

# Blackjack and Monte Carlo

Using Monte Carlo Methods to Enhance Blackjack Gameplay  
through Reinforcement Learning

Hayk Hovhannisan

BS in Data Science, American University of Armenia

18 May 2023



# Table of contents

1. Blackjack and Monte Carlo

2. Table of contents

  1. Blackjack

3. Environment

4. Monte Carlo Simulation

5. My implementation

  1. Cont'd

6. Results

  1. Optimal Policy

  2. State Score

7. Discussion

8. Thank You For Attention!

# Blackjack

## Rules

- ♣ Objective

- ♦ Get close to 21

- ♥ Values

- ♦ J Q K → 10

- ♥ A → 1 or 11

- ♠ Player 2 card

- ♣ Dealer 1 face up 1 face down

- ♠ No suits





# Environment

♣ Infinite Deck

## Action Space

Hit (1) or Stick (0)

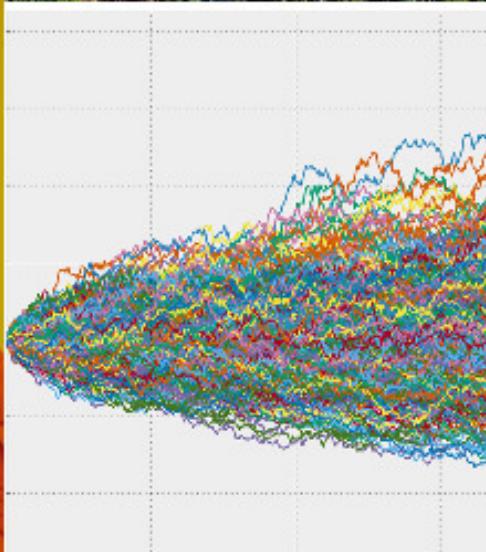
## Observation space

( Players Sum , Dealers Up card, Usable Ace )

( `int` , `int` , `bool` )

## Terminal states

Win (1) Lose (-1) Draw (0)



# Monte Carlo Simulation

♥ Prior knowledge Experience

♦ Transitions Probability Distribution

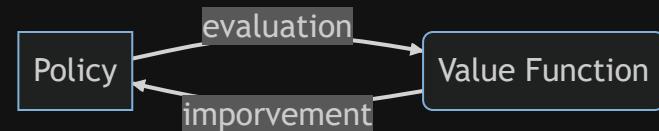
Sample Transitions

♣ Averaging Sampling Results

♠ Discrete finite Episodes → Experience

---

♣ Policy Control



# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop over each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to  $\text{Returns}(S_t, A_t)$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to  $\text{Returns}(S_t, A_t)$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21        eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 #####
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11     for _ in range(num_episodes):
12         episode = generate_episode(Q, eps)
13         G = 0
14         for t in range(len(episode)-1,-1,-1):
15             state, action, reward = episode[t]
16
17             G = gamma*G+reward
18             r[state, action].append(G)
19             Q[state][action] = np.mean(r[state,action])
20
21             eps = max(eps_min, eps*eps_decay)
22
23     return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop ove each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to  $\text{Returns}(S_t, A_t)$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

# My implementation

```
1 import gym
2 import numpy as np
3 from collections import defaultdict
4
5 ###
6
7 def mc_control(num_ep, gamma, eps, eps_min, decay):
8     Q = defaultdict(lambda: np.ones(2)*-np.inf)
9     r = defaultdict(lambda: list())
10
11    for _ in range(num_episodes):
12        episode = generate_episode(Q, eps)
13        G = 0
14        for t in range(len(episode)-1,-1,-1):
15            state, action, reward = episode[t]
16
17            G = gamma*G+reward
18            r[state, action].append(G)
19            Q[state][action] = np.mean(r[state,action])
20
21            eps = max(eps_min, eps*eps_decay)
22
23    return Q
```

## MC Control

### Initialize:

$$\pi(s) \in A(s) \text{ (arbitrarily), } \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily), } \forall s \in S, a \in A(s)$$

$$\text{Returns}(s, a) \leftarrow \text{empty list}, \forall s \in S, a \in A(s)$$

### Loop forever ( $\forall$ episode):

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly, s.d.  $P(S_0, A_0) > 0$

Generate episode from  $S_0, A_0$  following  $\pi$  :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T_1}, R_T$$

$$G \leftarrow 0$$

Loop over each episode,  $t = T - 1, T - 2, \dots, 0$  :

$$G \leftarrow \gamma G + R_{S_t, A_t}$$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  :

Append  $G$  to Returns( $S_t, A_t$ )

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$\pi(S_t) \leftarrow (\text{argmax})_a Q(S_t, a)$$

Source: Sutton et al

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1  def eps_greedy(Q, state, epsilon):
2      if state[0]<12:
3          return 1
4      elif np.random.random() < epsilon:
5          # explore
6          return np.random.choice([0,1])
7      else:
8          # exploit
9          return np.argmax(Q[state])
10
11  #
```

## EPISODE GENERATION

```
1  def generate_episode(Q, eps):
2      state = env.reset()[0]
3      episode = []
4
5      done = False
6      while not done:
7          action = eps_greedy(Q, state, eps)
8          next_state, reward, done, _, _ = env.step(action)
9          episode.append((state, action, reward))
10         state = next_state
11
12  return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

Cont'd

## GREEDY POLICY

```
1 def eps_greedy(Q, state, epsilon):
2     if state[0]<12:
3         return 1
4     elif np.random.random() < epsilon:
5         # explore
6         return np.random.choice([0,1])
7     else:
8         # exploit
9         return np.argmax(Q[state])
10
11 #
```

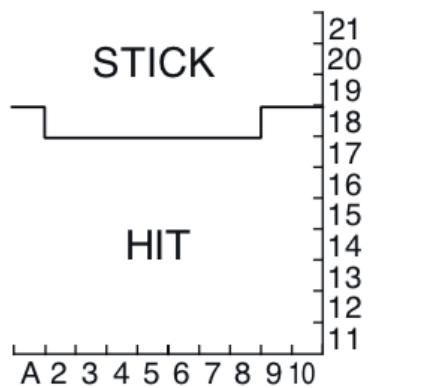
## EPISODE GENERATION

```
1 def generate_episode(Q, eps):
2     state = env.reset()[0]
3     episode = []
4
5     done = False
6     while not done:
7         action = eps_greedy(Q, state, eps)
8         next_state, reward, done, _, _ = env.step(action)
9         episode.append((state, action, reward))
10        state = next_state
11
12 return episode
```

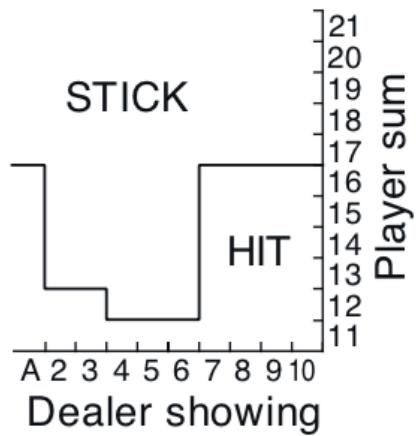


# Optimal Policy

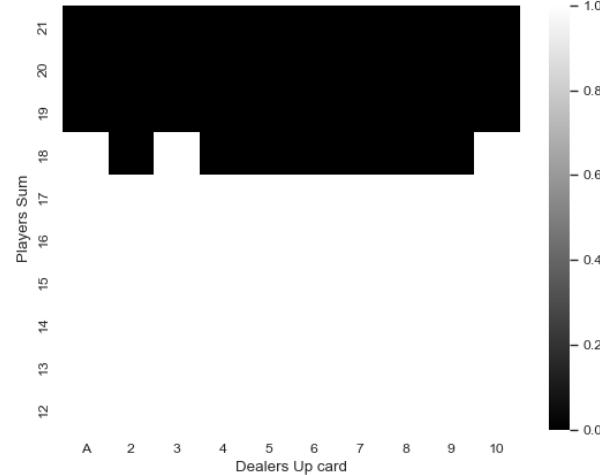
Usable  
ace



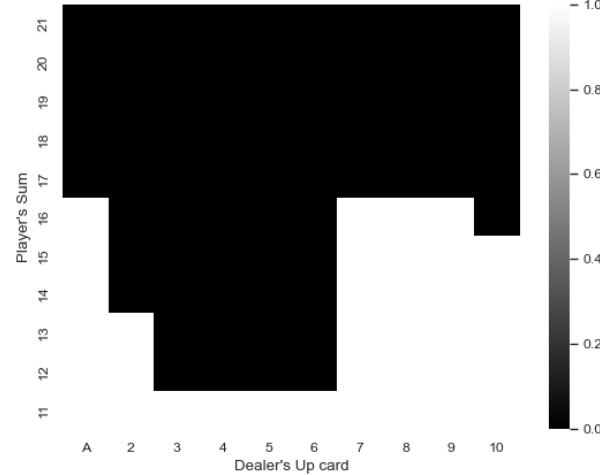
No  
usable  
ace

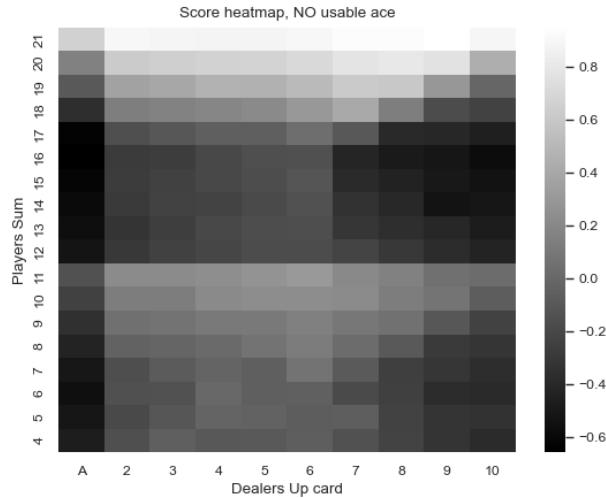
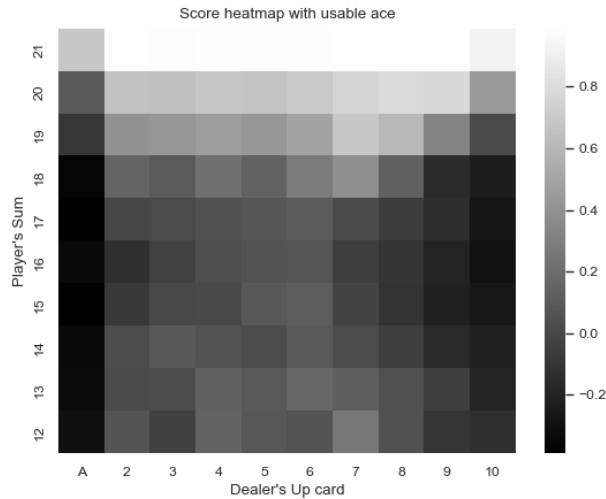


Optimal policy with usable ace

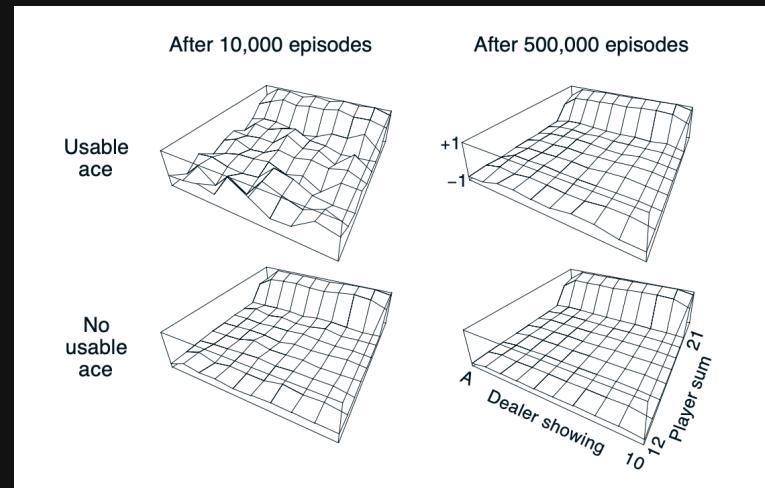


Optimal policy, NO usable ace





# State Score



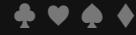
# Discussion

Explore-exploit trade-off

Convergence

About optimality

Off-policy learning



# Thank You For Attention!

Questions?



END