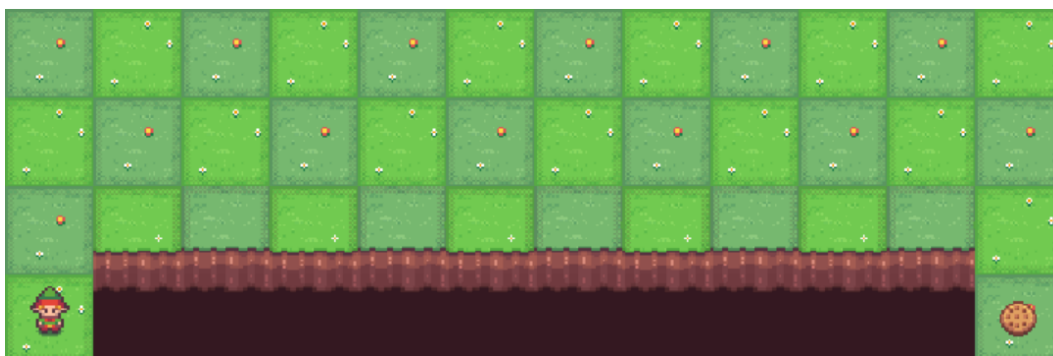


# فرآیند تصمیم مارکوف (MDP)



سید محمد حسین هاشمی ۴۰۲۲۳۶۳۱۴۳  
سید حسین حسینی ۴۰۱۲۳۶۳۲۳۸

آبان ۱۴۰۲

## فهرست مطالب

۱	لیست صخره ها	۲
۲	آماده سازی دیکشنری states	۲
۳	تابع محاسبه امتیاز خانه	۳
۴	تابع به روز کردن ارزش عمل	۴
۵	تابع بروز کردن مقادیر states	۵
۶	حلقه مارکوف	۶
۷	لیست صخره ها	۷

## ۱ لیست صخره ها

cliffwalking\_MDP.py

```
1 | traps = [float(i[0] * 12 + i[1]) for i in env.cliff_positions]
```

شماره خانه های صخره درون لیست traps ذخیره می شود.

## ۲ آماده سازی دیکشنری states

cliffwalking\_MDP.py

```
1 | states = dict()
2 | for i in range(0, 4):
3 |     for j in range(0, 12):
4 |         if i * 12 + j in traps:
5 |             states[i * 12 + j] = -100
6 |         elif i == 3 and j == 11:
7 |             states[i * 12 + j] = 1000
8 |         else:
9 |             states[i * 12 + j] = {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0}
```

از دیکشنری states برای محاسبات استفاده می کنیم. در اینجا مقادیر اولیه آن را تعریف کرده ایم.

## ۳ تابع محاسبه امتیاز خانه

cliffwalking\_MDP.py

```
1 | def calculateScore(inp):
2 |     global states
3 |     return max(states[inp].values()) if isinstance(states[inp], type(dict())) else states[inp]
```

این تابع ارزش خانه فراخوانی شده را برمی گرداند. بدین صورت که در صورت بودن هدف و یا صخره مستقیماً امتیاز آن برگشت داده می شود و در غیر این صورت مقدار بیشترین عمل برگشت داده می شود.

## ۴ تابع به روز کردن ارزش عمل

cliffwalking\_MDP.py

```
1 def changeScore(kp, ks, v1, v2, v3):  
2     global states  
3     score = (v1 + v2 + v3) / 3  
4     change = score - states[kp][ks]  
5     states[kp][ks] = score  
6     return abs(change)
```

از این تابع برای بروز مقدار ارزش هر عمل استفاده می‌شود و مقدار بازگشت داده شده مقدار تغییرات آن است.

## ۵ تابع بروز کردن مقادیر states

```
cliffwalking_MDP.py

1 def updateState():
2     global states
3     totalChanges = 0
4     for Key, Val in states.items():
5         if isinstance(Val, type(dict())):
6             row = Key // 12
7             col = Key % 12
8             up = row * 12 + col if row == 0 else (row - 1) * 12 + col
9             right = row * 12 + col if col == 11 else row * 12 + col + 1
10            down = row * 12 + col if row == 3 else (row + 1) * 12 + col
11            left = row * 12 + col if col == 0 else row * 12 + col - 1
12            for Ks, Vs in Val.items():
13                upScore = calculateScore(up)
14                rightScore = calculateScore(right)
15                leftScore = calculateScore(left)
16                downScore = calculateScore(down)
17
18                if Ks == 0:
19                    totalChanges += changeScore(Key, Ks, upScore, rightScore, leftScore)
20                elif Ks == 1:
21                    totalChanges += changeScore(Key, Ks, upScore, rightScore, downScore)
22                elif Ks == 2:
23                    totalChanges += changeScore(Key, Ks, leftScore, rightScore, downScore)
24                elif Ks == 3:
25                    totalChanges += changeScore(Key, Ks, upScore, rightScore, downScore)
26            return totalChanges
```

این تابع اصلی برای بروز کردن مقادیر states است. بدین صورت که با هر بار اجرا بصورت مستقل عمل می‌کند و هر اجرای آن هر خانه و عمل‌های آن را ارزش‌گذاری و مقدار کل تغییرات را برمی‌گرداند. پس از پیمایش states در صورتی که خانه صخره و یا هدف نبود، موقعیت بعدی آن را از هر سمت به دست می‌آوریم و سپس امتیاز اعمال را با استفاده از توابع بالا بروز رسانی می‌کنیم.

## ۶ حلقه مارکوف

```
cliffwalking_MDP.py

1 n = 1000
2 changes = updateState()
3
4 while n > 0 and changes > 0.005:
5     changes = updateState()
6     n -= 1
7
8 policies = dict()
9 for Kp, Vp in states.items():
10     if isinstance(Vp, type(dict())):
11         Kmax = 1
12         VMax = Vp[1]
13         for K, V in Vp.items():
14             if V > VMax:
15                 Kmax = K
16                 VMax = V
17         policies[Kp] = Kmax
18
19 policies[47] = 1
```

این حلقه همان حلقه معروف فرایند مارکوف است که در صورت کم نشدن مقدار تغییرات هزار بار تابع `updateState` را فراخوانی و مقدار `states` را بروز می‌کند. بعد از آن `policies` را با استفاده از `states` ایجاد کرده‌ایم، بدین صورت که با ارزش‌ترین عمل انتخاب می‌شود. در نهایت برای اینکه خانه ۴۷ خانه هدف است حساب نشده که برای رفع باگ مقدار ۱ (حرکت به راست) را برای آن قرار داده‌ایم.

## ۷ لیست صخره ها

```
cliffwalking_MDP.py

1 # Define the maximum number of iterations
2 max_iter_number = 1000
3 next_state = 36
4 winRate = 0
5 sumRewards = 0
6 for __ in range(max_iter_number):
7     # TODO: Implement the agent policy here
8     # Note: .sample() is used to sample random Action from the environment's Action space
9
10    # Choose an Action (Replace this random Action with your agent's policy)
11    # Action = env.action_space.sample()
12    action = policies[next_state]
13
14    # Perform the Action and receive feedback from the environment
15    next_state, reward, done, truncated, info = env.step(action)
16
17    sumRewards += reward
18
19    if done:
20        winRate += 1
21        sumRewards += 1000
22
23    if done or truncated:
24        observation, info = env.reset()
25
26 # Close the environment
27 env.close()
```

با قرار دادن policies در کد کار به اتمام می‌رسد.