# Agent Based Systems

**Paolo Turrini**

⌂ www.dcs.warwick.ac.uk/~pturrini ✉ p.turrini@warwick.ac.uk

Assignment Project Exam Help

We have seen MDPs and how to calculate the optimal policy (VIA).

However:

https://powcoder.com

- Maybe the state space is too big to do it

- Even if we do know the states we might not know how they are related.

Today we are going to see how to handle these cases, using Reinforcement Learning.

Add WeChat powcoder

What if we don't know what game we are playing?

Play anyway and see what happens! and play as much as possible!

| Game size | Board size N | 3^N | Percent legal | Maximum legal game position ([A094777][20]) |
|---|---|---|---|---|
| 1×1 | 1 | 3 | 33% | 1 |
| 2×2 | 4 | 81 | 70% | 57 |
| 3×3 | 9 | 19,683 | 64% | 12,675 |
| 4×4 | 16 | 43,046,721 | 56% | 24,318,165 |
| 5×5 | 25 | $8.4 \times 10^{11}$ | 49% | $1.1 \times 10^{11}$ |
| 9×9 | 81 | $4.4 \times 10^{38}$ | 23.4% | $1.039 \times 10^{38}$ |
| 13×13 | 169 | $4.3 \times 10^{80}$ | 8.66% | $3.72497923 \times 10^{79}$ |
| 19×19 | 361 | $1.74 \times 10^{172}$ | 1.196% | $2.08168199382 \times 10^{170}$ |
| 21×21 | 441 | $2.57 \times 10^{210}$ | | |

Figure: The complexity of Go

# Neural networks + tree search

Understanding the value of game positions using:

Neural Networks using pattern recognition from a database of previously played games.

Tree Search self-playing (a lot!) and estimating the value of moves.

📕 David Silver et al.
Mastering the game of Go
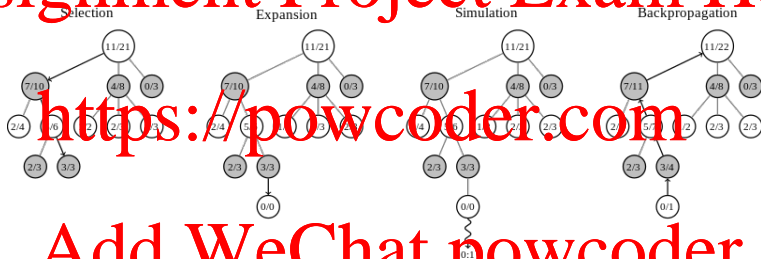with deep neural networks and tree search
Nature, 2016.

📕 David Silver et al.
Mastering the game of Go without human knowledge
Nature, 2017.

> I'm going to only focus on how to infer value
> **without** using pre-processed information

Selection      Expansion      Simulation      Backpropagation

To evaluate intermediate game positions we play a huge number of games from then on.

- Begin at the start state
- The game ends when we reach either goal state $+1$ or $-1$
- Rewards: $+1$ and $-1$ for terminal states respectively, $-0.04$ for all others

Assignment Project Exam Help

https://powcoder.com

0.3

0.1 0.1

Add WeChat powcoder

Stochastic actions (possibly different at each state), four directions

Stochastic actions (possibly different at each state), four directions

Stochastic actions (possibly different at each state), four directions

Stochastic actions (possibly different at each state), four directions

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Stochastic actions (possibly different at each state), four directions

Stochastic actions (possibly different at each state), four directions

Stochastic actions (possibly different at each state), four directions
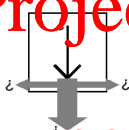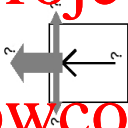
Assignment Project Exam Help

https://powcoder.com
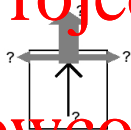
Stochastic actions (possibly different at each state), four directions

Add WeChat powcoder

Assignment Project Exam Help

https://powcoder.com

Stochastic actions (possibly different at each state), four directions

Add WeChat powcoder

Stochastic actions (possibly different at each state), four directions

Assignment Project Exam Help

This is what is known (by the agent) about the environment

https://powcoder.com

- Partially observable (we know where we are, not where we will end up being)
- Markovian (past doesn't matter)
- Stochastic actions (we are not in full control of our choices)
- Discounted rewards (we might be more or less patient)

Add WeChat powcoder

**Passive reinforcement learning**:

- I have a policy
- I don't know the probabilities
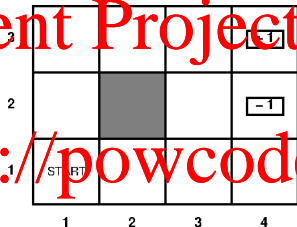- I don't know the values of states
- I don't know the value of actions

**Active reinforcement learning**:

- I don't even have a policy

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Assignment Project Exam Help

https://powcoder.com
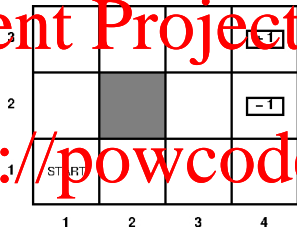


Add WeChat powcoder

- I don't know the values nor the rewards

- I don't know the probabilities

- I'm gonna play anyway

**The plan**: I execute a series of trials until the end states,
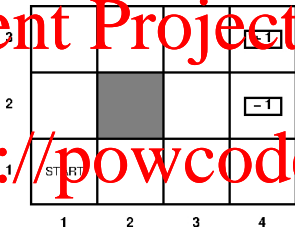just like Monte-Carlo Tree Search!

**Remember**: the expected utility is the expected sum of discounted rewards under the policy

**Assume**: $\gamma = 1$, just to make things simple

Assignment Project Exam Help

https://powcoder.com



Add WeChat powcoder

Suppose I get these trials

$(1,1)_{-0.04} \leadsto (2,1)_{-0.04} \leadsto (3,1)_{-0.04} \leadsto (2,1)_{-0.04} \leadsto (3,1)_{-0.04} \leadsto (3,2)_{-0.04} \leadsto (3,3)_{-0.04} \leadsto (3,4)_{+1}$

$(1,1)_{-0.04} \leadsto (2,1)_{-0.04} \leadsto (3,1)_{-0.04} \leadsto (3,2)_{-0.04} \leadsto (3,3)_{-0.04} \leadsto (2,3)_{-0.04} \leadsto (3,3)_{-0.04} \leadsto (3,4)_{+1}$

$(1,1)_{-0.04} \leadsto (1,2)_{-0.04} \leadsto (1,3)_{-0.04} \leadsto (2,3)_{-0.04} \leadsto (2,4)_{-1}$

**Idea**: Frequency is the key!
Each trial provides a sample of the expected rewards for each state visited.

When a transition occurs from state $s$ to state $s'$ we apply the following update:

$$v^\pi(s) = v^\pi(s) + \alpha(r(s) + \gamma v^\pi(s') - v^\pi(s))$$

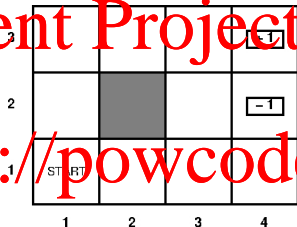where $\alpha \in [0, 1]$ is a confidence parameter: how much we value the new information.

$\alpha$ can be the inverse of the number of times we visited a state: the more we visited, the less we want to learn.

**Notice:** rare transitions? Well, they are rare.
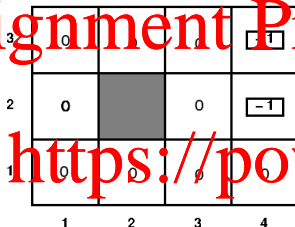
Initialise the values, for

- $\gamma = 1$
- deterministic agent
- $\alpha = \frac{1}{n+1}$ where $n$ is the number of times we visited a state
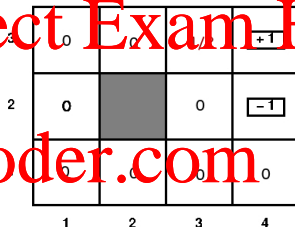- $r = 0$ everywhere but the terminal states

Suppose we can walk (Up, Up, Right, Right, Right)
Apply the update to states, as you walk along:

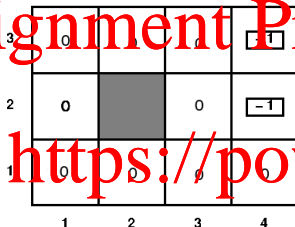$$v^\pi(s) = v^\pi(s) + \alpha(r(s) + \gamma v^\pi(s') - v^\pi(s))$$
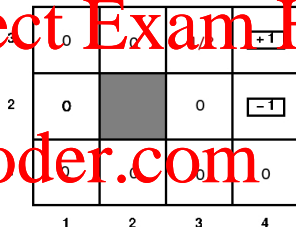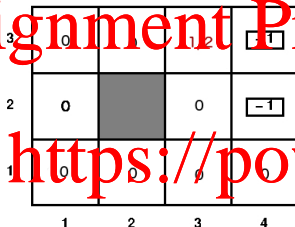
Suppose we can walk (*Up, Up, Right, Right, Right*)
Apply the update to states, as you walk along:

$$v^\pi(s) = v^\pi(s) + \alpha(r(s) + \gamma v^\pi(s') - v^\pi(s))$$

I keep walking the same way...

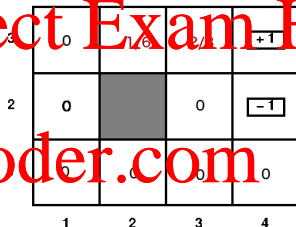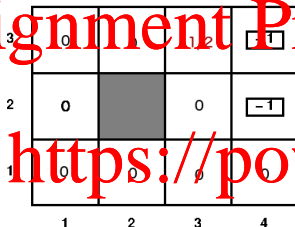$$v^\pi(s) = v^\pi(s) + \alpha(r(s) + \gamma v^\pi(s') - v^\pi(s))$$

Again (*Up*, *Up*, *Right*, *Right*, *Right*)....

$$v^\pi(s) = v^\pi(s) + \alpha(r(s) + \gamma v^\pi(s') - v^\pi(s))$$

Assignment Project Exam Help

https://powcoder.com

- We have a policy which we follow;
- We backpropagate the value with a Bellman-like adjustment;
- We can use a learning rate, depending on our confidence.

Add WeChat powcoder

Now we start without a fixed policy...

What the agent needs to learn is the values of the optimal policy

$$V(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V(s')$$

**Important:** we can't stick to our (locally optimal) habits,
we need to try new stuff!

Exploration vs Exploitation

$Q(s, a)$

the value of performing action $a$ in state $s$

$$Q(s, a) = r(s) + \gamma \sum_{s'} P(s'|(s, a)) \max_{a'} Q(s', a')$$

$$U(s) = \max_a Q(s, a)$$

is the value of performing action $a$ in state $s$
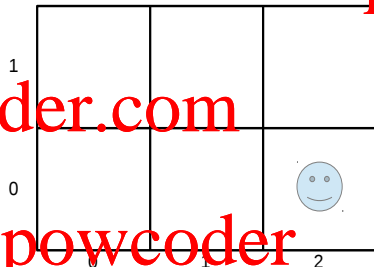
Assignment Project Exam Help

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

https://powcoder.com

It's a temporal difference learning, without fixed policy!

Add WeChat powcoder

- A 2 × 3 grid world
- A pit, an exit and some walls are known in this grid world, but their locations are unknown
- Arrive at the exit: win; fall in the pit: die; hit a wall, suffer
- **Goal:** Get out of this maze (i.e. safely arrive at the exit) as quickly as possible

1

0

0       1       2

- **State**: The agent's current location
- **Action**: *LEFT, RIGHT, UP, DOWN*
- **Environment Dynamics**:
  - Collision results in no movement
  - otherwise, move one square in the intended direction
- **Rewards**:
  - normal move: -1
  - hit a wall: -10
  - die: -100
  - exit: +100
- **Our Goal**: find the best route to exit

- $\alpha = 0.5$; $\gamma = 0.9$
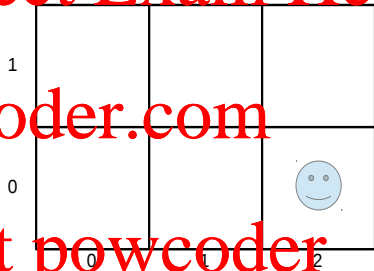- All Q-values are initialised as 0



Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

- $\alpha = 0.5$; $\gamma = 0.9$
- All Q-values are initialised as 0
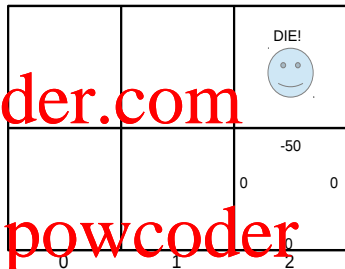- Choose *UP*, and receive -100

- $\alpha = 0.5$; $\gamma = 0.9$
- All Q-values are initialised as 0
- Choose *UP*, and receive -100
- update Q-value:

$$Q([0, 2], UP)$$
$$= (1 - 0.5) \times 0+$$
$$0.5 \times (-100 + 0.9 \times 0)$$
$$= -50$$

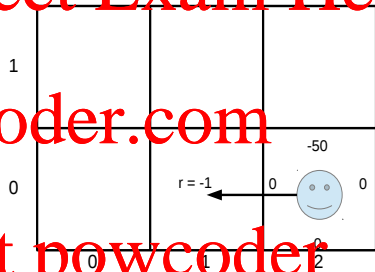- $\alpha = 0.5$, $\gamma = 0.9$
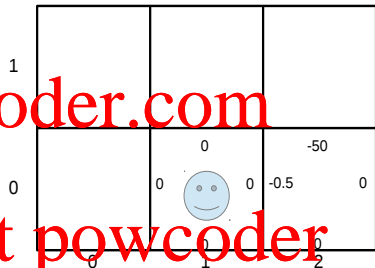- Choose *LEFT*, and receive -1

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *LEFT*, and receive -1
- update Q-value:
  $$Q(0,2),LEFT)$$
  $$= (1 - 0.5) \times 0+$$
  $$0.5 \times (-1 + 0.9 \times 0)$$
  $$= -0.5$$

Assignment Project Exam Help

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *UP*, and receive -10

https://powcoder.com

Add WeChat powcoder

Assignment Project Exam Help

- $\alpha = 0.5$, $\gamma = 0.9$

- Choose *UP*, and receive -10

- update Q-value:

https://powcoder.com

$$= (1 - 0.5) \times 0 +$$
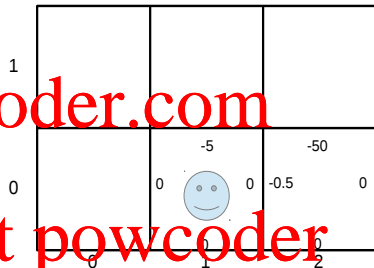$$0.5 \times (-10 + 0.9 \times 0)$$

Add WeChat powcoder

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *LEFT*, and receive -1

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *LEFT*, and receive -1
- update Q-value:
  $Q(0,1), LEFT)$
  $= (1 - 0.5) \times 0+$
  $\quad 0.5 \times (-1 + 0.9 \times 0)$
  $= -0.5$

Assignment Project Exam Help

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *UP*, and receive -1

https://powcoder.com

Add WeChat powcoder

Assignment Project Exam Help

- $\alpha = 0.5$, $\gamma = 0.9$

- Choose *UP*, and receive -1

- update Q-value:
  $$Q(0,0,UP)$$
  $$= (1 - 0.5) \times 0 +$$
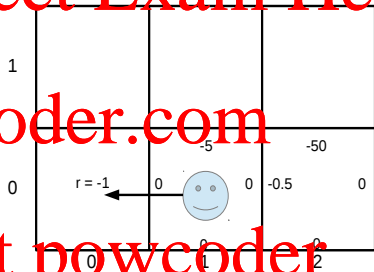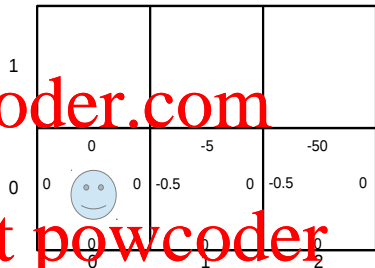  $$0.5 \times (-1 + 0.9 \times 0)$$
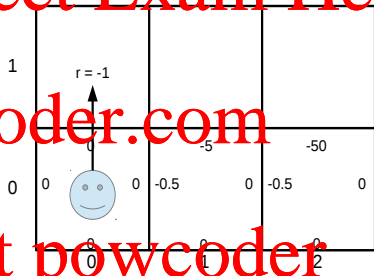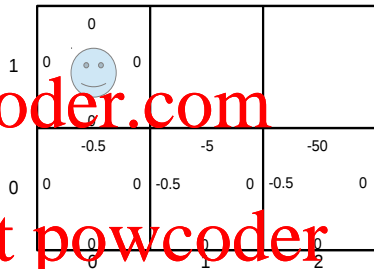  $$= -0.5$$

https://powcoder.com

Add WeChat powcoder

- $\alpha = 0.5$, $\gamma = 0.9$
  - Choose *RIGHT*, and receive 100

- $\alpha = 0.5$, $\gamma = 0.9$
- Choose *RIGHT*, and receive 100
- update Q-value:

$$Q([0,1], \textit{RIGHT})$$
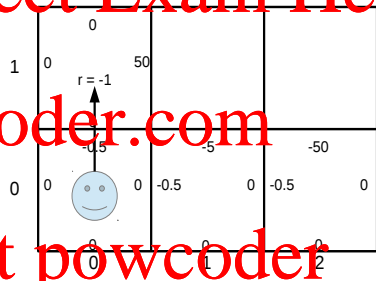$$= (1 - 0.5) \times 0 +$$
$$0.5 \times (100 + 0.9 \times 0)$$
$$= 50$$

- $\alpha = 0.5$, $\gamma = 0.9$
- The next time agent visits [0,0] and performs *UP*:

  $Q([0,0], UP)$
  $= (1 - 0.5) \times (-0.5) +$
  $0.5 \times (-1 + 0.9 \times 50)$
  $= 21.75$

Assignment Project Exam Help

- $\alpha = 0.5$, $\gamma = 0.9$
- The next time agent visits [0,0] and performs $UP$:

https://powcoder.com

$Q([0,0], UP)$
$= (1 - 0.5) \times (-0.5) +$
$0.5 \times (-1 + 0.9 \times 50)$
$= 21.75$

Add WeChat powcoder

| | 0 | | |
|---|---|---|---|
| 1 | 0 | 50 | |
| | 0 | | |
| | 21.75 | -5 | -50 |
| 0 | 0 | 0 -0.5 | 0 -0.5 0 |
| | 0 | 1 | 2 |

Assignment Project Exam Help

- Quick learning speed.
- Model-free: no need to explicitly compute probabilities, or record trajectory.
- Guarantee to converge.

https://powcoder.com

Add WeChat powcoder

- $\alpha$:
  - Learning step
  - balance between existing experiences (weight: $1 - \alpha$) and new observations (weight: $\alpha$)
- $\gamma$:
  - Future discount
  - balance between current reward (weight: $1$) and next $N$ step's reward (weight: $\gamma^N$)
- $\epsilon$:
  - indicating how 'bold' the agent is
  - balance between **exploitation** (take greedy action, $1 - \epsilon$ chance) and **exploration** (take random action, $\epsilon$ chance)

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Assignment Project Exam Help

- Decision making in sequential environments typical of AI practice
- Optimization techniques (VIA)
- … under incomplete information (Active/Passive Learning by Belmann updates)

https://powcoder.com

**What next?** Population dynamics and multi-agent reinforcement learning.

Add WeChat powcoder