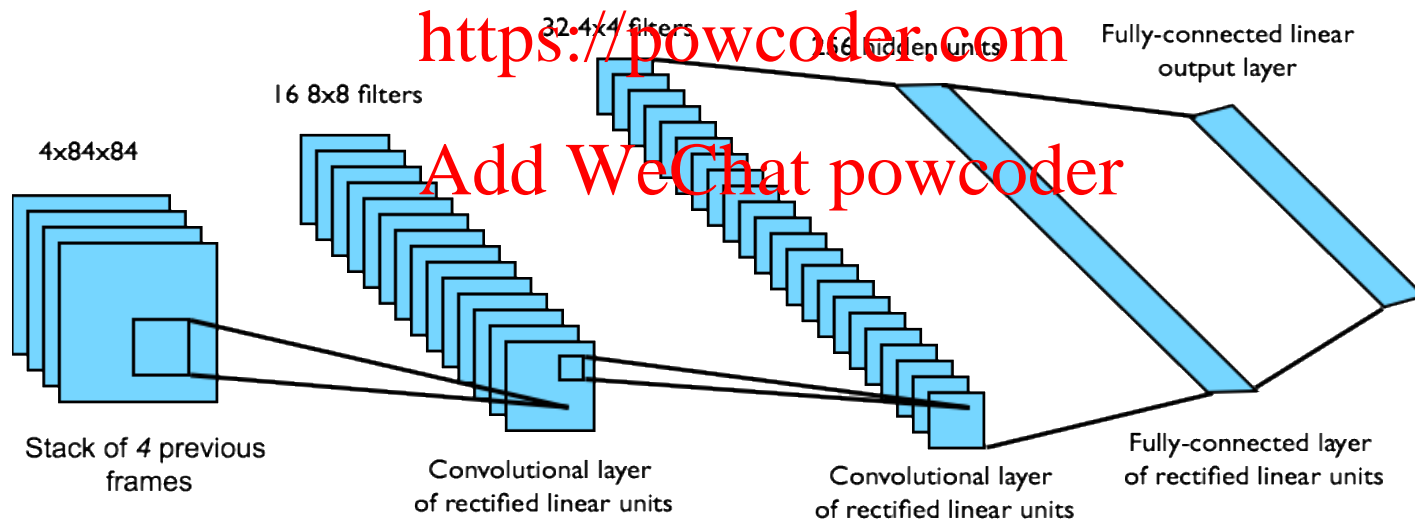
Add WeChat powcoder



DQN - Playing Atari

Assignment Project Exam Help
Add WeChat powcoder

- | End-to-end learning of values $Q(s, a)$ from pixels s
- | Input state s is stack of raw pixels from last 4 frames
- | Output is $Q(s, a)$ for 18 joystick/ button positions
- | Reward is change in score for that step



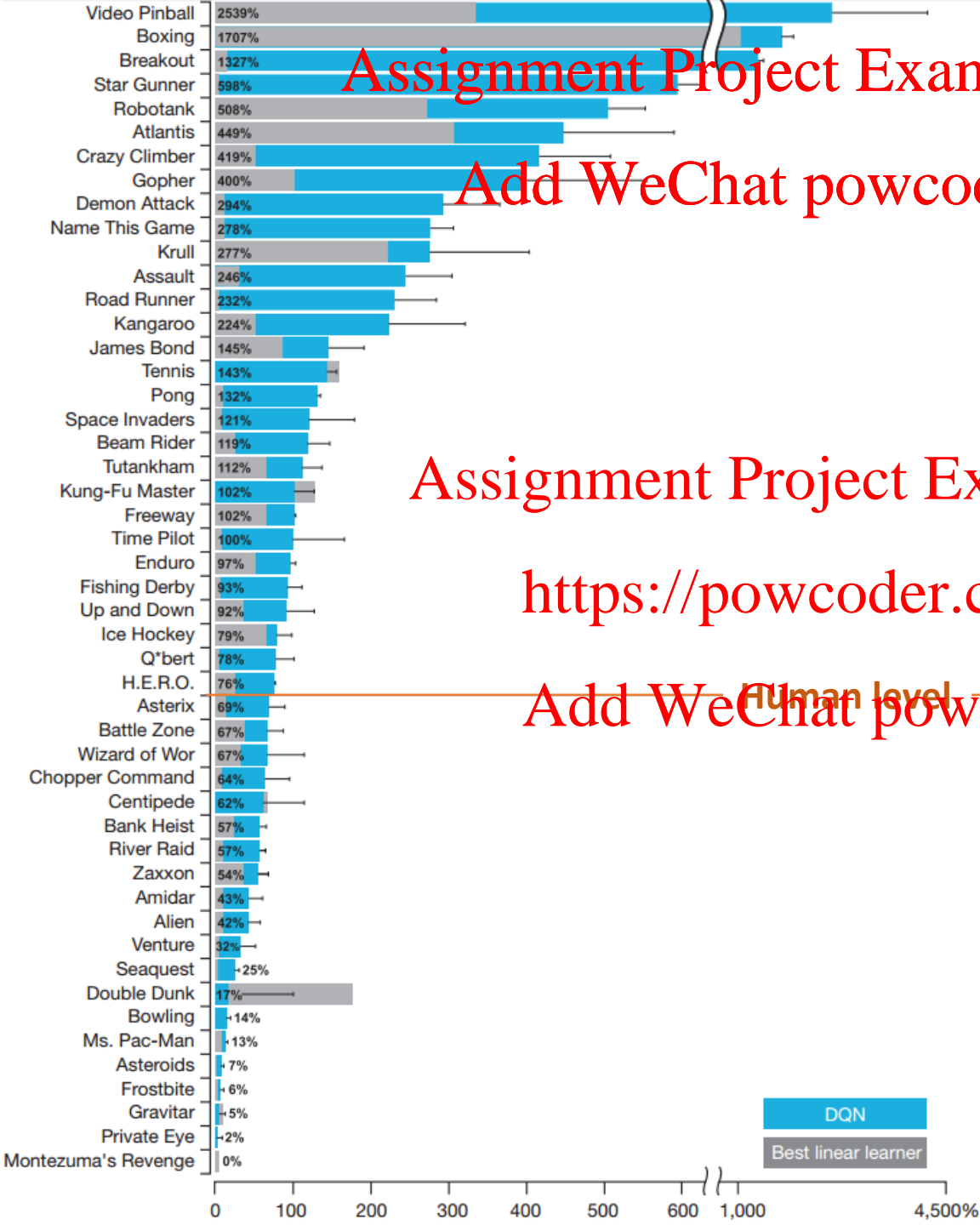
Network architecture and hyperparameters fixed across all games

DQN - Playing Atari

Assignment Project Exam Help
Add WeChat powcoder

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N
Initialize action-value function Q with random weights
for episode = 1, M **do**
 Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
 for $t = 1, T$ **do**
 With probability ϵ select a random action a_t
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
 Execute action a_t in emulator and observe reward r_t and image x_{t+1}
 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}
 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
 end for
end for



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Human level

DQN for Atari

Assignment Project Exam Help
Add WeChat powcoder

DQN paper:

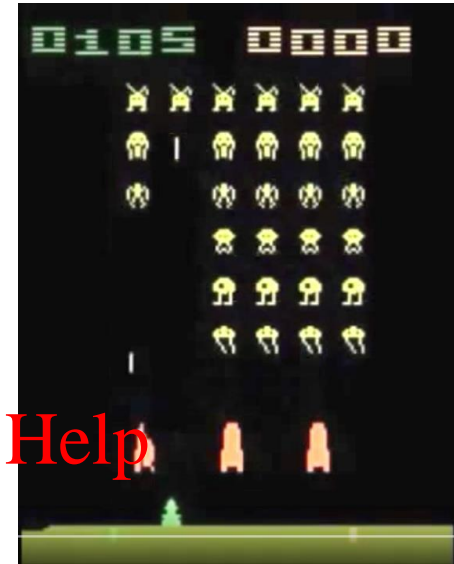
www.nature.com/articles/nature14236

DQN demo:

<https://www.youtube.com/watch?v=qXKQf2BOSE>

DQN source code:

www.sites.google.com/a/deepmind.com/dqn/



Assignment Project Exam Help

Downsides of RL

Add WeChat powcoder

- RL is less sampling efficient than supervised learning because it involves bootstrapping, which uses an estimate of the Q-value to update the Q-value predictor
- Rewards are usually sparse and learning requires to reach the goal by chance
- Therefore, RL might not find a solution at all if the state space is large or if the task is difficult

Assignment Project Exam Help

Summary

Add WeChat powcoder

- The goal of Reinforcement learning:
 - learn to predict actions that maximize future rewards
- Markov Decision Process
 - Formalizes the RL framework as
 - $MDP = (S, A, R, P_{sa}, \gamma)$
- Approaches to reinforcement learning:
 - Learn value function (offline)
 - Learn optimal policy (offline)
 - Learn Q-function (online)

Assignment Project Exam Help

References

Andrew Ng's Reinforcement Learning course, lecture 16

<https://www.youtube.com/watch?v=Rtxl449ZjSc>

Andrej Karpathy's blog post on policy gradient

<http://karpathy.github.io/2016/05/31/rl/>

Mnih et. al, Playing Atari with Deep Reinforcement Learning (DeepMind)

https://www.cs.toronto.edu/~vmnih/vl3m/atari_deepqn.pdf

Intuitive explanation of deep Q-learning

<https://www.nervanasys.com/demystifying-deep-reinforcement-learning/>

Assignment Project Exam Help

Next Class

Add WeChat powcoder

Reinforcement Learning II

Q-learning cont'd; deep Q-learning (DQN)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder