

# Assignment2

April 24, 2021

## 1 Assignment 2: Optimal Policies with Dynamic Programming

Welcome to Assignment 2. This notebook will help you understand: - Policy Evaluation and Policy Improvement. - Value and Policy Iteration. - Bellman Equations.

### 1.1 Gridworld City

Gridworld City, a thriving metropolis with a booming technology industry, has recently experienced an influx of grid-loving software engineers. Unfortunately, the city's street parking system, which charges a fixed rate, is struggling to keep up with the increased demand. To address this, the city council has decided to modify the pricing scheme to better promote social welfare. In general, the city considers social welfare higher when more parking is being used, the exception being that the city prefers that at least one spot is left unoccupied (so that it is available in case someone really needs it). The city council has created a Markov decision process (MDP) to model the demand for parking with a reward function that reflects its preferences. Now the city has hired you — an expert in dynamic programming — to help determine an optimal policy.

### 1.2 Preliminaries

You'll need two imports to complete this assignment: - `numpy`: The fundamental package for scientific computing with Python. - `tools`: A module containing an environment and a plotting function.

There are also some other lines in the cell below that are used for grading and plotting — you needn't worry about them.

In this notebook, all cells are locked except those that you are explicitly asked to modify. It is up to you to decide how to implement your solution in these cells, **but please do not import other libraries** — doing so will break the autograder.

```
[37]: %matplotlib inline
import numpy as np
import tools
import grader
```

In the city council's parking MDP, states are nonnegative integers indicating how many parking spaces are occupied, actions are nonnegative integers designating the price of street parking, the reward is a real value describing the city's preference for the situation, and time is discretized by

hour. As might be expected, charging a high price is likely to decrease occupancy over the hour, while charging a low price is likely to increase it.

For now, let's consider an environment with three parking spaces and three price points. Note that an environment with three parking spaces actually has four states — zero, one, two, or three spaces could be occupied.

```
[38]: # -----
# Discussion Cell
# -----
num_spaces = 3
num_prices = 3
env = tools.ParkingWorld(num_spaces, num_prices)
V = np.zeros(num_spaces + 1)
pi = np.ones((num_spaces + 1, num_prices)) / num_prices
```

The value function is a one-dimensional array where the  $i$ -th entry gives the value of  $i$  spaces being occupied.

```
[39]: V
```

```
[39]: array([0., 0., 0., 0.])
```

We can represent the policy as a two-dimensional array where the  $(i, j)$ -th entry gives the probability of taking action  $j$  in state  $i$ .

```
[40]: pi
```

```
[40]: array([[0.33333333, 0.33333333, 0.33333333],
          [0.33333333, 0.33333333, 0.33333333],
          [0.33333333, 0.33333333, 0.33333333],
          [0.33333333, 0.33333333, 0.33333333]])
```

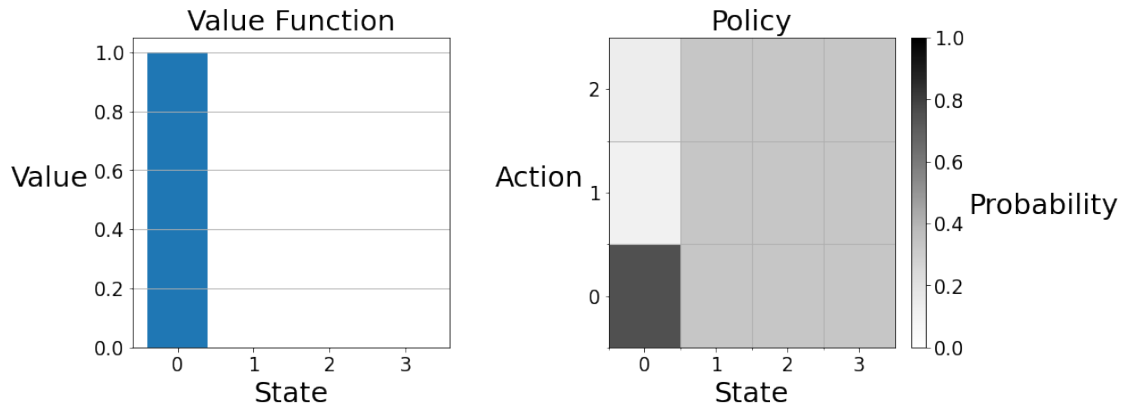
```
[41]: pi[0] = [0.75, 0.11, 0.14]

for s, pi_s in enumerate(pi):
    for a, p in enumerate(pi_s):
        print(f'pi(A={a}|S={s}) = {p.round(2)}', end='')
    print()
```

pi(A=0 S=0) = 0.75	pi(A=1 S=0) = 0.11	pi(A=2 S=0) = 0.14
pi(A=0 S=1) = 0.33	pi(A=1 S=1) = 0.33	pi(A=2 S=1) = 0.33
pi(A=0 S=2) = 0.33	pi(A=1 S=2) = 0.33	pi(A=2 S=2) = 0.33
pi(A=0 S=3) = 0.33	pi(A=1 S=3) = 0.33	pi(A=2 S=3) = 0.33

```
[42]: V[0] = 1

tools.plot(V, pi)
```



We can visualize a value function and policy with the `plot` function in the `tools` module. On the left, the value function is displayed as a barplot. State zero has an expected return of ten, while the other states have an expected return of zero. On the right, the policy is displayed on a two-dimensional grid. Each vertical strip gives the policy at the labeled state. In state zero, action zero is the darkest because the agent's policy makes this choice with the highest probability. In the other states the agent has the equiprobable policy, so the vertical strips are colored uniformly.

You can access the state space and the action set as attributes of the environment.

```
[43]: env.S
```

```
[43]: [0, 1, 2, 3]
```

```
[44]: env.A
```

```
[44]: [0, 1, 2]
```

You will need to use the environment's `transitions` method to complete this assignment. The method takes a state and an action and returns a 2-dimensional array, where the entry at  $(i, 0)$  is the reward for transitioning to state  $i$  from the current state and the entry at  $(i, 1)$  is the conditional probability of transitioning to state  $i$  given the current state and action.

```
[45]: state = 3
      action = 1
      transitions = env.transitions(state, action)
      transitions
```

```
[45]: array([[1.          , 0.12390437],
           [2.          , 0.15133714],
           [3.          , 0.1848436 ],
           [2.          , 0.53991488]])
```

```
[46]: for sp, (r, p) in enumerate(transitions):
      print(f'p(S\='{sp}', R={r} | S={state}, A={action}) = {p.round(2)}')
```

$p(S'=0, R=1.0 \mid S=3, A=1) = 0.12$   
 $p(S'=1, R=2.0 \mid S=3, A=1) = 0.15$   
 $p(S'=2, R=3.0 \mid S=3, A=1) = 0.18$   
 $p(S'=3, R=2.0 \mid S=3, A=1) = 0.54$

### 1.3 Section 1: Policy Evaluation

You're now ready to begin the assignment! First, the city council would like you to evaluate the quality of the existing pricing scheme. Policy evaluation works by iteratively applying the Bellman equation for  $v_\pi$  to a working value function, as an update rule, as shown below.

$$v(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v(s')]$$

This update can either occur “in-place” (i.e. the update rule is sequentially applied to each state) or with “two-arrays” (i.e. the update rule is simultaneously applied to each state). Both versions converge to  $v_\pi$  but the in-place version usually converges faster. **In this assignment, we will be implementing all update rules in-place**, as is done in the pseudocode of chapter 4 of the textbook.

We have written an outline of the policy evaluation algorithm described in chapter 4.1 of the textbook. It is left to you to fill in the `bellman_update` function to complete the algorithm.

```
[47]: # lock
def evaluate_policy(env, V, pi, gamma, theta):
    delta = float('inf')
    while delta > theta:
        delta = 0
        for s in env.S:
            v = V[s]
            bellman_update(env, V, pi, s, gamma)
            delta = max(delta, abs(v - V[s]))

    return V
```

```
[48]: # -----
# Graded Cell
# -----
def bellman_update(env, V, pi, s, gamma):
    """Mutate ``V`` according to the Bellman update equation."""
    # YOUR CODE HERE
    new_v = 0
    for action in env.A:
        for s_, (r,p) in enumerate(env.transitions(s, action)):
            new_v += pi[s,action]*( p*(r+(gamma*V[s_])) )
    V[s] = new_v
```

The cell below uses the policy evaluation algorithm to evaluate the city's policy, which charges a constant price of one.

```
[49]: # -----  
# Debugging Cell  
# -----  
# Feel free to make any changes to this cell to debug your code  
  
# set up test environment  
num_spaces = 10  
num_prices = 4  
env = tools.ParkingWorld(num_spaces, num_prices)  
  
# build test policy  
city_policy = np.zeros((num_spaces + 1, num_prices))  
city_policy[:, 1] = 1  
  
gamma = 0.9  
theta = 0.1  
  
V = np.zeros(num_spaces + 1)  
V = evaluate_policy(env, V, city_policy, gamma, theta)  
  
print(V)
```

```
[80.04173399 81.65532303 83.37394007 85.12975566 86.87174913 88.55589131  
90.14020422 91.58180605 92.81929841 93.78915889 87.77792991]
```

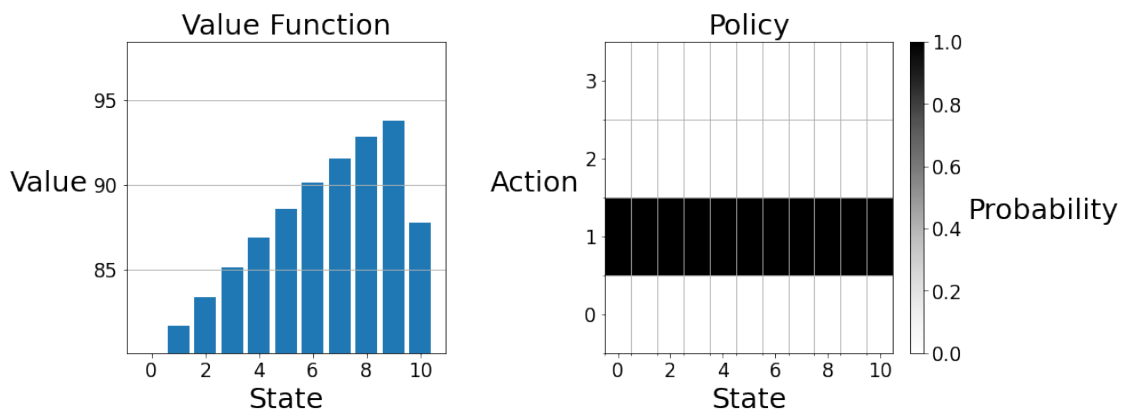
```
[50]: # -----  
# Tested Cell  
# -----  
# The contents of the cell will be tested by the autograder.  
# If they do not pass here, they will not pass there.  
  
# set up test environment  
num_spaces = 10  
num_prices = 4  
env = tools.ParkingWorld(num_spaces, num_prices)  
  
# build test policy  
city_policy = np.zeros((num_spaces + 1, num_prices))  
city_policy[:, 1] = 1  
  
gamma = 0.9  
theta = 0.1  
  
V = np.zeros(num_spaces + 1)  
V = evaluate_policy(env, V, city_policy, gamma, theta)
```

```
# test the value function
answer = [80.04, 81.65, 83.37, 85.12, 86.87, 88.55, 90.14, 91.58, 92.81, 93.78,
↪87.77]

# make sure the value function is within 2 decimal places of the correct answer
assert grader.near(V, answer, 1e-2)
```

You can use the `plot` function to visualize the final value function and policy.

```
[51]: # lock
tools.plot(V, city_policy)
```



Observe that the value function qualitatively resembles the city council’s preferences — it monotonically increases as more parking is used, until there is no parking left, in which case the value is lower. Because of the relatively simple reward function (more reward is accrued when many but not all parking spots are taken and less reward is accrued when few or all parking spots are taken) and the highly stochastic dynamics function (each state has positive probability of being reached each time step) the value functions of most policies will qualitatively resemble this graph. However, depending on the intelligence of the policy, the scale of the graph will differ. In other words, better policies will increase the expected return at every state rather than changing the relative desirability of the states. Intuitively, the value of a less desirable state can be increased by making it less likely to remain in a less desirable state. Similarly, the value of a more desirable state can be increased by making it more likely to remain in a more desirable state. That is to say, good policies are policies that spend more time in desirable states and less time in undesirable states. As we will see in this assignment, such a steady state distribution is achieved by setting the price to be low in low occupancy states (so that the occupancy will increase) and setting the price high when occupancy is high (so that full occupancy will be avoided).

## 1.4 Section 2: Policy Iteration

Now the city council would like you to compute a more efficient policy using policy iteration. Policy iteration works by alternating between evaluating the existing policy and making the policy greedy with respect to the existing value function. We have written an outline of the policy iteration algorithm described in chapter 4.3 of the textbook. We will make use of the policy evaluation algorithm you completed in section 1. It is left to you to fill in the `q_greedify_policy` function, such that it modifies the policy at  $s$  to be greedy with respect to the  $q$ -values at  $s$ , to complete the policy improvement algorithm.

```
[52]: def improve_policy(env, V, pi, gamma):
    policy_stable = True
    for s in env.S:
        old = pi[s].copy()
        q_greedify_policy(env, V, pi, s, gamma)

        if not np.array_equal(pi[s], old):
            policy_stable = False

    return pi, policy_stable

def policy_iteration(env, gamma, theta):
    V = np.zeros(len(env.S))
    pi = np.ones((len(env.S), len(env.A))) / len(env.A)
    policy_stable = False

    while not policy_stable:
        V = evaluate_policy(env, V, pi, gamma, theta)
        pi, policy_stable = improve_policy(env, V, pi, gamma)

    return V, pi
```

```
[53]: # -----
# Graded Cell
# -----
def q_greedify_policy(env, V, pi, s, gamma):
    """Mutate ``pi`` to be greedy with respect to the q-values induced by ``V``.
    ↪ """
    # YOUR CODE HERE
    G = np.zeros_like(env.A, dtype=float)
    for a in env.A:
        transitions = env.transitions(s, a)
        for s_, (r, p) in enumerate(transitions):
            G[a] += p * (r + gamma * V[s_])

    greed_actions = np.argmax(G == np.amax(G))
    for a in env.A:
```

```

    if a in greed_actions:
        pi[s, a] = 1 / len(greed_actions)
    else:
        pi[s, a] = 0

```

```

[54]: # -----
# Debugging Cell
# -----
# Feel free to make any changes to this cell to debug your code

gamma = 0.9
theta = 0.1
env = tools.ParkingWorld(num_spaces=6, num_prices=4)

V = np.array([7, 6, 5, 4, 3, 2, 1])
pi = np.ones((7, 4)) / 4

new_pi, stable = improve_policy(env, V, pi, gamma)

# expect first call to greedify policy
expected_pi = np.array([
    [0, 0, 0, 1],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
])
assert np.all(new_pi == expected_pi)
assert stable == False

# the value function has not changed, so the greedy policy should not change
new_pi, stable = improve_policy(env, V, new_pi, gamma)

assert np.all(new_pi == expected_pi)
assert stable == True

```

```

[55]: # -----
# Tested Cell
# -----
# The contents of the cell will be tested by the autograder.
# If they do not pass here, they will not pass there.
gamma = 0.9
theta = 0.1
env = tools.ParkingWorld(num_spaces=10, num_prices=4)

```



```

V, pi = policy_iteration(env, gamma, theta)

V_answer = [81.60, 83.28, 85.03, 86.79, 88.51, 90.16, 91.70, 93.08, 94.25, 95.
→25, 89.45]
pi_answer = [
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [1, 0, 0, 0],
    [0, 0, 0, 1],
    [0, 0, 0, 1],
]

# make sure value function is within 2 decimal places of answer
assert grader.near(V, V_answer, 1e-2)
# make sure policy is exactly correct
assert np.all(pi == pi_answer)

```

When you are ready to test the policy iteration algorithm, run the cell below.

```

[56]: env = tools.ParkingWorld(num_spaces=10, num_prices=4)
      gamma = 0.9
      theta = 0.1
      V, pi = policy_iteration(env, gamma, theta)

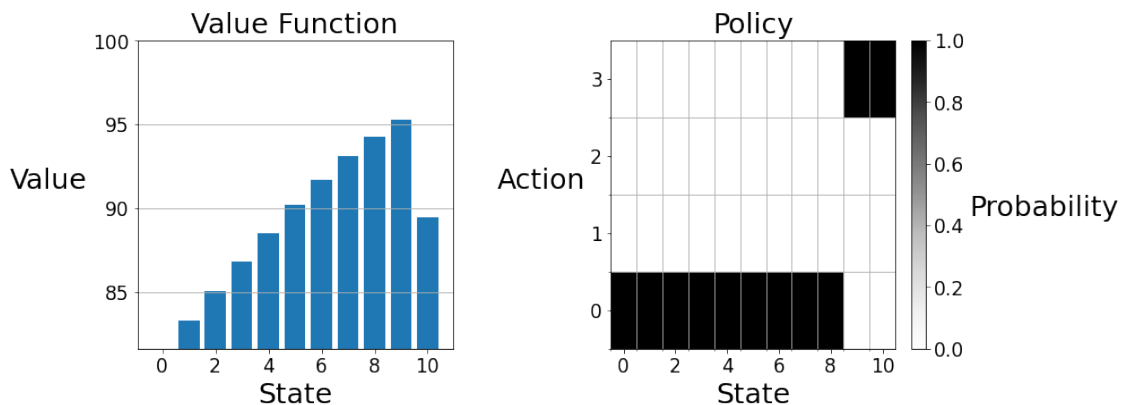
```

You can use the plot function to visualize the final value function and policy.

```

[57]: tools.plot(V, pi)

```



You can check the value function (rounded to one decimal place) and policy against the answer below:

State	Value	Action	0	1	2	3
0	86.8	0	88.5	90.2	91.7	93.1
1	81.6	0	85.0	89.5	93.1	95.3
2	83.3	0	88.5	90.2	91.7	93.1
3	85.0	0	88.5	90.2	91.7	93.1

## 1.5 Section 3: Value Iteration

The city has also heard about value iteration and would like you to implement it. Value iteration works by iteratively applying the Bellman optimality equation for  $v_*$  to a working value function, as an update rule, as shown below.

$$v(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v(s')]$$

We have written an outline of the value iteration algorithm described in chapter 4.4 of the textbook. It is left to you to fill in the `bellman_optimality_update` function to complete the value iteration algorithm.

```
[58]: def value_iteration(env, gamma, theta):
    V = np.zeros(len(env.S))
    while True:
        delta = 0
        for s in env.S:
            v = V[s]
            bellman_optimality_update(env, V, s, gamma)
            delta = max(delta, abs(v - V[s]))
        if delta < theta:
            break
    pi = np.ones((len(env.S), len(env.A))) / len(env.A)
    for s in env.S:
        q_greedy_policy(env, V, pi, s, gamma)
    return V, pi
```

```
[59]: # -----  
# Graded Cell  
# -----  
def bellman_optimality_update(env, V, s, gamma):  
    """Mutate ``V`` according to the Bellman optimality update equation."""  
    # YOUR CODE HERE  
    vmax = - float('inf')  
    for action in env.A:  
        va = 0  
        for s_, (r,p) in enumerate(env.transitions(s, action)):  
            va += p*(r+(gamma*V[s_]))  
        vmax = max(va, vmax)  
    V[s] = vmax
```

```
[60]: # -----
# Debugging Cell
# -----
# Feel free to make any changes to this cell to debug your code

gamma = 0.9
env = tools.ParkingWorld(num_spaces=6, num_prices=4)

V = np.array([7, 6, 5, 4, 3, 2, 1])

# only state 0 updated
bellman_optimality_update(env, V, 0, gamma)
assert list(V) == [5, 6, 5, 4, 3, 2, 1]

# only state 2 updated
bellman_optimality_update(env, V, 2, gamma)
assert list(V) == [5, 6, 7, 4, 3, 2, 1]
```

```
[61]: # -----
# Tested Cell
# -----
# The contents of the cell will be tested by the autograder.
# If they do not pass here, they will not pass there.
gamma = 0.9
env = tools.ParkingWorld(num_spaces=10, num_prices=4)

V = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11])

for _ in range(10):
    for s in env.S:
        bellman_optimality_update(env, V, s, gamma)

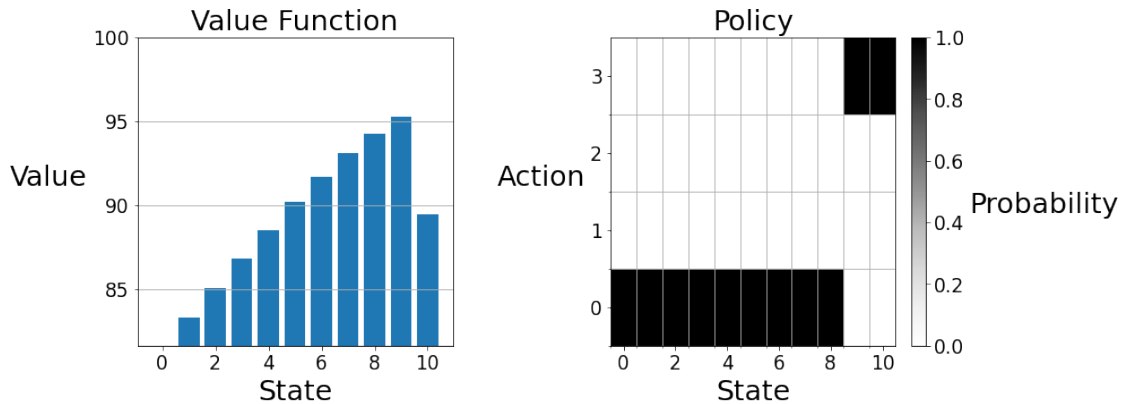
# make sure value function is exactly correct
answer = [61, 63, 65, 67, 69, 71, 72, 74, 75, 76, 71]
assert np.all(V == answer)
```

When you are ready to test the value iteration algorithm, run the cell below.

```
[62]: env = tools.ParkingWorld(num_spaces=10, num_prices=4)
gamma = 0.9
theta = 0.1
V, pi = value_iteration(env, gamma, theta)
```

You can use the plot function to visualize the final value function and policy.

```
[63]: tools.plot(V, pi)
```



You can check your value function (rounded to one decimal place) and policy against the answer below:

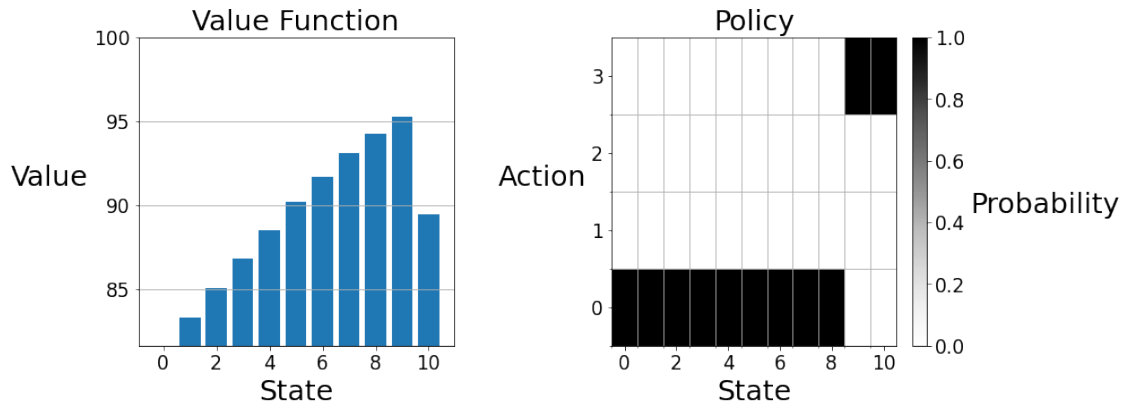
State	Value	Action
0	83.1	0
1	85.0	0
2	86.8	0
3	88.5	0
4	90.2	0
5	91.7	0
6	93.1	0
7	94.3	0
8	95.3	0
9	95.0	3
10	89.5	3

In the value iteration algorithm above, a policy is not explicitly maintained until the value function has converged. Below, we have written an identically behaving value iteration algorithm that maintains an updated policy. Writing value iteration in this form makes its relationship to policy iteration more evident. Policy iteration alternates between doing complete greedifications and complete evaluations. On the other hand, value iteration alternates between doing local greedifications and local evaluations.

```
[64]: def value_iteration2(env, gamma, theta):
    V = np.zeros(len(env.S))
    pi = np.ones((len(env.S), len(env.A))) / len(env.A)
    while True:
        delta = 0
        for s in env.S:
            v = V[s]
            q_greedify_policy(env, V, pi, s, gamma)
            bellman_update(env, V, pi, s, gamma)
            delta = max(delta, abs(v - V[s]))
        if delta < theta:
            break
    return V, pi
```

You can try the second value iteration algorithm by running the cell below.

```
[65]: env = tools.ParkingWorld(num_spaces=10, num_prices=4)
gamma = 0.9
theta = 0.1
V, pi = value_iteration2(env, gamma, theta)
tools.plot(V, pi)
```



## 1.6 Wrapping Up

Congratulations, you've completed assignment 2! In this assignment, we investigated policy evaluation and policy improvement, policy iteration and value iteration, and Bellman updates. Gridworld City thanks you for your service!