

# Assignment Project Exam Help

## Announcements

Add WeChat powcoder

**Reminder:** pset5 self-grading form and pset6 out Thursday, due 11/19 (1 week)

Assignment Project Exam Help

- No lab this week!
- <https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

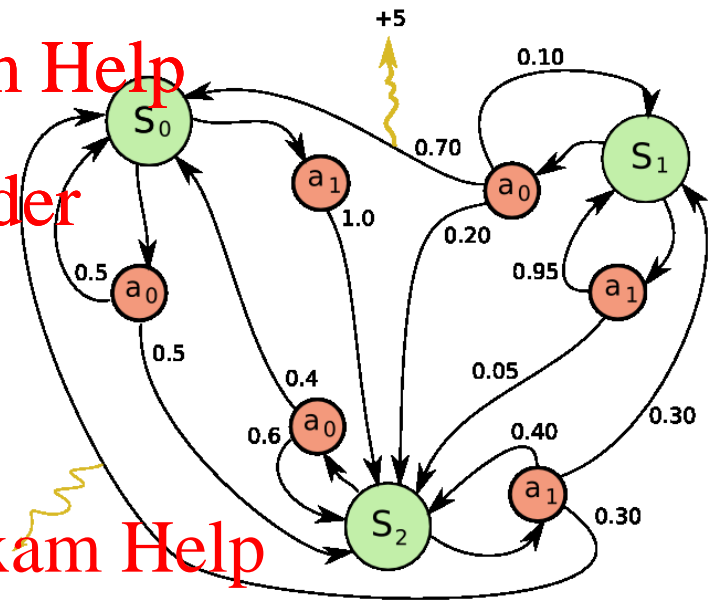
Add WeChat powcoder

# Reinforcement Learning II

# Assignment Project Exam Help

## Recall: MDP notation

Add WeChat powcoder



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

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



# MDP (Simple example)


Assignment Project Exam Help  
Add WeChat powcoder



# 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  $R: S \rightarrow \mathbb{R}$
- 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

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  $\pi$
  - ▶ with discount factor  $\gamma$

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

- ▶ 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-table, 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

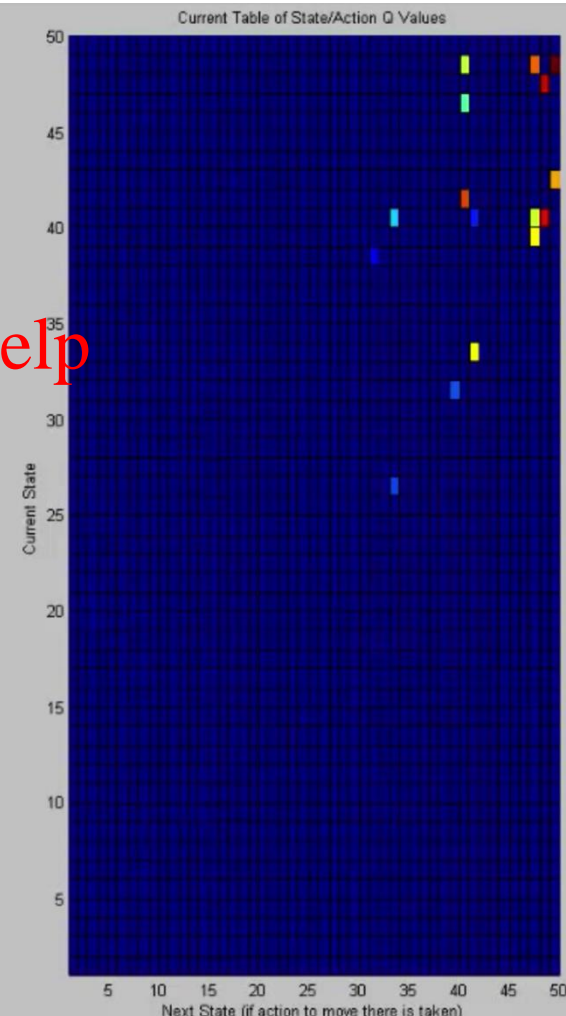
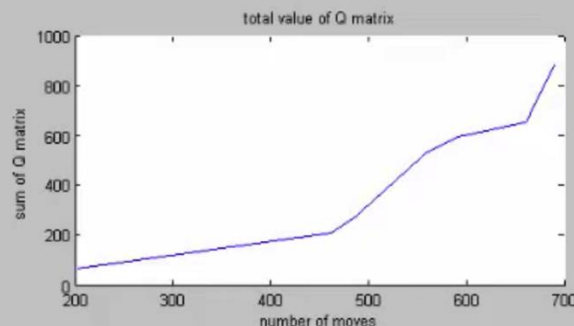
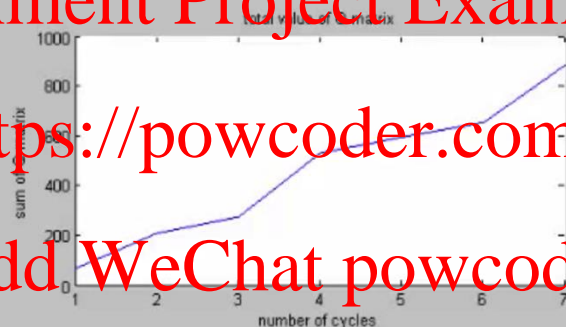
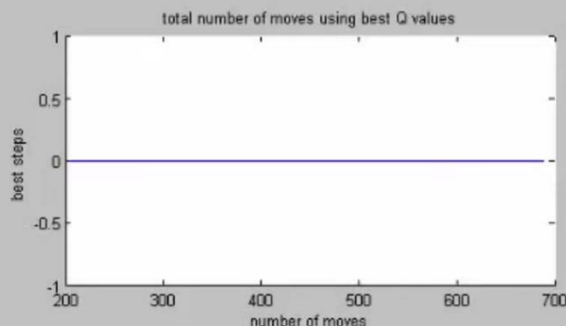
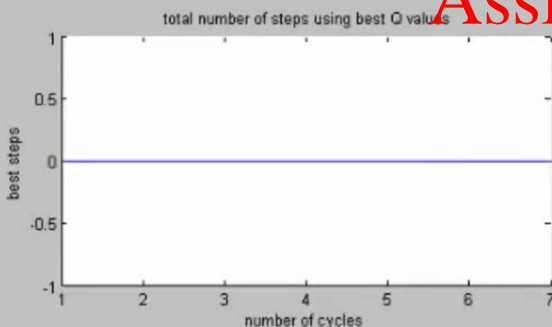
discount

largest increase over all  
possible actions in new state



# Q-learning example

Goal: get from bottom left to top right



<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  $\epsilon$ -greedy approach: at each step with small probability  $\epsilon$ , 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  $\epsilon$  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 ( $p(s, a, s') \in \mathbb{R}$ )

<https://powcoder.com>

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

# Assignment Project Exam Help

## Deep RL

Add WeChat powcoder

- $V$ ,  $Q$  or  $\pi$  can be approximated with deep network

- Deep Q-Learning

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

- Alternative: learn a Policy Network

- Input: state
- Output: distribution over actions

# Assignment Project Exam Help

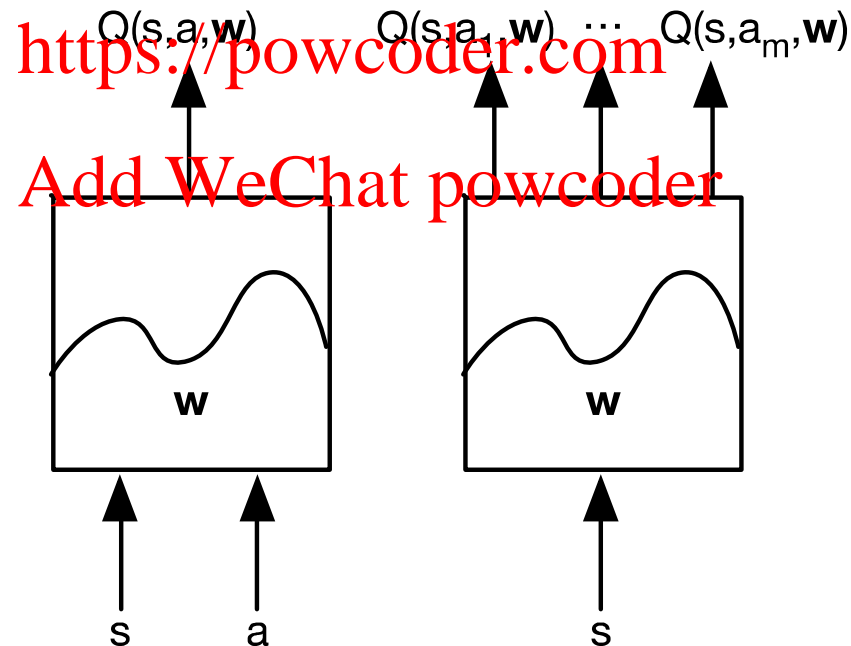
## Q-value network

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

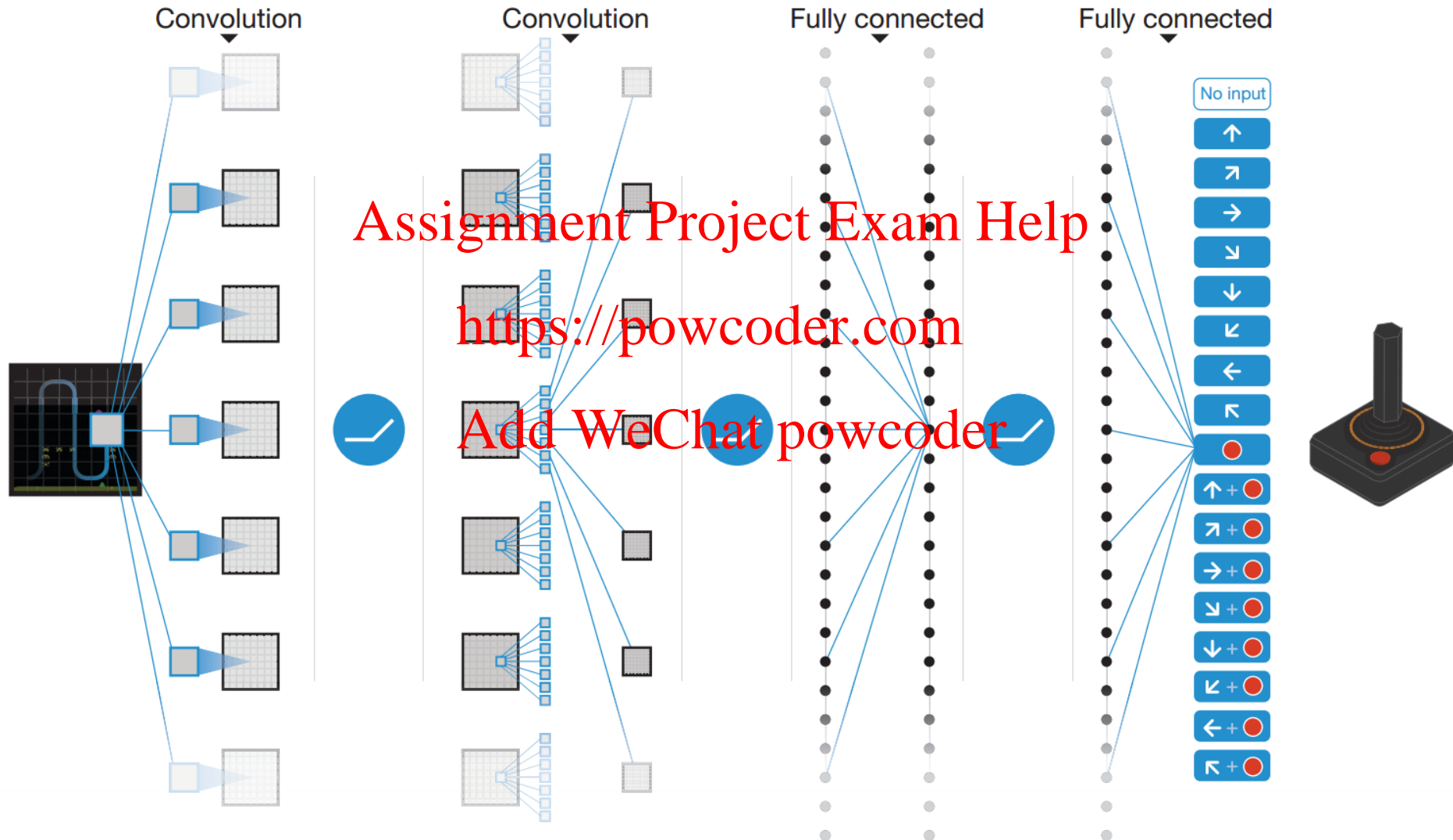
$$l = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

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

# Assignment Project Exam Help

## DQN - Playing Atari

Add WeChat powcoder

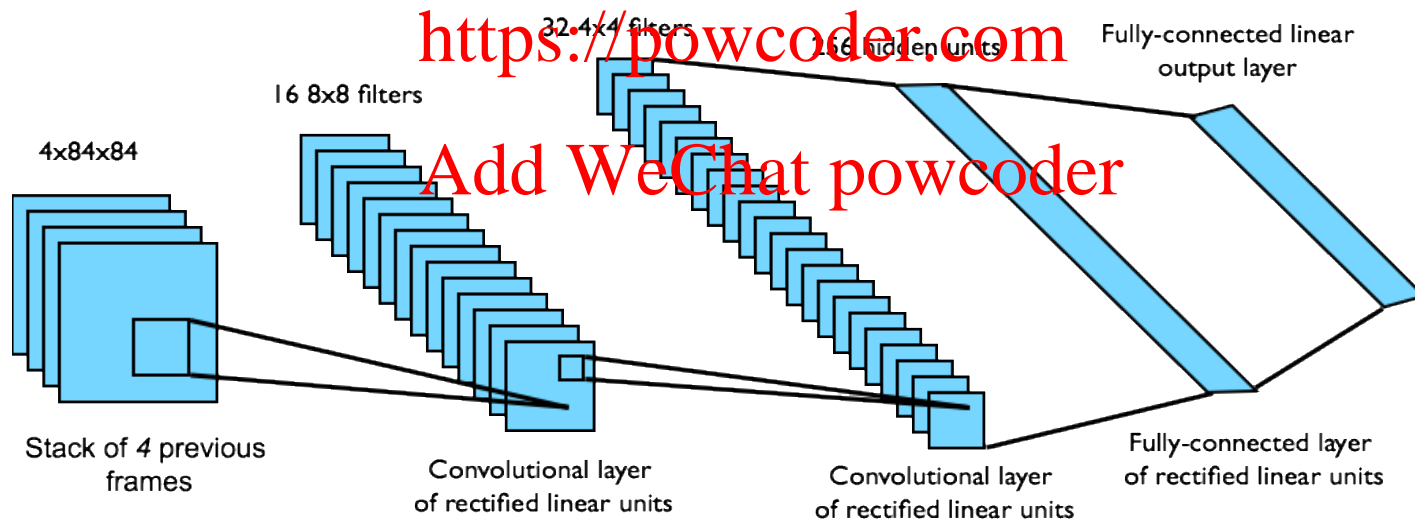




# DQN - Playing Atari

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

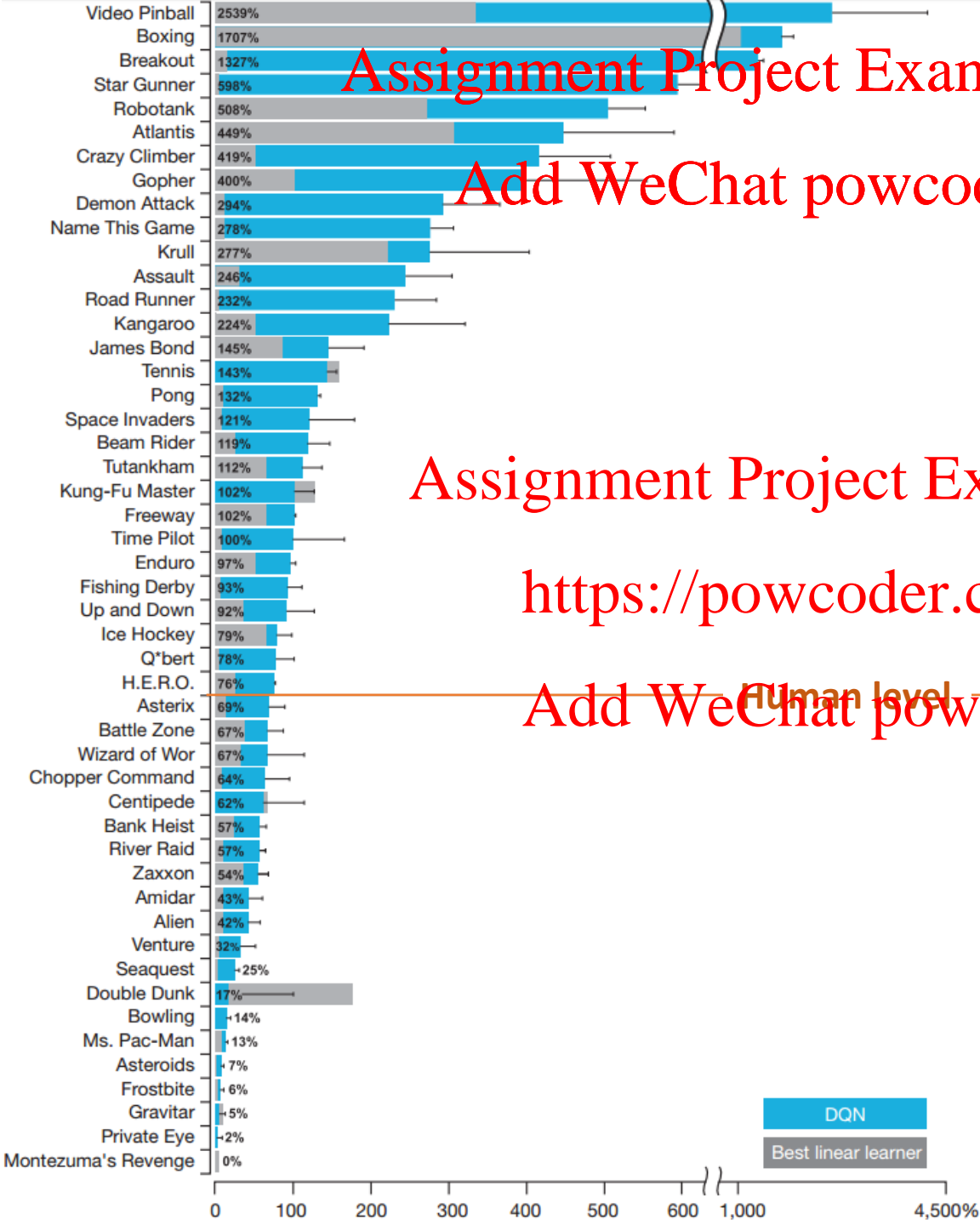
---

## 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  $\epsilon$  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

Best linear learner



# DQN for Atari

Assignment Project Exam Help  
Add WeChat powcoder

DQN paper:

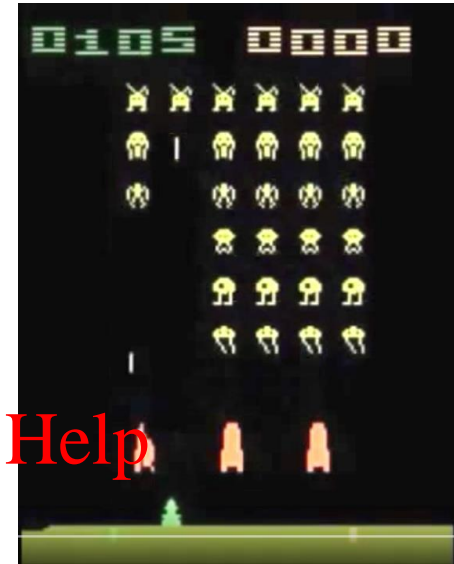
[www.nature.com/articles/nature14236](http://www.nature.com/articles/nature14236)

DQN demo:

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

DQN source code:

[www.sites.google.com/a/deepmind.com/dqn/](http://www.sites.google.com/a/deepmind.com/dqn/)



# Assignment Project Exam Help

## Deep RL

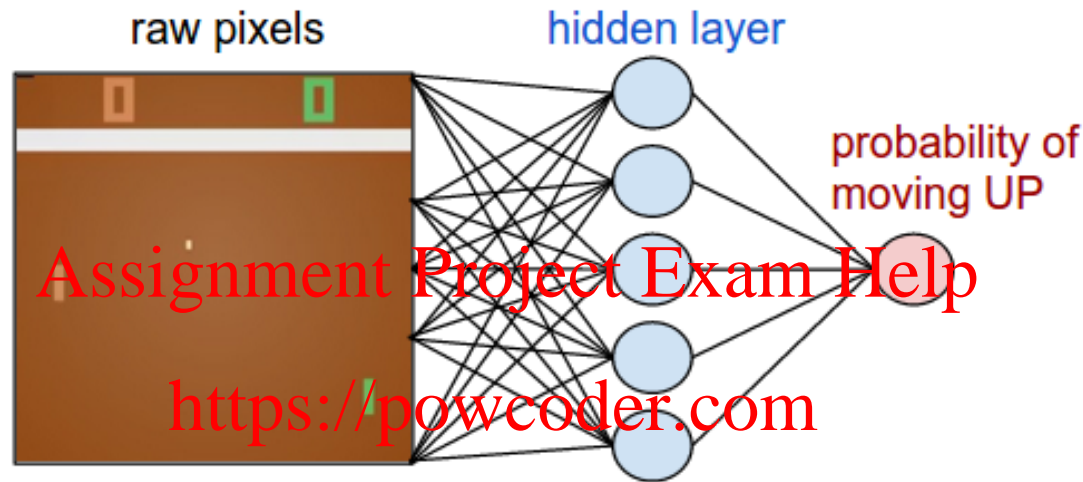
Add WeChat powcoder

- $V$ ,  $Q$  or  $\pi$  can be approximated with deep network
- Deep Q-Learning
  - Input: state, action
  - Output: Q-value
- Alternative: learn a Policy Network
  - Input: state
  - Output: distribution over actions

# Assignment Project Exam Help

## Policy network for pong

Add WeChat powcoder



- define a *policy network* that implements the player
- takes the state of the game and decides what to do (move UP or DOWN)
- 2-layer neural network that takes the raw image pixels\* (100,800 = 210x160x3), outputs the probability of going UP

\*feed at least 2 frames to the policy network so that it can detect motion.

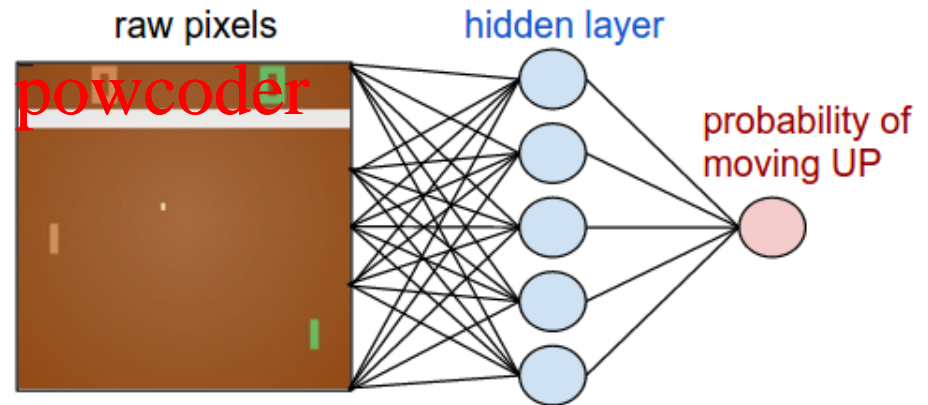
<http://karpathy.github.io/2016/05/31/rl/>



# Assignment Project Exam Help

## Policy gradient

Add WeChat powcoder



- Suppose network predicts

$$p(\text{UP}) = 30\%$$

$$p(\text{DOWN}) = 70\%$$

Assignment Project Exam Help

- Can sample an action from this distribution and execute it
- Can immediately use gradient of 1.0 for DOWN and backprop to find the gradient vector that would encourage the network to predict DOWN

<https://powcoder.com>

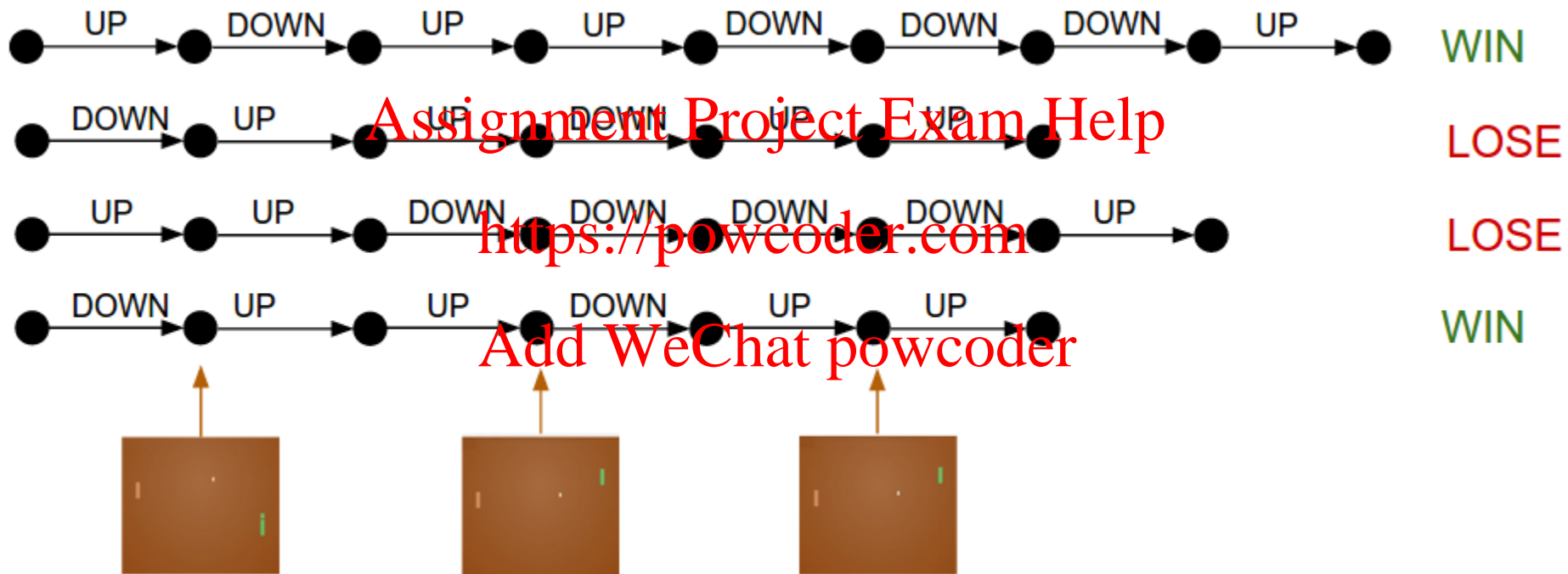
Add WeChat powcoder

**Problem:** do not yet know if going DOWN is good!

**Solution:** simply wait until the end of the game, then take the reward we get (either +1 if we won or -1 if we lost), and enter that as the gradient for taken actions

# Assignment Project Exam Help

## Policy gradient



# Assignment Project Exam Help

## Problems with this?

Add WeChat powcoder

- what if we made a good action in frame 50 (bouncing the ball back correctly), but then missed the ball in frame 150?

Assignment Project Exam Help

- If every single action is now labeled as bad (because we lost), wouldn't that discourage the correct bounce on frame 50?

<https://powcoder.com>

Add WeChat powcoder

- Yes, but after thousands/millions of games, network will learn a good policy



# Assignment Project Exam Help

## Policy gradient

Add WeChat powcoder

Want to maximize

$$E_{x \sim p(x|\theta)}[f(x)]$$

Assignment Project Exam Help

<https://powcoder.com>

$f(x)$  is the reward function

$p(x)$  is the policy network with parameters  $\theta$

Add WeChat powcoder

(i.e. change the network's parameters so that action samples get higher rewards)

# 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

## 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/docs/dqn.pdf>

Intuitive explanation of deep Q-learning

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

Assignment Project Exam Help

Next Class

Add WeChat powcoder

## Unsupervised Learning III: Anomaly Detection

Anomaly detection methods: density estimation, reconstruction based method, One Class SVM; evaluating anomaly detection

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder