

# Announcements

**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



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Reinforcement Learning

# Deep Mind's bot playing Atari Breakout

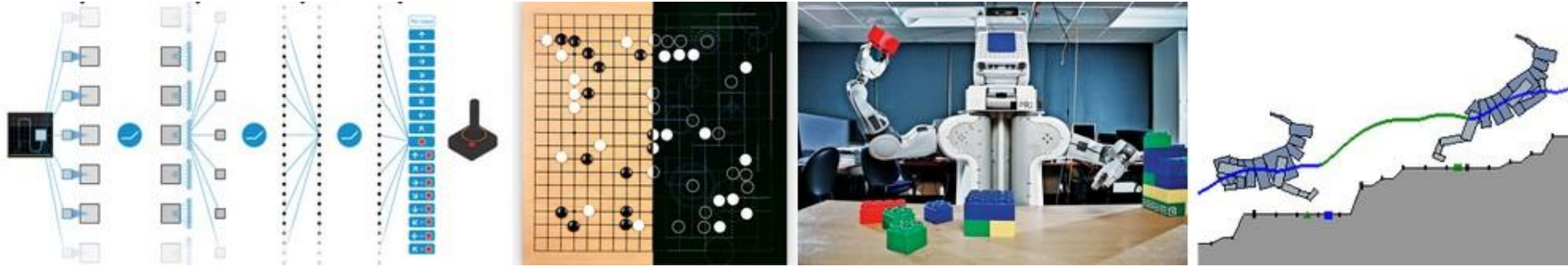
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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





# Reinforcement Learning

Assignment Project Exam Help

- 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

<https://powcoder.com>

Add WeChat powcoder

# What is reinforcement learning?

Reinforcement Learning



# Types of learning



Supervised



Unsupervised



Reinforcement

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Supervised learning

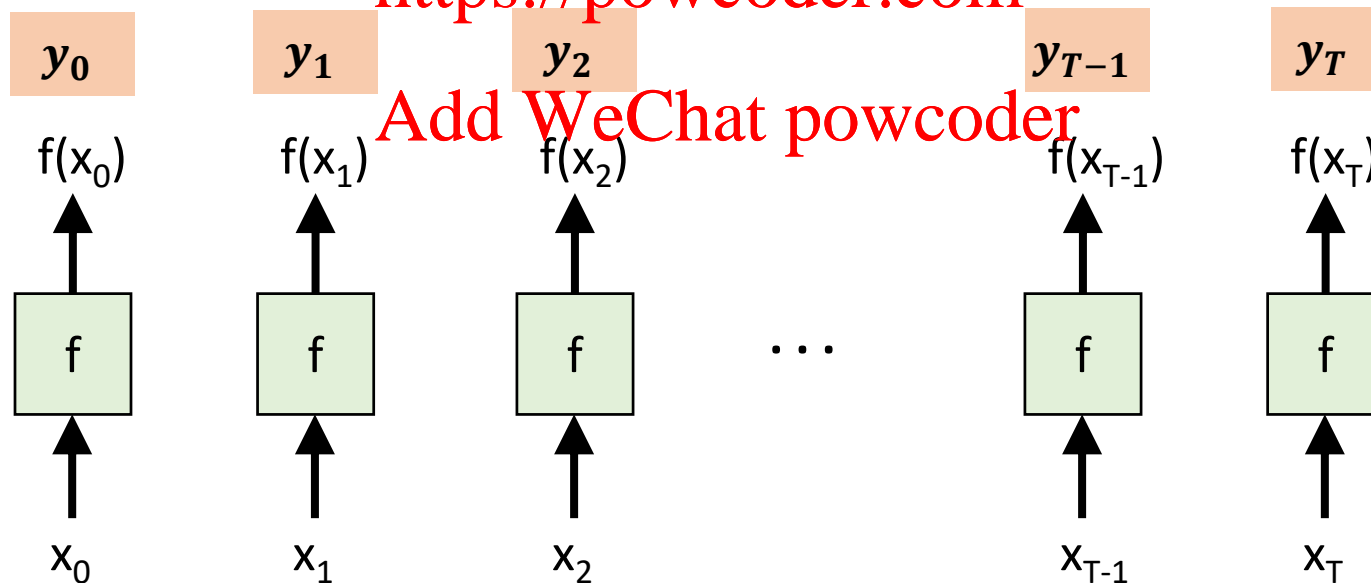
- model  $f$  receives input  $x$
- also gets correct output  $y$
- predictions do not change future inputs

Assignment Project Exam Help

**Supervised learning:** (in arbitrary order of examples)

<https://powcoder.com>

Add WeChat powcoder



# This is not how humans learn!

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder





# Reinforcement learning

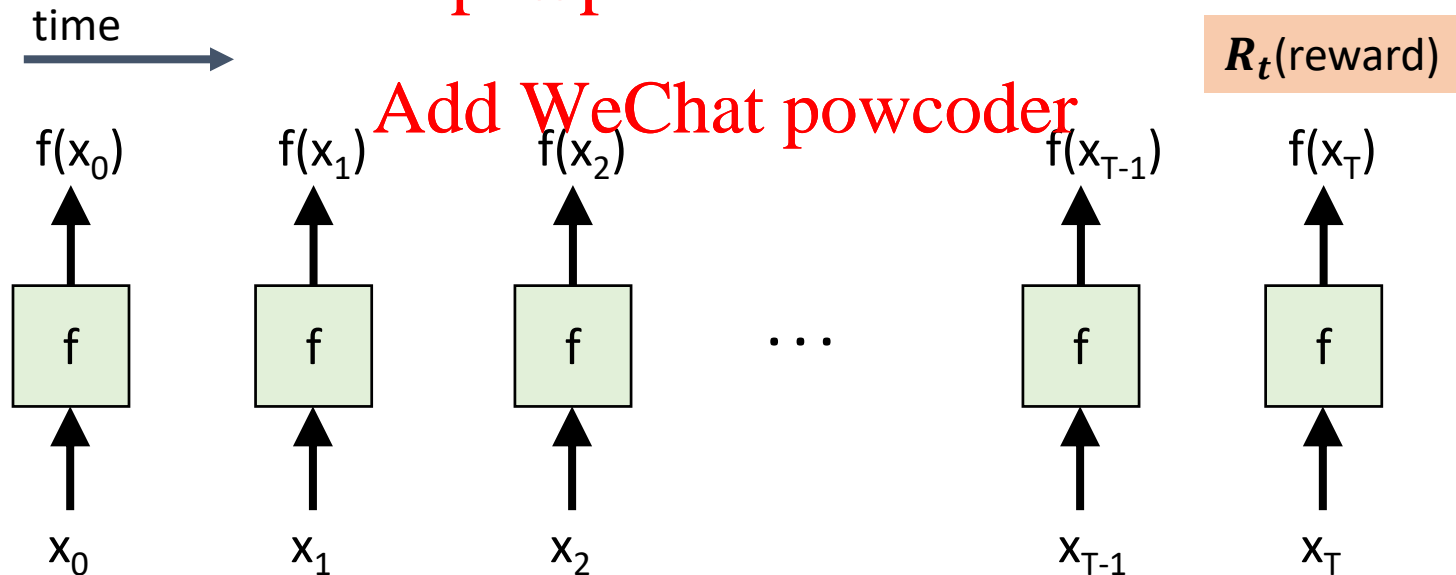
- agent receives input  $x$ , chooses action
- gets  $R$  (reward) after  $T$  time steps
- actions affect the next input (state)

Assignment Project Exam Help

**Reinforcement learning:**

<https://powcoder.com>

Add WeChat powcoder



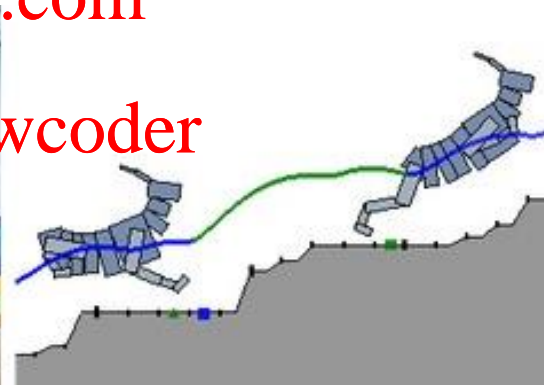
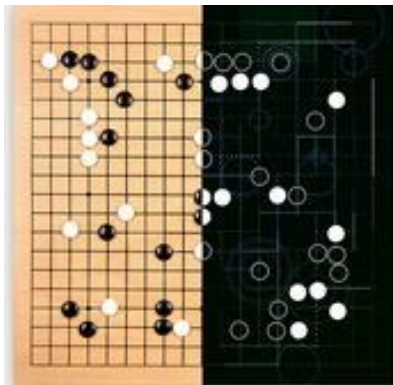
# Input is the “world’s” state

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Output is an action

- 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



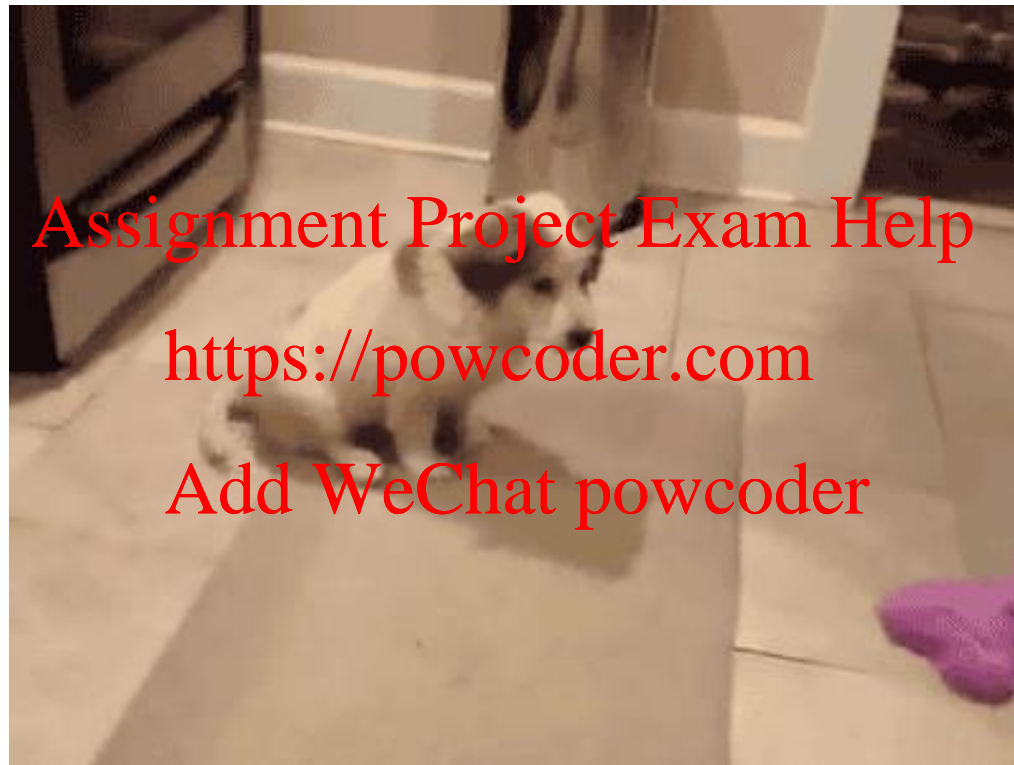
Actions affect state!



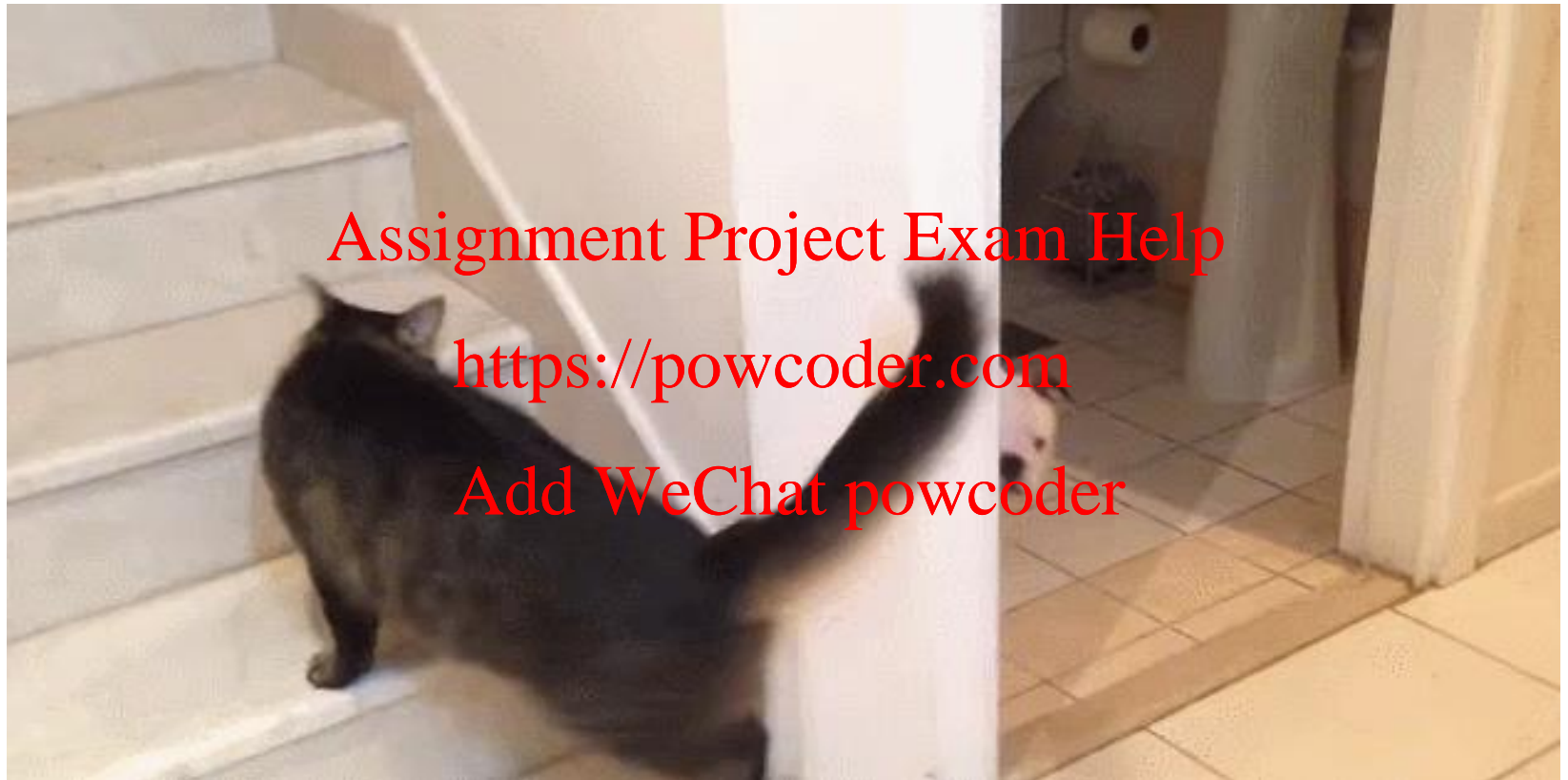
action→reward



# Only some actions lead to rewards



# Some rewards are negative





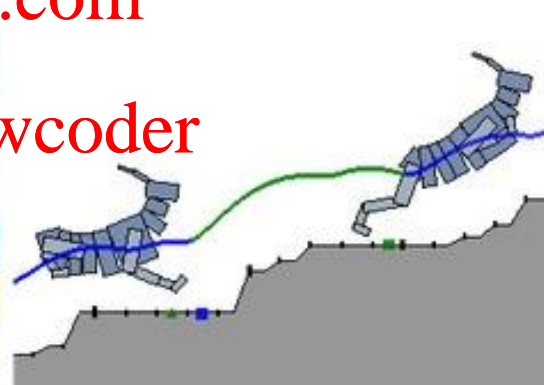
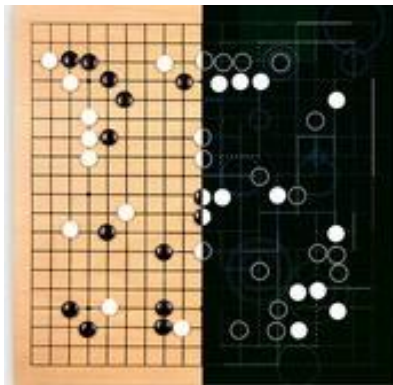
# Reward examples

- 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

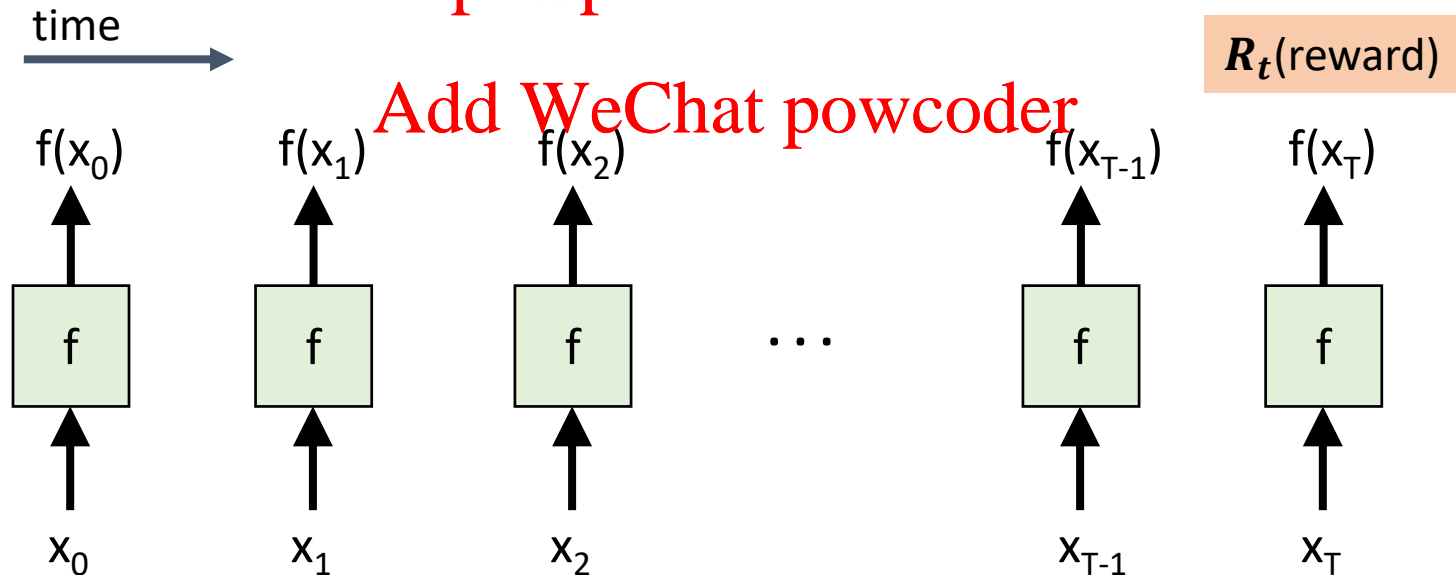
- Learn to predict actions that maximize future rewards
- Need a new mathematical framework

Assignment Project Exam Help

**Reinforcement learning:**

<https://powcoder.com>

Add WeChat powcoder







Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Markov Decision Process

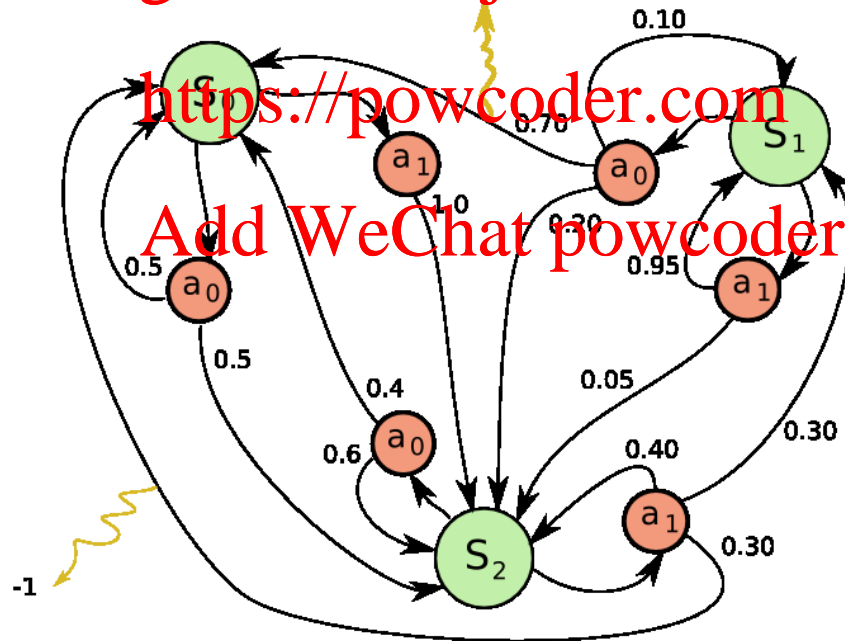
Reinforcement Learning



# Markov Decision Process (MPD)

**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

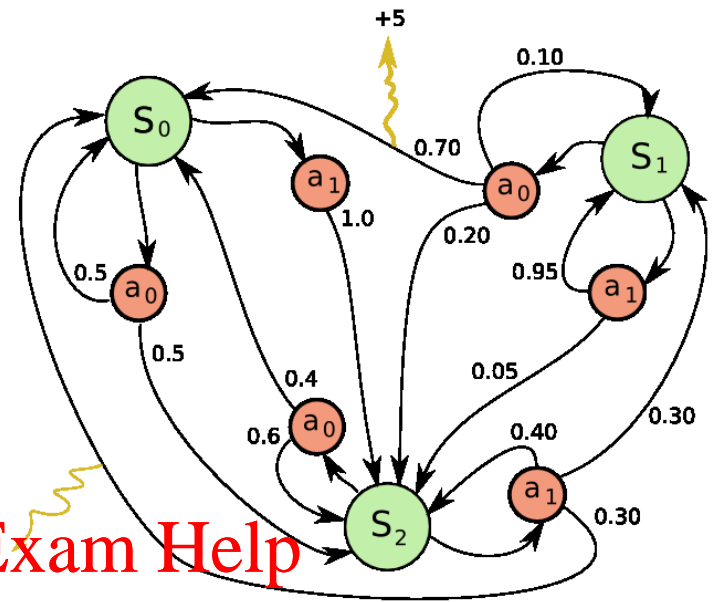


[https://en.wikipedia.org/wiki/Markov\\_decision\\_process](https://en.wikipedia.org/wiki/Markov_decision_process)

# MDP notation

- $S$  – 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

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



# Discount factor $\gamma$

- 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  $\gamma$  value controls which the path the agent should prefer



# MDP (Simple example)



# MDP (Simple example)

	1	2	3	4
1				
2				
3				

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# MDP (Simple example)

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

	1	2	3	4	
1					✓
2					✗
3					



# 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



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# 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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# 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

$$\begin{aligned}
 P_{(3,3),\uparrow}((2,3)) &= 0.8 \\
 P_{(3,3),\uparrow}((3,4)) &= 0.1 \\
 P_{(3,3),\uparrow}((3,2)) &= 0.1 \\
 P_{(3,3),\uparrow}((1,3)) &= 0 \\
 &\vdots
 \end{aligned}$$



# MDP - Dynamics

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

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

- Continue... <https://powcoder.com>

-0.02

-0.02

-0.02

Add WeChat powcoder

- Total payoff:

-0.02

-0.02

-0.02

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

# How do we choose good actions?



# Choosing actions in MDP

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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

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

- In our example:

$R$  – get to charge station, avoid stairs

$\gamma$  – discourage long paths, how much to delay reward



# MDP – Policy $\pi$

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	→	←	↑	↓

Assignment Project Exam Help

<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

<https://powcoder.com>

Add WeChat powcoder

# Policy Value and Q functions

Reinforcement Learning



# MDP – Policy value function

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	→	←	↑	↓

- Value function for policy  $\pi: S \rightarrow \mathbb{R}$

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

Expected sum of discounted rewards



# MDP – Policy value function

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

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

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

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

# MDP – Policy value function

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

Assignment Project Exam Help

*solve*  
*|S|*  
eqs.  $\left\{ \begin{array}{l} V_1 = R(1) + \gamma P_{1,\downarrow}(2) V_2 + \dots \\ V_2 = \dots \\ \dots \\ V_{10} = R(10) + \gamma (P_{10,\uparrow}(6) V_6 + P_{10,\uparrow}(9) V_9 + P_{10,\uparrow}(11) V_{11}) \\ \dots \end{array} \right.$

<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  $\pi$

# Optimal value function

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

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

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

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

The **optimal policy**  $\pi^*$  is the policy that corresponds to the optimal value function

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

# Value function vs. Q-function

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

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

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



# Optimal Q-function

If we know the optimal Q-function  $Q^*(s, a)$ , the optimal policy  $\pi^*$  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

<https://powcoder.com>

Add WeChat powcoder

# RL Approaches

Reinforcement Learning



# Reinforcement learning approaches

## 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

# Value Iteration Algorithm

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

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

```
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
```



# Policy Iteration Algorithm

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

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

```
Initialize a policy  $\pi'$  arbitrarily
Repeat
     $\pi \leftarrow \pi'$ 
    Compute the values using  $\pi$  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  $\pi = \pi'$ 
```

# Reinforcement learning approaches

Optimal value function

need  $P_{sa}$

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

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

Assignment Project Exam Help

<https://powcoder.com>

Optimal policy

need  $\pi, P_{sa}$

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

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

<https://powcoder.com>

Add WeChat powcoder

# Q-Learning (discrete)

Reinforcement Learning

# Q-value function

- ▶ 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

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



# Optimal Q-value function

- ▶ 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.

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

- ▶ Optimal value maximises over all decisions. Informally:

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

- ▶ Formally, optimal values decompose into a Bellman equation

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

# Q-learning algorithm

The agent interacts with the environment, updates Q recursively

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

```
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

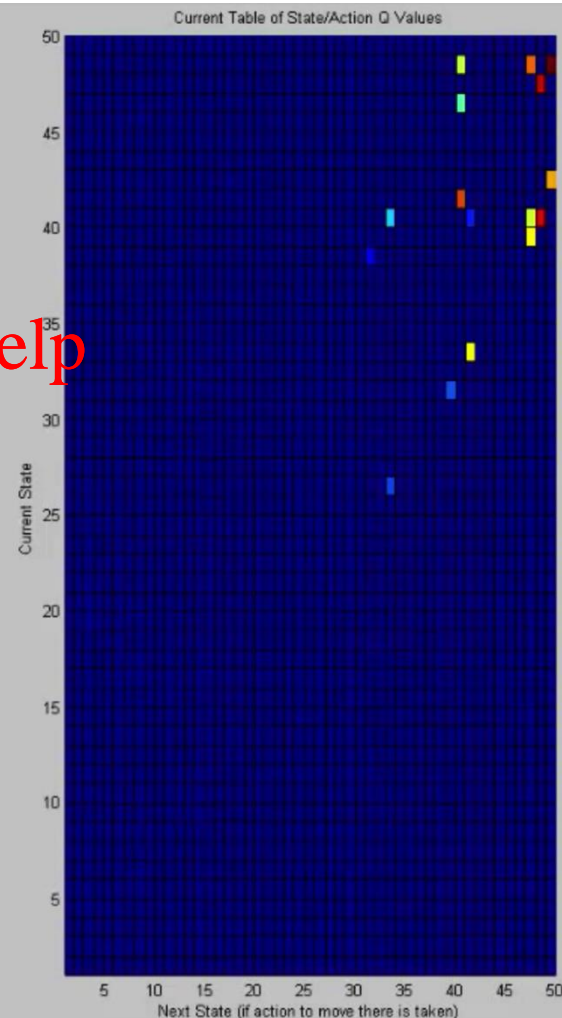
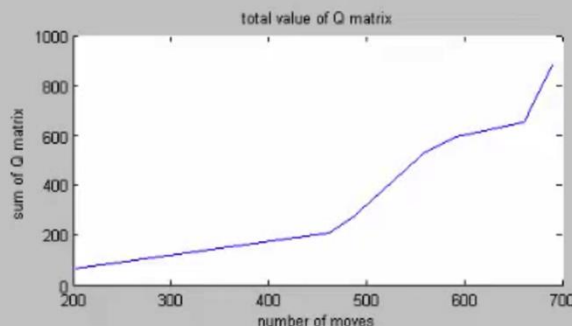
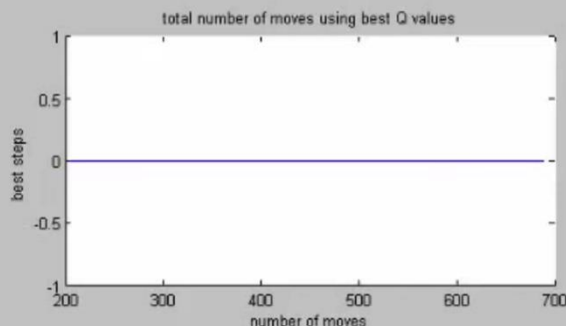
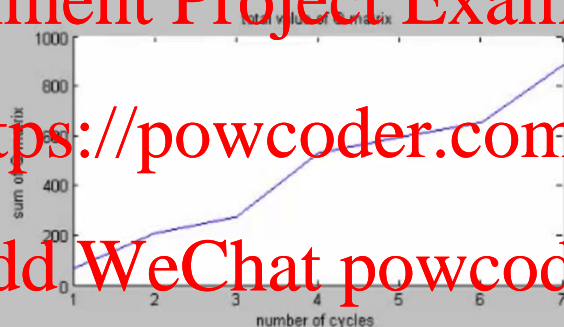
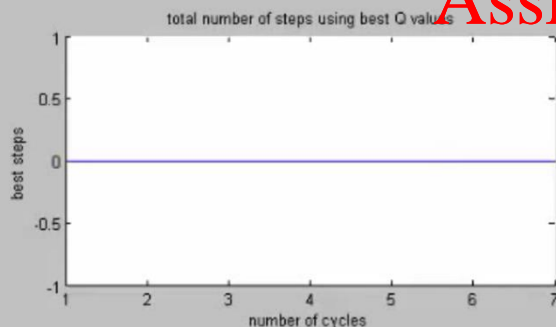
Goal: get from bottom left to top right



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



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

# Exploration vs exploitation

- 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

<https://powcoder.com>

Add WeChat powcoder

# Continuous state

Reinforcement Learning



# Continuous state - Pong

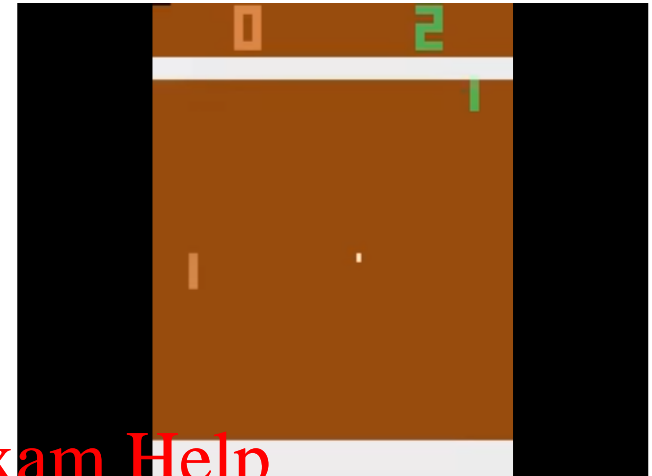
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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

# MDP for Pong



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

<https://powcoder.com>

Add WeChat powcoder

# Deep RL

Reinforcement Learning



# Deep RL playing DOTA



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

[https://www.youtube.com/watch?v=eHipy\\_j29Xw](https://www.youtube.com/watch?v=eHipy_j29Xw)

# Deep RL

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

- Deep Q-Learning

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

Assignment Project Exam Help

Cover today

- Alternative: learn a Policy Network

- Input: state
  - Output: distribution over actions

<https://powcoder.com>  
Add WeChat powcoder



# Q-value network

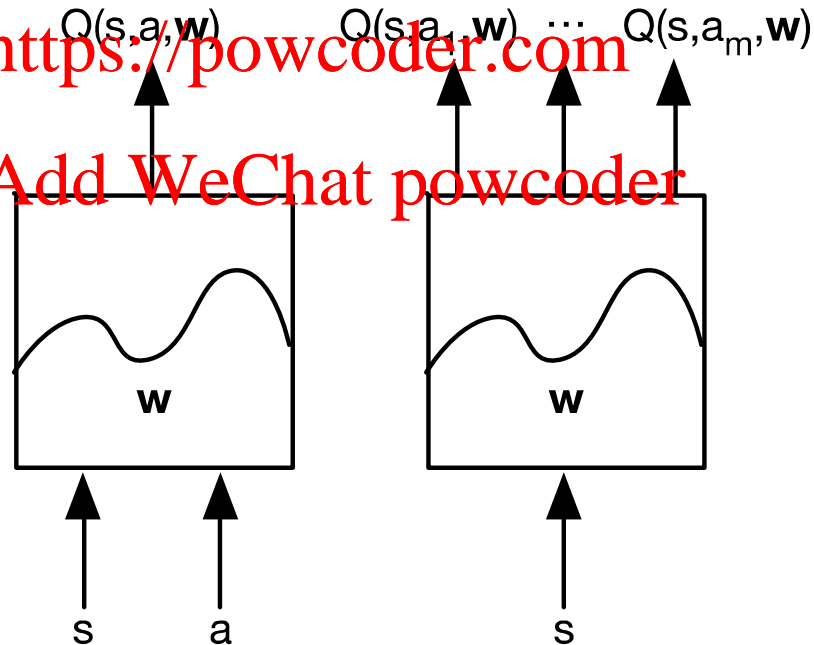
Represent value function by **Q-network** with weights **w**

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Q-value network

- ▶ Optimal Q-values should obey Bellman equation

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

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

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

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

# Deep Q-network (DQN)

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

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

Assignment Project Exam Help

<https://powcoder.com>

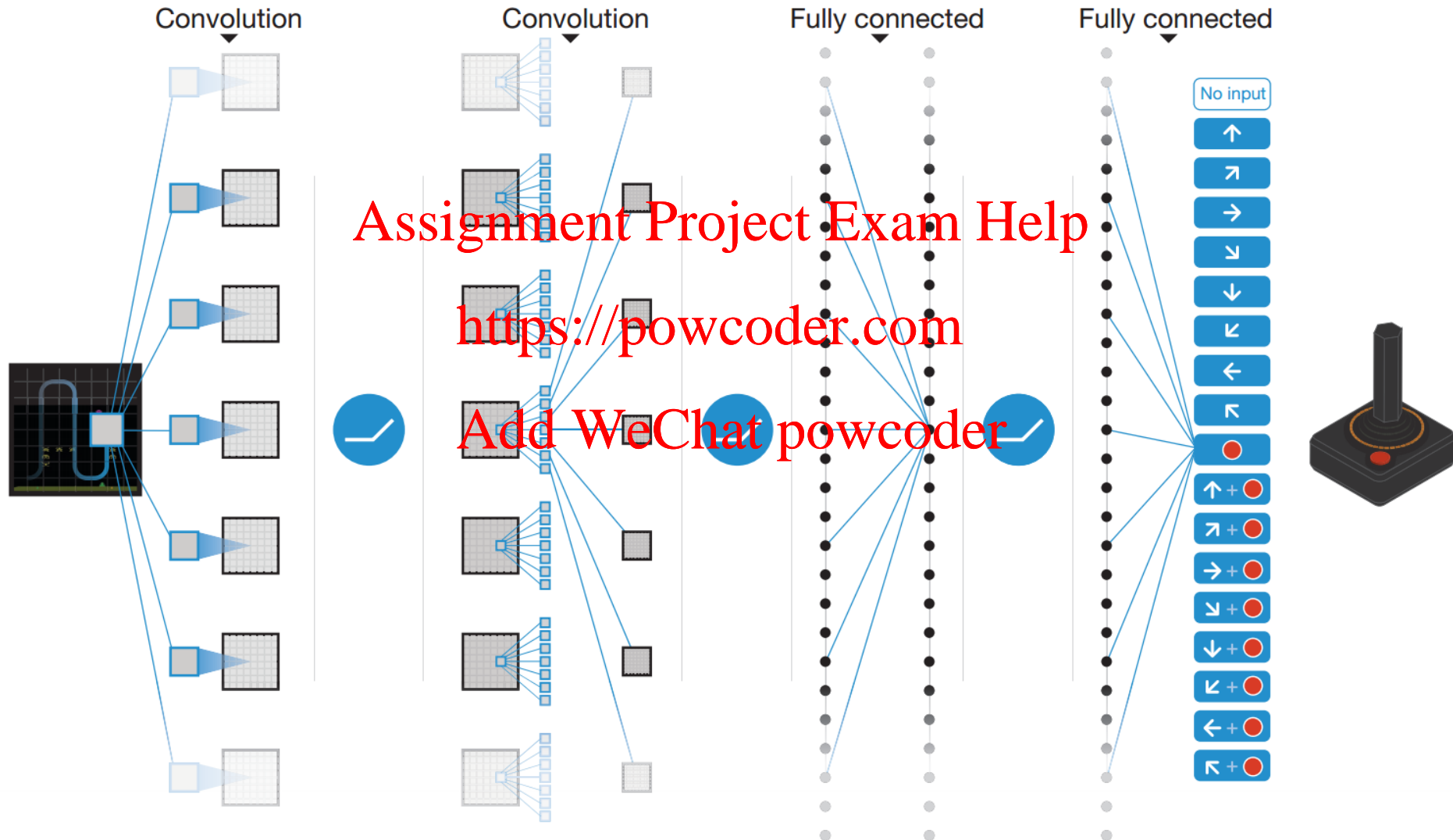
Add WeChat powcoder

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

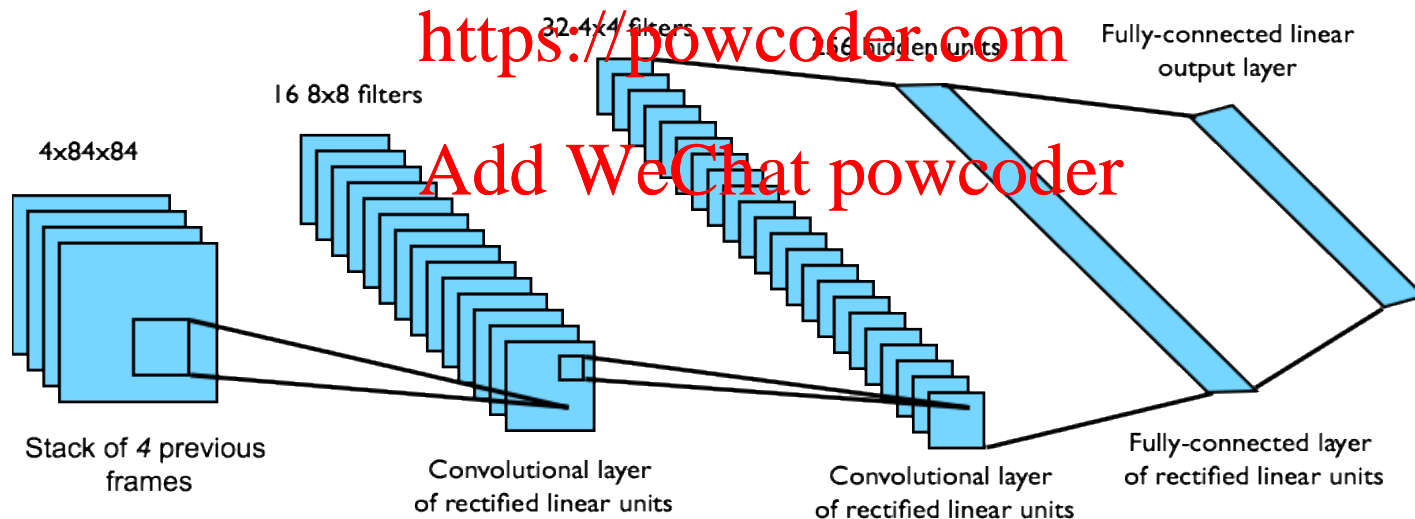
# DQN - Playing Atari





# DQN - Playing Atari

- | 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

---

**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**

    Initialize 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  $\mathcal{D}$

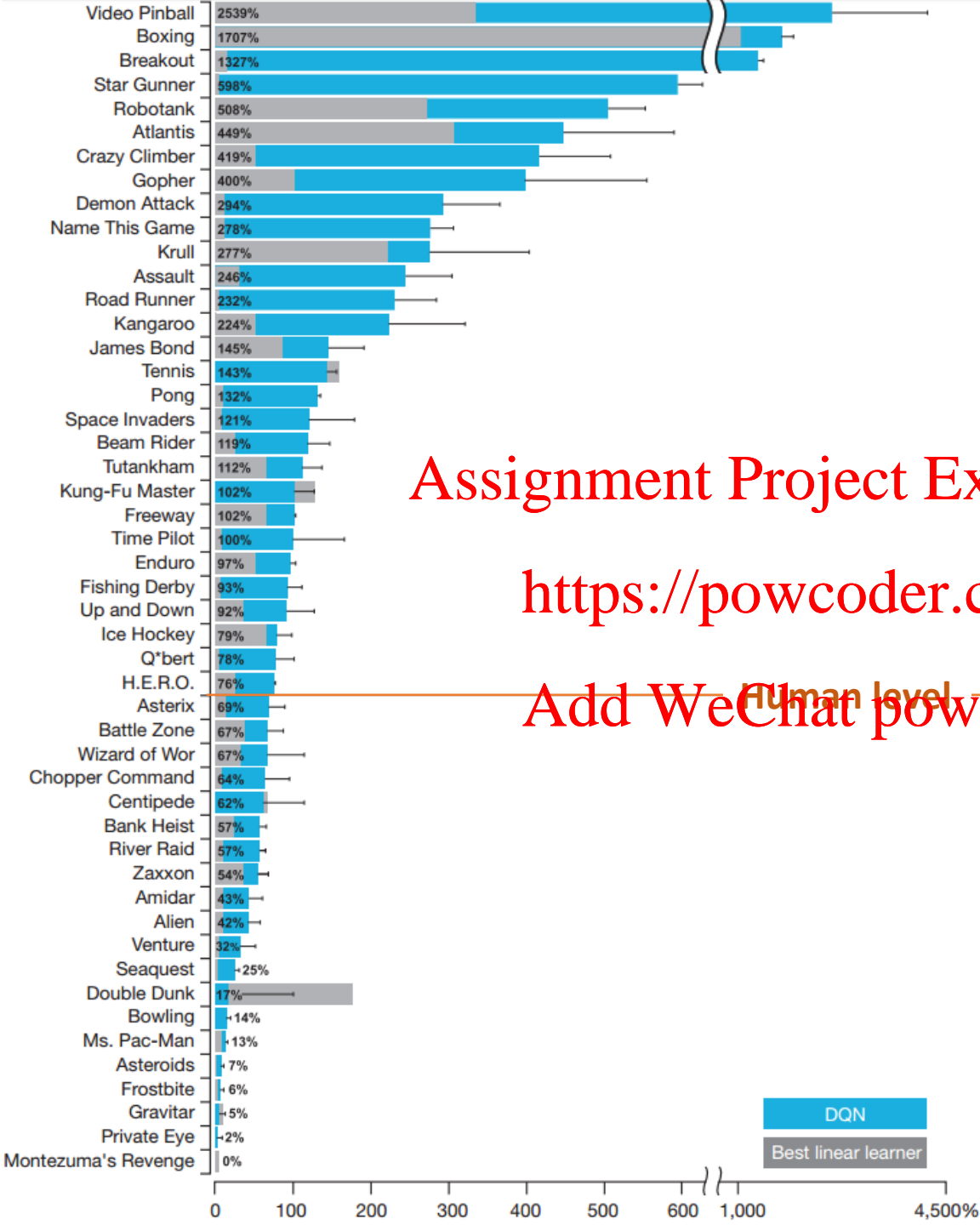
        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        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

<https://powcoder.com>

Add WeChat powcoder

# DQN for Atari

DQN paper:

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

DQN demo:

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

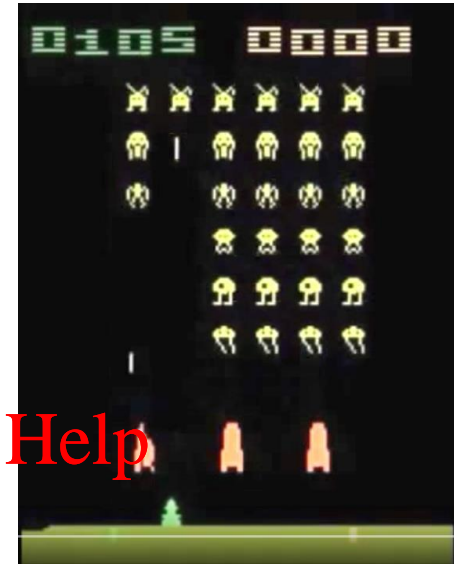
DQN source code:

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder





# Downsides of RL

- 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

# Summary

- 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)

# References

Andrew Ng's Reinforcement Learning course, lecture 16

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

Assignment Project Exam Help

Andrej Karpathy's blog post on policy gradient

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

<https://powcoder.com>

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

<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>

Add WeChat powcoder

Intuitive explanation of deep Q-learning

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

# Next Class

## Reinforcement Learning II

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder